



2013

Descriptors for DNA Sequences Based on Joint Diagonalization of Their Feature Matrices from Dinucleotide Physicochemical Properties

Hongjie Yu

Department of Mathematics, School of Science, Anhui Science and Technology University, Fengyang 233100, China

Deshuang Huang

Machine Learning and Systems Biology Laboratory, Tongji University, Shanghai 201804, China

Follow this and additional works at: <https://tsinghuauniversitypress.researchcommons.org/tsinghua-science-and-technology>



Part of the [Computer Sciences Commons](#), and the [Electrical and Computer Engineering Commons](#)

Recommended Citation

Hongjie Yu, Deshuang Huang. Descriptors for DNA Sequences Based on Joint Diagonalization of Their Feature Matrices from Dinucleotide Physicochemical Properties. *Tsinghua Science and Technology* 2013, 18(5): 446-453.

This Research Article is brought to you for free and open access by Tsinghua University Press: Journals Publishing. It has been accepted for inclusion in *Tsinghua Science and Technology* by an authorized editor of Tsinghua University Press: Journals Publishing.

Descriptors for DNA Sequences Based on Joint Diagonalization of Their Feature Matrices from Dinucleotide Physicochemical Properties

Hongjie Yu and Deshuang Huang*

Abstract: Numerical characterizations of DNA sequence can facilitate analysis of similar sequences. To visualize and compare different DNA sequences in less space, a novel descriptors extraction approach was proposed for numerical characterizations and similarity analysis of sequences. Initially, a transformation method was introduced to represent each DNA sequence with dinucleotide physicochemical property matrix. Then, based on the approximate joint diagonalization theory, an eigenvalue vector was extracted from each DNA sequence, which could be considered as descriptor of the DNA sequence. Moreover, similarity analyses were performed by calculating the pair-wise distances among the obtained eigenvalue vectors. The results show that the proposed approach can capture more sequence information, and can jointly analyze the information contained in all involved multiple sequences, rather than separately, whose effectiveness was demonstrated intuitively by constructing a dendrogram for the 15 beta-globin gene sequences.

Key words: descriptors; approximate joint diagonalization; dendrogram; physicochemical property; similarity analysis

1 Introduction

Numerical characterizations of DNA sequences offer a visual means for inspection of data^[1]. Thus, the description of the sequence is of great importance. However, as complex objects may be similar in one aspect and quite different in another, the space of similarity for such objects turns to be multidimensional. Recently, many numerical characterizations for DNA or protein sequences have been introduced, with most of them being extracted from string representations and graphical representations. The string representations also allowed

the extraction of some simpler and more important features, which had been initially used for comparison of the genome sequence^[2], and later for alignment-free comparison of regulatory sequences^[3].

A survey of previous reports^[4,5] revealed many frequency-based algorithms for sequence comparisons and analyses have been based on single as well as dinucleotides. Wu et al. [6] investigated the analysis approaches using adjacent nucleotides of a DNA sequence that revealed their hidden biological information whereas Randić^[7] introduced a condensed representation of DNA based on pairs of nucleotides.

Some researchers^[8,9] have recently introduced a graphical representation of DNA sequences based on the neighboring dinucleotides, as another example of a linear representation for strings of sequences, or have used sparse representation approach to depict objects^[10]. Dual nucleotides may also be categorized into groups according to their physicochemical properties^[11] (for reviews up to 2011).

In general, DNA sequences may be converted into

• Hongjie Yu is with the Department of Mathematics, School of Science, Anhui Science and Technology University, Fengyang 233100, China. E-mail: yhj70@mail.ustc.edu.cn.

• Deshuang Huang is with Machine Learning and Systems Biology Laboratory, Tongji University, Shanghai 201804, China. E-mail: dshuang@tongji.edu.cn.

* To whom correspondence should be addressed.

Manuscript received: 2013-06-15; revised: 2013-09-04; accepted: 2013-09-05

numerical signals through binarization processing of the sequences. Thus, the binary values 0-1 can be used to describe each two adjacent positions irrespective of their dinucleotide combination^[12]. Certainly, the binary representation may be regarded as one of the earliest and the most popular transformations of DNA, although other different transformation methods^[13-22] or sparse representation have also been explored.

Among these transformation methods, some are not intended for biological problems, while others have no simple numerical interpretation. Moreover, some of the representations neglect the sequence structure and the transformations tend to be irreversible. Until now, there is no ideal approach that is able to transform every type of DNA/protein sequences so as to numerically analyze the relationship among them with efficiency.

In this study, we introduce a novel method for numerical characterization of DNA sequences and apply the method to similarity analysis of these sequences, where we derived descriptors from sequences. The application of the numerical characterization of DNA sequence is illustrated by examining the relationship among different species. The results from comparison with earlier approaches demonstrated that our proposed approach is effective and is capable of analyzing the similarities among multiple sequences.

2 Descriptors of DNA Sequences

Numerical characterization of DNA sequences can facilitate us to perform similarity analysis of multiple sequences. In this section, we will introduce a novel method to transform each DNA sequence into a symmetric matrix, through which feature vectors may be readily derived.

2.1 Construction of feature matrix for sequence

Numerical characterization of 2-D or 3-D graphical representations for DNA sequences have been widely used, where some feature matrices were transformed from the primary sequences. Thus, the descriptors could be derived from these matrices. These descriptors characterizing sequence can be used as the components of similarity measures between a pair of sequences^[11].

Considering a DNA sequence $S = "S_1S_2 \cdots S_L"$, where $S_i \in \{A, T, G, C\}$, $i = 1, 2, \dots, L$, and L denotes the length of the sequence, there are 16 kinds of dinucleotides in total (shown in Table 1). It is known that the four nucleic acids A, T, G, and C can be grouped

Table 1 16 kinds of dinucleotides.

$S_j \setminus S_{j+1}$	A	T	G	C
A	AA	AT	AG	AC
T	TA	TT	TG	TC
G	GA	GT	GG	GC
C	CA	CT	CG	CC

in line with the distributions of purine-pyrimidine (R/Y), amino-keto (M/K), and weak-strong (W/S) bonds. The classifications are as follows: R = (A, G) and Y = (C, T), M = (A, C) and K = (G, T), W = (A, T) and S = (C, G).

According to Ref. [11], many different binary techniques have been assigned the values 0-1 to Y, K, S and to R, M, W, respectively. Considering graph theory, DNA/protein sequences may be regarded as node-edge-node models, where there are four types of nodes, A-T-G-C, and 16 kinds of edges (shown in Table 1). However, all the 16 kinds of edges may be grouped into 12 kinds of dinucleotide bonds by combining every two adjacent loci.

Table 2 exhibits the decision criteria for converting dinucleotides into 12-tuple row vectors. Therefore, by scanning every adjacent site successively, such as pairs of loci $(S_1, S_2), (S_2, S_3), \dots, (S_{L-1}, S_L)$, a $(L - 1)$ by 12 adjacency matrix transformed from the primary sequence may be obtained via the relationships of all the adjacent dinucleotides, referred to as m :

$$m = (\alpha_i)_{(L-1) \times 12} \tag{1}$$

where

$$\alpha_i = \text{the } i\text{-th row vector} \tag{2}$$

Table 2 Decision criterion for mapping dinucleotides into 0-1 values according to the combination of bonds.

$S_j S_{j+1}$	RR	RY/YR	YY	MM	MK/KM	KK	WW	WS/SW	SS
AA	1			1			1		
AT		1			1		1		
AG	1				1			1	
AC		1		1				1	
TA		1			1		1		
TT			1			1	1		
TG		1				1		1	
TC			1		1			1	
GA	1				1			1	
GT		1				1		1	
GG	1					1			1
GC		1			1				1
CA	1		1					1	
CT		1			1			1	
CG		1			1				1
CC			1	1					1

if $S_i S_{i+1}$ is the i -th kind of dinucleotide, $i = 1, 2, \dots, L - 1$.

Generally, the primary DNA/protein sequences could be considered as symbolic signals which may have a rich statistical structure. Many signal processing algorithms are focused on the organization of these statistical structures, such as the least square regression for utilizing the data correlation structure^[22].

Since symmetric matrix has many merits^[23], we can use the sparse matrix $\mathbf{m}_{(L-1) \times 12}$ mapped from the primary sequence to obtain a symmetric feature matrix $\mathbf{M}_{12} = \mathbf{m}^T \times \mathbf{m}$ for representing each sequence. Furthermore, as each sequence has a different length, \mathbf{M}_{12} could be divided by the sum of all its elements to eliminate the negative influence from non-equal length among multiple sequences. Thus, all the transformed feature matrices could become comparable.

According to signal processing, the symmetric matrix \mathbf{M}_{12} may be interpreted as the observations from a two-layer ensemble of "sensors". The virtual sensors have an input via sixteen dinucleotides while the output is a 12×12 symmetric matrix for each DNA sequence (as depicted in Tables 1 and 2). Thus, the matrix analysis approach may be generalized from the field of signal processing to that of bioinformatics, for performing multiple sequences similarity analyses.

2.2 Feature extraction from multiple sequences via approximate joint diagonalization upon matrices

For all the multiple sequences, we can use Approximate Joint Diagonalization (AJD) of their corresponding transformed matrices $\mathcal{M} = \{\mathbf{M}^{(1)}, \mathbf{M}^{(2)}, \dots, \mathbf{M}^{(N)}\}$, which has been successfully applied in the dataset with the first exon from sequences of 11 beta globin genes^[24].

In brief, AJD corresponds to the problem of seeking a matrix \mathbf{V} , which will lead $\mathbf{V}^H \mathbf{M}^{(n)} \mathbf{V}$ to be the diagonal as possible for all n , where \mathbf{V} is a unitary matrix. This is based on the premise that a set of matrices $\{\mathbf{M}^{(1)}, \mathbf{M}^{(2)}, \dots, \mathbf{M}^{(N)}\}$ consists of common statistical information of the observations which are the estimates of matrices in the form $\mathbf{V}^H \mathbf{M}^{(n)} \mathbf{V}$.

In general, for any $n \times n$ matrix \mathbf{V} , the AJD criterion may be defined as the following non-negative function of \mathbf{V} :

$$J(\mathbf{V}, \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}) \stackrel{\text{def}}{=} \sum_{i=1}^N \|\mathbf{A}^{(i)} - \mathbf{V}^H \mathbf{M}^{(i)} \mathbf{V}\|^2 \quad (3)$$

Usually, AJD does not require that the involved matrix set \mathcal{M} be exactly simultaneously diagonalized by a common unitary matrix. Mostly, the criterion for AJD, as indicated by Eq. (3), cannot be zeroed, and the matrices can only be approximately jointly diagonalized. Thus, AJD deals with a kind of an "average eigen-structure", which is particularly convenient for statistically inferring the structural information extracted from sample statistics.

Considering two transformations:

$$\tau_1 : \text{Sequence}^{(i)} \mapsto \mathbf{M}^{(i)}.$$

Sequence⁽ⁱ⁾ denotes the i -th sequence, where the length of the sequence is L , and $i = 1, 2, \dots, N$, while $\mathbf{M}^{(i)} \in \mathbb{R}^{12 \times 12}$ stands for the corresponding matrices mapped from each primary DNA sequences, and $\mathbf{M}^{(i)}$ is a symmetric matrix, which can be determined by scanning along the sequence $S^{(i)}$ via the decision criterion listed in Table 2.

$$\tau_2 : \mathbf{M}^{(i)} \mapsto (\lambda_1^{(i)}, \lambda_2^{(i)}, \dots, \lambda_{12}^{(i)}).$$

The feature vector $\mathbf{F}_{12}^{(i)} = (\lambda_1^{(i)}, \lambda_2^{(i)}, \dots, \lambda_{12}^{(i)})$ is a 12-tuple vector consisting of all the eigenvalues extracted by AJD upon $\mathbf{M}^{(i)}$. Thus, compound transformation may be obtained as follows:

$$\tau_2 \circ \tau_1 : \text{Sequence}^{(i)} \mapsto (\lambda_1^{(i)}, \lambda_2^{(i)}, \dots, \lambda_{12}^{(i)}) \quad (4)$$

From Formula (4), we can freely extract the features of the DNA sequence. From the viewpoint of algebra space, the transformation may also be presented as

$$\text{Ker } f : \mathcal{S} \xrightarrow{\tau} \mathbf{F}^{1 \times 12} \quad (5)$$

where \mathcal{S} denotes the original sequence space comprising of primary DNA sequence having the length L , while $\mathbf{F}^{1 \times 12}$ indicates the *objective feature space* that is transformed from the original space. Further, the diagonal elements of \mathbf{A} are simply the eigenvalues of the dinucleotide PhysicoChemical Matrix (PCM) via AJD.

According to the results from previous work^[24], it was determined that the AJD-PCM algorithm has the property of distance-preserving. Thus, we can calculate all the Eigenvalue Vectors (EVs) for each obtained dinucleotide PCM, such as $\mathbf{F}_{12}^{(i)} = (\lambda_1^{(i)}, \lambda_2^{(i)}, \dots, \lambda_{12}^{(i)})$, $i = 1, 2, \dots, N$. Further, the N corresponding 12-tuple vectors may be obtained that can be regarded as features extracted from the original DNA sequence. The AJD-PCM algorithm is shown in Algorithm 1.

3 Numerical Characterization of Sequences

Most numerical characterization methods can represent each sequence only *separately* rather than *jointly*. In this

Algorithm 1 Descriptors for DNA sequences via AJD-PCM

Input: Multiple DNA sequences with different length $L, S^{(1)}, S^{(2)}, \dots, S^{(N)}$.

Initialize: Tol - An imposed tolerance on the change in objective function for a stopping condition

begin

- **for** $n = 1$ **to** N **do**
 Transform original sequences $S^{(n)}$ into 12 by 12 symmetric matrix $M^{(n)}$

end for

- Consider the obtained matrix set $\mathcal{M} = \{M^{(1)}, M^{(2)}, \dots, M^{(N)}\}$ and objective function

$$J(V, \Lambda^{(1)}, \Lambda^{(2)}, \dots, \Lambda^{(N)}) \stackrel{\text{def}}{=} \sum_{i=1}^N \|\Lambda^{(i)} - V^H M^{(i)} V\|^2.$$

while $J(V, \Lambda^{(1)}, \Lambda^{(2)}, \dots, \Lambda^{(N)}) \geq \text{tolerance}$ **do**

Update V using AC-DC algorithm^[25]

end while

- **for** $n = 1$ **to** N **do**
 $F^{(n)} \leftarrow \text{diag}(\Lambda^{(n)})$
 Plot and categorize the N feature curves with $F^{(n)}$
- end for**
- Calculate pairwise distances using all these N vectors $F^{(n)}$
- Draw the dendrogram using the pairwise distances matrix
- end**

section, a novel method of numerical characterization and graphical representation for DNA sequences is presented, which can *jointly* consider the mutual information of all the involved sequences. We selected a dataset comprising of 15 sequences of the first exon in the *beta-globin gene*^[26], as shown in Table 3.

3.1 Calculation of eigenvalue vectors

To enable comparison of multi-sequences, all matrices

Table 3 The concise information from beta-globin genes of 15 species.

Species	AC (GeneBank)	Location	Length (nt)
Human	U01317	62 187-63 610	1424
Chimpanzee	X02345	4189-5532	1344
Gorilla	X61109	4538-5881	1344
Lemur	M15734	154-1595	1442
Rat	X06701	310-1505	1196
Mouse	V00722	275-1462	1188
Goat	M15387	279-1749	1471
Sheep	DQ352470	238-1708	1471
Mouflon	DQ352468	238-1706	1469
Bovine	X00376	278-1741	1464
Rabbit	V00882	277-1419	1143
European hare	Y00347	1485-2620	1136
Opossum	J03643	467-2488	2022
Gallus	V00409	465-1810	1346
Muscovy duck	X15739	291-1870	1580

extracted from each corresponding sequence were normalized by dividing each M_{12} with the sum of values for all their elements. According to the procedure for AJD-PCM algorithm depicted in Section 2.2, all the 12-tuple EVs are calculated through AJD of all the eleven neighboring nucleotide matrices (PCM). The results are listed in Table 4, where each column denotes a 12-tuple EV for each sequence. If each vector is connected in an orderly manner as head and tail using the walk strategy, 15 curves may be plotted, as shown in Fig. 1.

3.2 Convergence speed analysis

Based on the convergence analysis of AJD, the convergence speed of the optimization scheme (see Eq. (3)) should be considered. The diagonalizer matrix V of the objective function at the optimum point depends upon an unknown least error. Thus, a simplified expression for the least error is as follows:

$$\text{Err}(j) \stackrel{\text{def}}{=} \sum_{i=1}^N \|\Lambda_j^{(i)} - V_j^H M^{(i)} V_j\|_F^2 \quad (6)$$

where j ranges from 2 to a presupposed maximal number of iterations.

The predefined error threshold is denoted as ϵ when $|\text{Err}(j+1) - \text{Err}(j)| < \epsilon$, which indicates that when the error does not change drastically from the j -th iteration to the $(j+1)$ -th one, we can get either the obtained diagonalizer V or the diagonalized $\Lambda^{(i)}$, which just occurs at j -th step. Figure 2 shows a plot of the errors with each iteration, from which it could be found that the AJD algorithm converged at just the third iteration or so.

4 Similarity Analysis of Sequences Based on Their Extracted Descriptors

Based on sequence descriptors extracted via AJD-PCM, multiple sequences can be compared.

4.1 Calculation of pair-wise distances

Consequently, numerical characterization was applied for investigating the similarity of multiple sequences from the dataset listed in Table 3. As described in Section 2.2, the genetic distance was calculated in the present study. For every two EVs: $F^{(i)}$ and F , which have been arranged within the corresponding columns in Table 4, it was found that the degree of dissimilarity may be determined through Euclidean distance between every two column vectors. The Euclidean distance between the i -th and the j -th sequences may be

Table 4 Descriptors extracted via AJD-PCM from the beta-globin genes of 15 species.

No.	Human	Chimpanzee	Gorilla	Lemur	Rat	Mouse	Goat	Sheep	Mouflon	Bovine	Rabbit	European hare	Opossum	Gallus	Muscovy duck
1	0.01991	0.01941	0.01965	0.02109	0.02271	0.02186	0.02550	0.02551	0.02567	0.02563	0.02134	0.02168	0.02312	0.03092	0.03209
2	0.02429	0.02415	0.02413	0.02407	0.02347	0.02335	0.02355	0.02355	0.02350	0.02285	0.02418	0.02480	0.02042	0.02419	0.02340
3	0.02706	0.02733	0.02743	0.02630	0.02523	0.02516	0.02435	0.02428	0.02431	0.02403	0.02531	0.02589	0.02378	0.01984	0.01844
4	0.03013	0.03115	0.03093	0.03071	0.03176	0.03274	0.02885	0.02873	0.02839	0.02862	0.03213	0.03088	0.03451	0.01962	0.01949
5	0.02377	0.02211	0.02211	0.02248	0.02644	0.02598	0.02541	0.02574	0.02587	0.02629	0.02416	0.02381	0.02636	0.03274	0.03350
6	0.02581	0.02612	0.02633	0.02553	0.02851	0.02884	0.02978	0.02932	0.02924	0.02869	0.02644	0.02731	0.02949	0.02648	0.02670
7	0.02147	0.02122	0.02168	0.02201	0.02422	0.02565	0.02706	0.02671	0.02685	0.02672	0.02346	0.02408	0.02448	0.02840	0.02919
8	0.03647	0.03788	0.03653	0.04367	0.03489	0.03309	0.02695	0.02746	0.02661	0.02702	0.04114	0.03594	0.02218	0.03033	0.02430
9	0.04168	0.04436	0.04452	0.03447	0.02542	0.02707	0.01895	0.01940	0.01952	0.02077	0.02937	0.02921	0.03577	-0.01266	-0.01326
10	0.02849	0.02850	0.02847	0.02954	0.03132	0.03046	0.03159	0.03172	0.03157	0.03096	0.03085	0.03116	0.03030	0.03171	0.03188
11	0.02896	0.02864	0.02887	0.02941	0.03042	0.02978	0.03097	0.03110	0.03110	0.03078	0.03104	0.03175	0.02974	0.03334	0.03319
12	0.01845	0.01779	0.01732	0.02002	0.01993	0.02132	0.02191	0.02111	0.02114	0.02174	0.01955	0.01937	0.01730	0.03827	0.03751

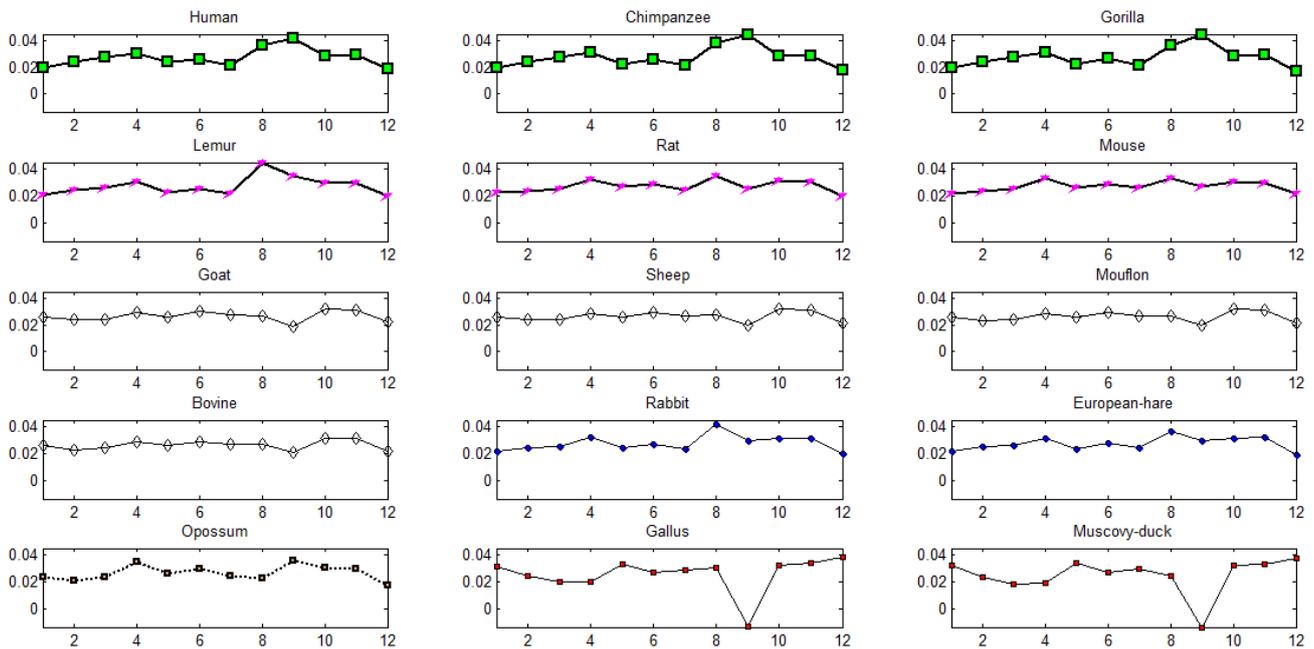


Fig. 1 Graphical representation of the exon in the beta-globin gene from 15 species based on 12-tuple EVs via AJD upon all the corresponding PCP matrices. The y-axis indicates the values of each element in feature vectors $(\lambda_1, \lambda_2, \dots, \lambda_{12})$.

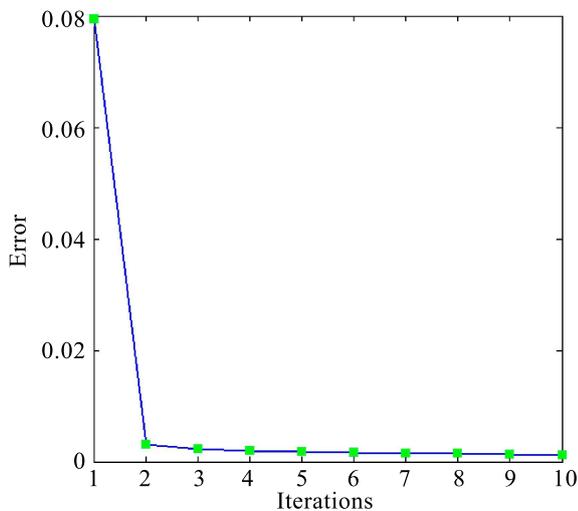


Fig. 2 The performance curve for the convergence speed of AJD for the dataset listed in Table 3.

calculated as:

$$D(S^{(i)}, S^{(j)}) \stackrel{\text{def}}{=} \|M^{(i)} - M^{(j)}\|_F = \|F^{(i)} - F^{(j)}\|_F \tag{7}$$

where $F^{(i)} = (\lambda_1^{(i)}, \lambda_2^{(i)}, \dots, \lambda_{12}^{(i)})$ denotes the feature vectors through distance-preserving transformations from the primary sequences with different lengths, while $\|\cdot\|_F$ indicates the Frobenius norm of a matrix or vector. Obviously, a higher distance value indicates more dissimilarity between the two sequences. The data for comparison of dissimilarity among these 15 coding sequences was obtained by calculating the Euclidean distance via Eq. (7). The pair-wise distance results are listed in Table 5.

4.2 Phylogeny of 15 beta-globin genes

Other than the alphabet representation of biological sequences, which is easily processed on a computer, it

Table 5 Pairwise distance for the beta-globin genes of 15 species.

Species $i \setminus j$	Chimpanzee	Gorilla	Lemur	Rat	Mouse	Goat	Sheep	Mouflon	Bovine	Rabbit	European hare	Opossum	Gallus	Muscovy duck
Human	0.0037	0.0036	0.0106	0.0178	0.0169	0.0269	0.0262	0.0264	0.0252	0.0141	0.0136	0.0181	0.0619	0.0636
Chimpanzee		0.0016	0.0120	0.0208	0.0199	0.0301	0.0294	0.0297	0.0285	0.0164	0.0165	0.0204	0.0652	0.0670
Gorilla			0.0128	0.0207	0.0197	0.0297	0.0290	0.0292	0.0280	0.0168	0.0165	0.0193	0.0652	0.0669
Lemur				0.0141	0.0146	0.0247	0.0240	0.0245	0.0235	0.0068	0.0103	0.0234	0.0562	0.0588
Rat					0.0036	0.0117	0.0109	0.0116	0.0107	0.0082	0.0055	0.0172	0.0461	0.0478
Mouse						0.0118	0.0112	0.0117	0.0106	0.0095	0.0060	0.0151	0.0472	0.0488
Goat							0.0013	0.0014	0.0025	0.0193	0.0153	0.0196	0.0385	0.0392
Sheep								0.0010	0.0021	0.0185	0.0145	0.0192	0.0391	0.0398
Mouflon									0.0019	0.0192	0.0151	0.0190	0.0391	0.0397
Bovine										0.0182	0.0142	0.0180	0.0398	0.0405
Rabbit											0.0056	0.0212	0.0512	0.0536
European hare												0.0171	0.0499	0.0518
Opossum													0.0570	0.0571
Gallus														0.0065

becomes difficult to observe other differences without aid, for which a phylogenetic tree provides a simpler means to inspect various biological sequences that facilitates sequence comparison with the intuitive patterns or pictures. Thus, we examined our proposed approach (AJD-PCM) through phylogenetic analysis, which may be summarized as follows:

(1) Initially, the 12-tuple EVs of each biological sequence were calculated through AJD-PCM;

(2) Then, the similarity distance was obtained with the Euclidean metric to form a pair-wise distance matrix;

(3) Finally, based on the pair-wise distance matrix, the dendrogram was plotted using MATLAB code.

The results of the experiment have been listed in Table 5 and Fig. 3, from which it was found that

the 15 species were clearly separated.

Figure 3 displays the intuitive categories of the 15 species into four subgroups according to the divergence between themselves and humans as follows:

(1) The first subgroup [(Gorilla, Chimpanzee), Human] is closest to the humans.

(2) The next subgroup includes the species [Bovine, (Goat, (Mouflon, Sheep))].

(3) The third ((Muscovy duck, Gallus), Opossum) is far away from the humans in the light of evolutionary relationships. Generally, these three species may be regarded as outgroup.

(4) The remaining five species fall into the last group.

These results are consistent with the evolutionary facts.

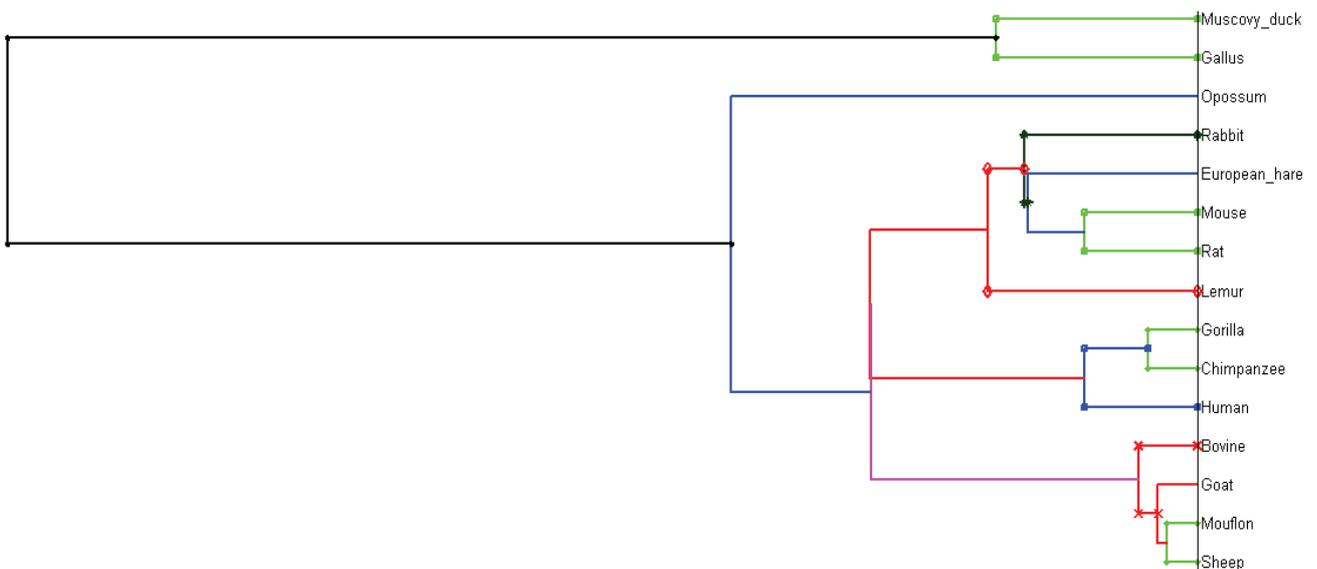


Fig. 3 The dendrogram for 15 sequences according to the pairwise distance listed in Table 3.

4.3 Comparison with related representative work

To comprehend the effectiveness of our approach, we compared the pair-wise distances obtained from the present study (Table 5) with those in the seventh table from another related representative work^[27]. Considering the similarity among Humans and the remaining 14 species, we extracted the first row (or column) from the two tables mentioned above, and calculated the correlation coefficient for these two reduced vectors, that is, the first row of Table 5 in this study, and the first column from the seventh table from Ref. [27]. The correlation degree was determined to be 98.76%. Remarkably, comparison of the two whole tables above-mentioned shows the correlation degree also reaches up to 93.39%.

5 Conclusions and Future Work

The proposed approach (AJD-PCM) allows the fusion of the physicochemical properties with the sequential property of the biological sequence with consideration of the sequential property at the stage of mapping each dinucleotide into a 12-tuple vector. Moreover, at the second stage, AJD could extract the features among multiple sequences *jointly* rather than *separately*, facilitating the simultaneous discovery of sub-groups of organisms having a common structure at the molecular level. The clustering results were consistent with the evolutionary facts demonstrating the rationality of the proposed numerical characterization of DNA sequence. Further work would involve enhancing our descriptor extraction algorithms for application towards genomic/proteomic sequence datasets.

Acknowledgements

This work was supported by the Key Project from Education Department of Anhui Province (No. KJ2013A076), the PhD Programs Foundation of Ministry of Education of China (No. 20120072110040), the National Natural Science Foundation of China (Nos. 61133010, 31071168, and 61005010), and the China Postdoctoral Science Foundation (No. 2012T50582).

References

- [1] A. Nandy, M. Harle, and S. C. Basak, Mathematical descriptors of DNA sequences: Development and applications, *ARKIVOC*, vol. ix, pp. 211-238, 2006.
- [2] B. E. Blaisdell, A measure of the similarity of sets of sequences not requiring sequence alignment, *Proceedings of the National Academy of Sciences*, vol. 83, pp. 5155-5159, 1986.
- [3] M. R. Kantorovitz, G. E. Robinson, and S. Sinha, A statistical method for alignment-free comparison of regulatory sequences, *Bioinformatics*, vol. 23, no. 13, pp. i249-i255, 2007.
- [4] G. E. Sims, S. R. Jun, G. A. Wu, and S. H. Kim, Alignment-free genome comparison with feature frequency profiles (FFP) and optimal resolutions, *Proceedings of the National Academy of Sciences*, vol. 106, no. 8, pp. 2677-2682, 2009.
- [5] S. R. Jun, G. E. Sims, G. A. Wu, and S. H. Kim, Whole-proteome phylogeny of prokaryotes by feature frequency profiles: An alignment-free method with optimal feature resolution, *Proceedings of the National Academy of Sciences*, vol. 107, no. 1, pp. 133-138, 2009.
- [6] Y. Wu, A. W-C Liew, H. Yan, and M. S. Yang, DB-Curve: A novel 2D method of DNA sequence visualization and representation, *Chemical Physics Letters*, vol. 367, pp. 170-176, 2003.
- [7] M. Randić, Condensed representation of DNA primary sequences, *Journal of Chemistry Information & Computer Science*, vol. 40, no. 1, pp. 50-56, 2000.
- [8] Z. B. Liu, B. Liao, and W. Zhu, A new method to analyze the similarity based on dual nucleotides of the DNA sequence, *MATCH*, vol. 61, pp. 541-552, 2009.
- [9] Z. B. Liu, B. Liao, W. Zhu, and G. H. Huang, A 2D graphical representation of DNA sequence based on dual nucleotides and its application, *International Journal of Quantum Chemistry*, vol. 109, no. 5, pp. 948-958, 2009.
- [10] C. Y. Lu, H. Min, J. Gui, L. Zhu, and Y. K. Lei, Face recognition via weighted sparse representation, *Journal of Visual Communication and Image Representation*, vol. 24, no. 2, pp. 111-116, 2013.
- [11] D. Bielinska-Waz, Graphical and numerical representations of DNA sequences: Statistical aspects of similarity, *Journal of Mathematical Chemistry*, vol. 49, no. 10, pp. 2345-2407, 2011.
- [12] R. F. Voss, Evolution of long-rang fractal correlations and 1/f noise in DNAbase sequences, *Physical Review Letter*, vol. 68, pp. 3805-3808, 1992.
- [13] M. Akhtar, J. Epps, and E. Ambikairajah, On DNA numerical representation for period-3 based exon prediction, in *5th International Workshop on Genomic Signal Processing and Statistics, Tuusula*, Piscataway, NJ, USA, 2007.
- [14] H. J. Jeffrey, Chaos game representation of gene structure, *Nucleic Acids Research*, vol. 18, no. 8, pp. 2163-2170, 1990.
- [15] C. Y. Lu and D. S. Huang, Optimized projections for sparse representation based classification, *Neurocomputing*, vol. 113, pp. 213-219, 2013.
- [16] R. Zhang and C. T. Zhang, Z curves, an intuitive tool for visualizing and analyzing the DNA sequences, *Journal of Biomolecular Structure & Dynamics*, vol. 11, no. 4, pp. 767-782, 1994.
- [17] M. Randić, Another look at the chaos-game representation of DNA, *Chemical Physics Letters*, vol. 456, no. 1, pp. 84-88, 2008.

- [18] S. Wang, F. Tian, W. Feng, and X. Liu, Applications of representation method for DNA sequences based on symbolic dynamics, *Journal of Molecular Structure: THEOCHEM*, vol. 909, pp. 33-42, 2009.
- [19] A. K. Brodzik and O. Peters, Symbol-balanced quaternionic periodicity transform for latent pattern detection in DNA sequences, in *Proceedings of IEEE ICASSP*, Philadelphia, PA, USA, 2005, pp. 373-376.
- [20] B. Liao, M. Tan, and K. Ding, Application of 2-D graphical representation of DNA sequence, *Chemical Physics Letters*, vol. 414, pp. 296-300, 2005.
- [21] W. Wang and D. H. Johnson, Computing linear transforms of symbolic signals, *IEEE Transactions on Signal Processing*, vol. 50, no. 3, pp. 628-634, 2002.
- [22] C. Y. Lu, H. Min, Z. Z. Zhao, L. Zhu, D. S. Huang, and S. C. Yan, Robust and efficient subspace segmentation via least squares regression, *European Conference on Computer Vision ECCV*, vol. 7578, no. 7, pp. 347-360, 2012.
- [23] G. H. Golub and C. F. V. Loan, *Matrix Computations*, 3rd Ed. Baltimore and London: Johns Hopkins University Press, 1996.
- [24] H. J. Yu and D. S. Huang, Graphical representation for DNA sequences via joint diagonalization of matrix pencil, *IEEE Journal of Biomedical and Health Informatics*, vol. 17, no. 3, pp. 503-511, 2013.
- [25] A. Yeredor, Non-orthogonal joint diagonalization in the least-squares sense with application in blind source separation, *IEEE Transactions on Signal Processing*, vol. 50, no. 7, pp. 1545-1553, 2002.
- [26] Q. Dai, X. Q. Liu, Y. H. Yao, and F. K. Zhao, Sequence comparison via polar coordinates representation and curve tree, *Journal of Theoretical Biology*, vol. 292, pp. 78-85, 2011.
- [27] C. Li, H. Ma, Y. Zhou, X. L. Wang, and X. Q. Zheng, Similarity analysis of DNA sequences based on the weighted pseudo-entropy, *Journal of Computational Chemistry*, vol. 32, no. 4, pp. 675-680, 2011.



Hongjie Yu received his BS degree in mathematics from Anhui University, China in 1996, his MS degree in mathematics from Hefei University of Technology, China in 2004, and his PhD degree in pattern recognition and intelligent system from University of Science and Technology of China in 2013. Since

1996, he has been with the Anhui Science and Technology University, Fengyang, Anhui Province, China, where he is now an associate professor in the Department of Mathematics. His current research interests include bioinformatics, biostatistics, machine learning, and intelligent computing.



Deshuang Huang received his MS and PhD degrees in electronic engineering from the National Defense University of Science and Technology and Xidian University, China, in 1989 and 1993, respectively. He is a chaired professor in Department of Computer Science and the director of Machine Learning and Systems

Biology Laboratory at Tongji University, China. His primary research interests include neural networks, pattern recognition, and bioinformatics. He is currently a Bioinformatics and Bioengineering/Neural Networks Technical Committee member of IEEE CIS, a member of the International Neural Network Society (INNS), a board member of the INNS Governors, and an associated editor of several mainstream international journals such as *Neural Networks*. He has published more than 300 papers in international journals, international conferences proceedings, and book chapters. He also published three monographs (in Chinese), one of which, entitled *Systematic Theory of Neural Networks for Pattern Recognition*, won the Second-Class Prize of the Eighth Excellent High Technology Books of China in 1997. He is currently a senior member of the IEEE.