

The Influence of Different Structure Representations on the Clustering of an RNA Nucleotides Data Set

T. H. Reijmers, R. Wehrens, and L. M. C. Buydens*

Laboratory for Analytical Chemistry, University of Nijmegen, Toernooiveld 1,
6525 ED Nijmegen, The Netherlands

Received February 3, 2001

The last couple of years an overwhelming amount of data has emerged in the field of biomolecular structure determination. To explore information hidden in these structure databases, clustering techniques can be used. The outcome of the clustering experiments largely depends, among others, on the way the data is represented; therefore, the choice how to represent the molecular structure information is extremely important. This article describes what the influence of the different representations on the clustering is and how it can be analyzed by means of a dendrogram comparison method. All experiments are performed using a data set consisting of RNA trinucleotides. Besides the most basic structure representation, the Cartesian coordinates representation, several other structure representations are used.

INTRODUCTION

Due to the recent developments in computer technology and the increasing popularity of the Internet, the size of publicly available biomolecular structure databases has increased considerably.¹ To facilitate the assessment of information in these large databases, many different techniques are used. Some techniques try to visualize the relationships between the 3D-objects by projecting the structures of interest in a low dimensional space. This reduces the dimensionality of the data and is done either by defining new variables describing some global features of the structures² or by making linear combinations of the original variables.^{3,4} Also clustering techniques can be used to explore the underlying structure of the data.^{5,6} The resulting clusters can reveal interesting information about the structural behavior of the considered molecules. Which structures are grouped, and how does this relate to correspondences/differences in chemical behavior? Especially for biomacromolecules which have certain preferred molecular conformations (nucleic acids: A-DNA, B-DNA; proteins: α -helix, β -sheets), this is thoroughly investigated.

Clustering works by seeking groups of objects that form natural clusters by examining the pairwise similarities systematically.^{7,8} The clustering results are defined by the agglomeration strategy, the similarity measure, and the representation of the data.⁹ In this article the influence of different structure representations on the outcome of clustering experiments is examined. Many different representations can be used to describe the structural information in the data. In the biomacromolecular databases that are accessible through the Internet (e.g. the NDB¹⁰ (URL: ndbserver.rutgers.edu/) and PDB¹¹ (URL: www.rcsb.org/pdb/) databases for nucleic acid and protein structures, respectively) the Cartesian coordinates representation is used to fix the spatial position of the atoms within the molecule. The

Cartesian coordinates representation is the most basic representation, and other representations can be deduced from this basic representation. A disadvantage of the Cartesian coordinates representation is the rather large number of variables, even for relative small molecules. When clustering experiments are performed, an additional disadvantage of the coordinates representation shows up. Before different structures can be clustered, they have to be aligned in order to determine their (dis)similarity. Therefore, instead of Cartesian coordinates, internal coordinates (torsion angles) are often used to represent molecular structures. Not only does this lead to a coordinate system that is independent of the location and orientation of the molecule, it also leads to a significant reduction in the number of variables. This is justified because bond lengths and angles normally do not change significantly in the structures. Besides the torsion angle representation, several other representations can be applied such as application-dependent measures like pseudo-torsion angles and distance matrices.

In this article, the effect of these different representations is examined on the outcome of the clustering results by using a dendrogram comparison technique developed by Fowlkes and Mallows.¹² Both disadvantages and advantages of the usage of these structure representations are discussed on the basis of clustering experiments concerning a database consisting of RNA trinucleotides (of course other biomacromolecule databases could be used as well).

EXPERIMENTAL SECTION

Data. The starting point of all calculations is a data set containing molecular structure information of 121 RNA trinucleotides. It is modeled after part of the data set described by Duarte and Pyle.² On the basis of groupings of similar objects in the so-called pseudotorsion angles space, Duarte and Pyle defined eight RNA classes: helical, C2 bend, base twist, chi switch, cross strand stack, flip turn, stacked turn, and stack switching classes. To make the data

* Corresponding author phone: +31-24-3653192; e-mail: L.Buydens@sci.kun.nl.

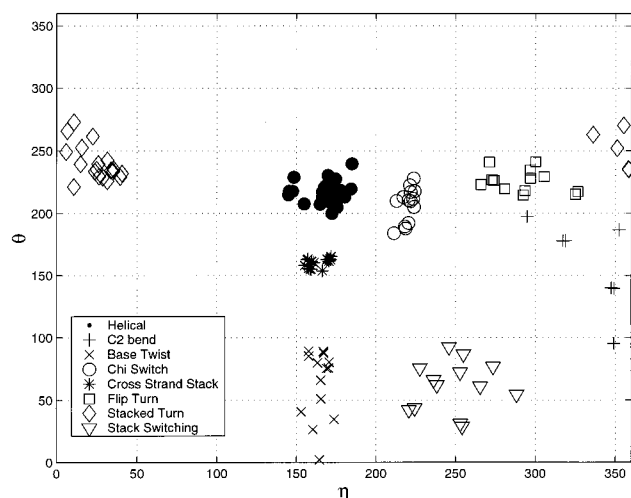


Figure 1. An overview is given of a sub set of the complete molecular data set of Duarte and Pyle of the RNA structures in the η and θ pseudotorsion space. Eight different RNA classes are depicted.

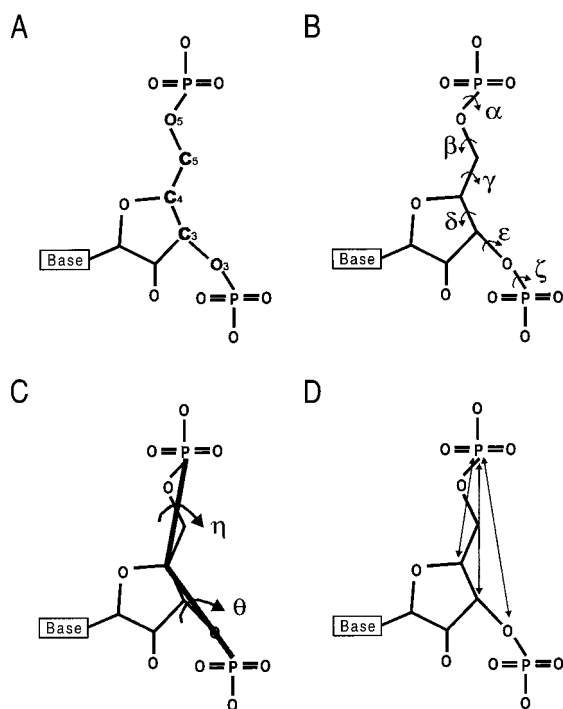


Figure 2. Overview of the four different representations used to fix the molecular structure information of RNA trinucleotides. A: Cartesian coordinates representation, B: torsion angles representation, C: pseudotorsion angles representation, D: distances representation.

set more manageable the number of objects in the original database is reduced. To remove any bias for a certain class (over 70% of the original data set contains helical objects), of each class at least five objects are taken along in the data set. Figure 1 depicts the eight different RNA classes in the pseudotorsion angles space.

The structure information is fixed by means of four different representations (see Figure 2). For the first representation, each trinucleotide is represented by the Cartesian coordinates of the backbone atoms P, O5, C5, C4, C3, and O3 (18 variables \times 3 = 54 variables). This representation is considered to be the golden standard or reference representation, since all others can be derived from it. In

Table 1. For the Four Different Structure Representations, an Overview Is Given of the Number of Variables Used To Represent Structure Information of a RNA Trinucleotide Backbone

structure representation	no. of variables
Cartesian coordinates	54
torsion angles	18
pseudotorsion angles	2
distances	153

the second representation the 3D-conformation of the trinucleotides is expressed by the backbone torsion angles (α , β , γ , δ , ϵ , and ζ) of each mononucleotide (6 variables \times 3 = 18 variables). The third representation is based on the pseudotorsion angles, η and θ , defined by Duarte and Pyle (2 variables).² Finally, a representation is used where the molecular structure information is fixed by means of listing all existing distances between the backbone atoms (P, O5, C5, C4, C3, and O3) in the trinucleotide structure ($\sum_{n=1}^{18} (n-1) = 153$ variables). Table 1 gives an overview of the different structure representations.

METHODS

Clustering. The first step in a clustering is the calculation of a dissimilarity matrix. Because four different representations are used, the calculation of the dissimilarities between the objects in the data set also differs. The calculation of the distance matrix for the Cartesian coordinates representation data set is the most computer-intensive because all RNA structures have to be optimally translated and rotated before a dissimilarity measure can be determined. After alignment of two structures by means of Procrustes analysis,^{13,14} the dissimilarity is given by the summation of the distances between the atoms. Because the torsion angles are circular, the maximum difference between two torsion angle values is half the range used. Distances larger than this value are mapped back into the permitted range (this is applied to both regular and pseudotorsion angles). Dissimilarities for the distances-based representation are obtained by simply summing the squared differences between the listed distances, the Fröbenius norm. Ward's method is used to cluster the objects on the basis of their dissimilarities.¹⁵ Of course other clustering algorithms may be used as well.

Comparison of the Clusterings. To express mathematically which representation gives the most similar dendrogram to the Cartesian coordinates dendrogram (the reference clustering), the dendrograms are compared with each other using the B_k value defined by Fowlkes and Mallows.¹² To calculate B_k for two dendrograms, they are split up into k clusters, and the similarity between the clusters of the different dendrograms is determined. For that, a matching matrix $[m_{ij}]$ is filled ($i = 1, \dots, k$; $j = 1, \dots, k$), where m_{ij} is the number of objects in common between the i th cluster of dendrogram 1 and the j th cluster of dendrogram 2. When all values in the matching matrix are added, the total number of objects, n , should be obtained. The distribution of n over the matching matrix gives an indication of how well the dendrograms are associated with each other. For similar dendrograms the matching matrix will contain some relative high values and many zeros. For dissimilar dendrograms, objects that are grouped in one cluster in dendrogram 1 will be scattered over different clusters in the other dendrogram.

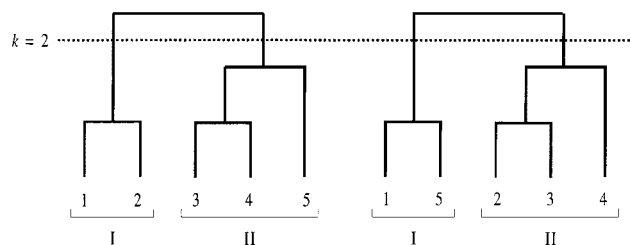


Figure 3. Two dendrograms for which the B_2 value is determined.

As a result, the matching matrix will contain many low values; the objects are scattered over the matching matrix instead of being concentrated in certain cells of the matrix. Given the matching matrix $[m_{ij}]$, B_k can be calculated as follows

$$B_k = T_k / \sqrt{P_k Q_k}$$

where

$$T_k = \sum_{i=1}^k \sum_{j=1}^k m_{ij}^2 - n$$

$$P_k = \sum_{i=1}^k m_{i\cdot}^2 - n$$

with

$$m_{i\cdot} = \sum_{j=1}^k m_{ij}$$

and

$$Q_k = \sum_{j=1}^k m_{\cdot j}^2 - n$$

with

$$m_{\cdot j} = \sum_{i=1}^k m_{ij}$$

Ultimately, for exactly identical clusterings, B_k gets a value of 1, while for completely different clusterings, a minimal value of 0 is obtained. To clarify the formulas given above, the B_2 value for the dendrograms in Figure 3 is calculated ($k = 2$ and $n = 5$). First the matching matrix is composed:

$$M = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

Notice that the summation of all matrix elements equals 5, the total number of objects. With

$$m_{i\cdot} = [2 \quad 3]$$

$$m_{\cdot j} = [2 \quad 3]$$

$$P_k = 2^2 + 3^2 - 5 = 8$$

$$Q_k = 2^2 + 3^2 - 5 = 8$$

$$T_k = 1^2 + 1^2 + 1^2 + 2^2 - 5 = 2$$

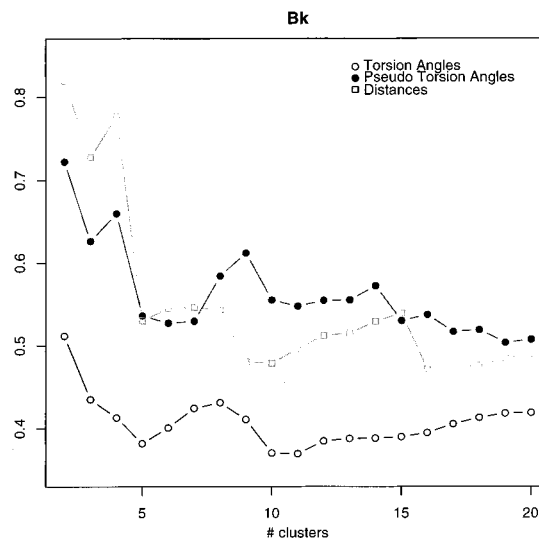


Figure 4. B_k values for the torsion angle (open circles), the pseudotorsion angle (closed circles), and the distances representation (open cubes). All dendrograms are compared with the Cartesian coordinates dendrogram.

B_2 becomes equal to 0.25, indicating that the dendrograms from Figure 3 are quite dissimilar. When the different values of k are plotted against the corresponding B_k values, the similarity of the two clusterings can be followed during the whole clustering process.

All calculations in this article are performed using Matlab for Unix Workstations, version 5.3, by Mathworks (URL: www.mathworks.com/) and R, version 1.1.1 (URL: www.r-project.org/).¹⁶

An R library for calculating similarities between partitionings is available at the website of the author (URL: www.sci.kun.nl/cac/people/rwehrens/software/compare-hc.tar.gz).

RESULTS AND DISCUSSION

To examine the influence of the different representations on the clustering of the molecular structure data set, all results are compared with a reference clustering, the Cartesian coordinates representation, using the B_k value defined by Fowlkes and Mallows.

In Figure 4 the B_k values are depicted for the torsion angle (open circles), the pseudotorsion angle (closed circles), and the distances (open cubes) representation. For each clustering level k (the horizontal axis) the similarity with the Cartesian coordinates dendrogram at the same clustering level (the vertical axis) is calculated. At all clustering levels the similarity of the torsion angle dendrogram with the coordinates dendrogram is the smallest. The pseudotorsion angle and distances dendrograms give relatively high B_k values, especially for small k 's. This indicates that at the start of the clustering, the splitting up of the objects into clusters for these representations and the coordinates representation is roughly the same. When the total number of objects of the data set is grouped into five clusters, the similarity with the reference dendrogram decreases dramatically. Beyond this point the similarity of both torsion angle based representations increases, reaching a second optimum at $k = 8$ and 9. While the similarity of the regular torsion angle representation slowly increases for larger k , the B_k values

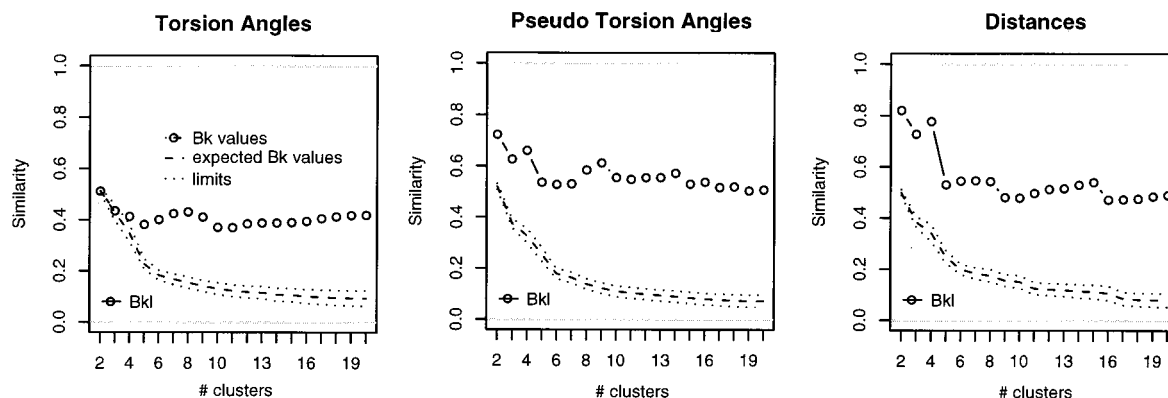


Figure 5. Calculated and expected B_k values for the torsion angles representation, the pseudotorsion angles representation and the distances-based representation.

for the pseudotorsion angle dendrogram slowly decreases. The clustering behavior of the distances-based representation after k equals 5 is the opposite of the behavior of the pseudotorsion angle representation. At the point where the pseudotorsion angle representation reaches its second optimum (k is 9), the distances representation is situated in a local minimum. For larger k the similarity slowly increases, until another drop is observed at k is 16.

To examine if the clusterings from Figure 4 are significantly different from a random clustering, for each representation at each cluster strength also expected B_k values are calculated under the null hypothesis of no cluster structure (formulas to calculate expected B_k values with confidence limits are given in ref 12).

In Figure 5 for all representations both the calculated B_k values (open circles) as well as the expected B_k values with its confidence limits (dotted lines) are plotted. For the most dissimilar representation, the torsion angle representation, for k equals 4 the real B_k curve rises above the expected B_k values curve, indicating that at that stage of the clustering the dendrogram has become significantly dissimilar from the randomized dendrogram and shows some agreement with the reference classification. For the other representations, for all k 's the corresponding clusterings differ significantly from a random clustering.

As already is described in the Experimental Section, the molecular data set is composed on the basis of eight classes defined by Duarte and Pyle. Besides the local optima for the torsion angle based representations and the sudden decrease of the similarity for the distances representation round k is 8, no clear evidence can be found that the data set indeed is composed of objects that can be split up into eight distinct classes.

In Figure 6 both the Cartesian coordinates (6A) and the pseudotorsion angles (6B) dendrograms are depicted. The objects are labeled according to the eight classes defined in Duarte and Pyle (defined in pseudotorsion space). In the pseudotorsion angle dendrogram clearly eight clusters can be discerned, each containing mainly objects of one of the eight different classes. One C2 bend RNA trinucleotide (labeled 1) is placed in the cluster with the flip turn objects (labeled 5). The coordinates dendrogram reflects to some extent clustering for the helical, base twist, stacked turn, stacked switching, and the flip turn objects (objects labeled with 6, 2, 7, 8, and 5). Also a large cluster exists containing mainly chi switch and cross strand stack objects (labeled 3

and 4). The C2 bend objects (1) are the only trinucleotides that are not grouped in the coordinates dendrogram. The existence of clusters for the other classes clearly demonstrate that some of the class definitions made by Duarte and Pyle hold. However, because Duarte and Pyle have defined the classes according to the properties of the majority of the objects in a certain area of the pseudotorsion angles plot, objects exist that were falsely assigned to a certain class. For example, according to Duarte and Pyle only 80% of the nucleotides in the flip turn region truly apply to the definitions. This is also seen in the coordinates dendrogram. Several objects, labeled according a certain class, are not located in the cluster containing the majority of the objects of that certain class. For example, not all objects labeled with 7 are clustered in the stacked turn cluster. Three stacked turn objects are clustered somewhere else.

To get an indication what the influence of the different representations is on the clustering, the mean and variance of the structures of each cluster are considered at a clustering level of 8. Figure 7 shows the mean structures (the backbone atoms of the trinucleotides) of the eight clusters of the coordinates dendrogram. Only the x- and y-coordinates of each atom are depicted. The variance of the xyz-coordinates of each atom is visualized by means of the magnitude of the point that correspond with the mean position of the atom. In almost all mean structures the larger size of the atoms at the start and end of the structure indicates that there is much more flexibility in that part of the structure than in the rest of the molecule. For the structures of cluster 1 and 4 the variances of all atoms are relatively small. This implies that the structures in these clusters are very similar and can easily be discerned from the other structures in the data set. When the clustering level is increased, these clusters are the last eligible for further splitting up into several new clusters.

The plots in Figure 8 correspond to the mean structures of the clusters where the torsion angle representation is used. The order in which the structures are plotted in the figure is, as much as possible, based on the order of the structures in the previous figure. The first mean structure in Figure 8 is most similar to the first mean structure of Figure 7, etc. Between brackets the summed distances between the atoms of the mean structures in Figures 7 and 8 are given. Again, the variances at the tips of the molecules are larger than the variances in other parts of the structure but in comparison with the structures of the previous figure the variances are now much bigger. Despite the large difference in variances,

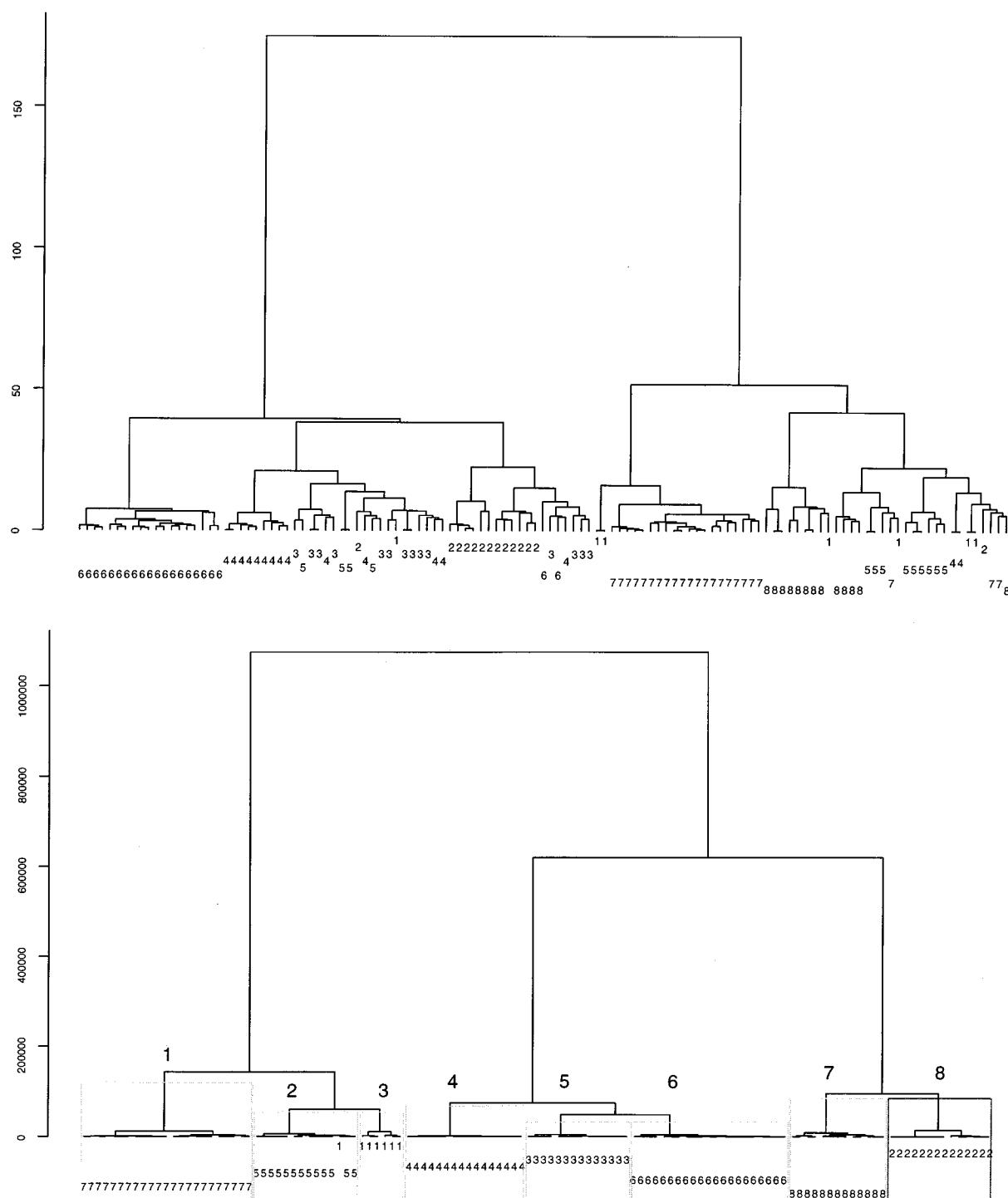


Figure 6. A: Cartesian coordinates dendrogram. The objects are labeled according to the class definitions of Duarte and Pyle: 1 = C2 bend, 2 = base twist, 3 = chi switch, 4 = cross strand stack, 5 = flip turn, 6 = helical, 7 = stacked turn, 8 = stacked switching. Figure 6. B: Pseudotorsion angles dendrogram. Cartesian coordinates dendrogram. The objects are labeled according to the class definitions of Duarte and Pyle: 1 = C2 bend, 2 = base twist, 3 = chi switch, 4 = cross strand stack, 5 = flip turn, 6 = helical, 7 = stacked turn, 8 = stacked switching.

some structures in Figure 8 look very similar to the structures in the Figure 7. Particularly, the mean structures of cluster 1 and 4, per chance the structures with the smallest overall variance in Figure 8, resemble the corresponding structures of the previous figure. Parts of the remaining structures are very similar to the reference structures. The torsion angle representation probably focuses on particular structural motifs in the RNA nucleotides and leaves the other parts undisturbed.

Figure 9 visualizes the mean structures of the pseudotorsion angle representation. In reference to variances, the same conclusions can be drawn as for the regular torsion angle representation: more flexibility at the tips of the structures and larger variances than is the case for the reference structures. However, the choice to describe the structure more globally by means of pseudotorsion angles instead of regular torsion angles resulted in structures much more similar to the Cartesian coordinates structures. Besides structures 1 and

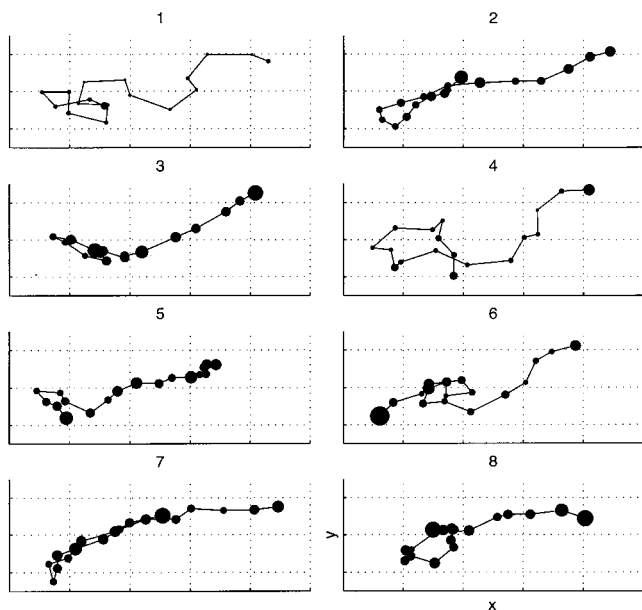


Figure 7. Mean structures of the eight clusters in the Cartesian coordinates dendrogram. The variance of the position of an atom in a particular cluster is indicated by the size of the plotting symbol.

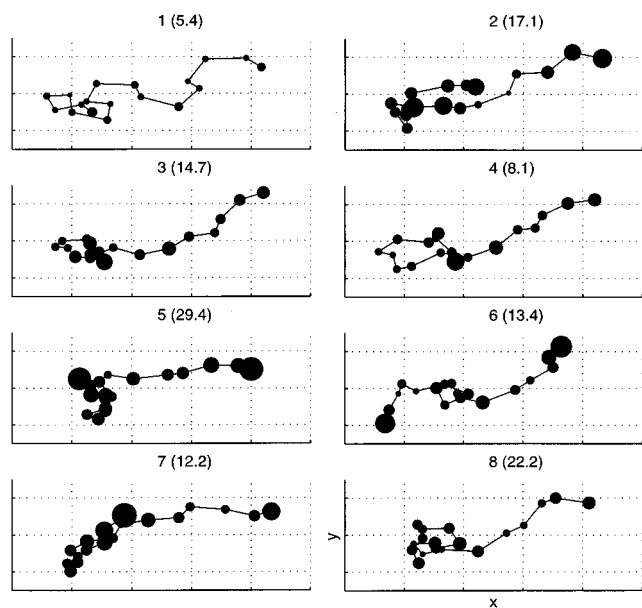


Figure 8. Mean structures and variances of the eight clusters from the torsion angles dendrogram.

4, also structures 5, 6, and 8 resemble the corresponding structures in Figure 7. The usage of the pseudotorsions instead of regular torsions does not have an effect on all the mean structures of the clusters; the mean structure of cluster 7 in Figure 9 does not differ much from the structure of cluster 7 in Figure 8.

Figure 10 shows the mean structures of the eight clusters of the distance-based representation. Similar variance values are found as in the Cartesian coordinates representation. Once more, analogous structures to the reference structures of cluster 1 and 4 can be detected. Additionally the mean structures of cluster 2 and 7 in Figure 10, the two most elongated structures, look very similar to their counterparts in Figure 7. This is probably caused by the nature of the used representation. Contrary to the other representations,

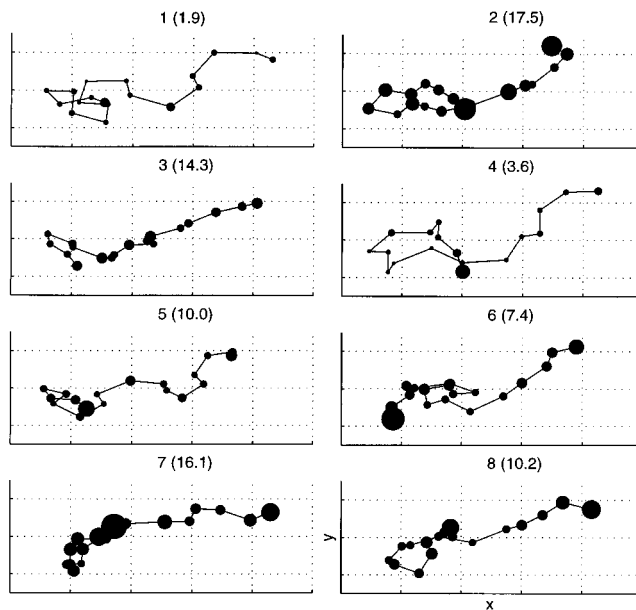


Figure 9. Mean structures and variances of the eight clusters from the pseudotorsion angles dendrogram.

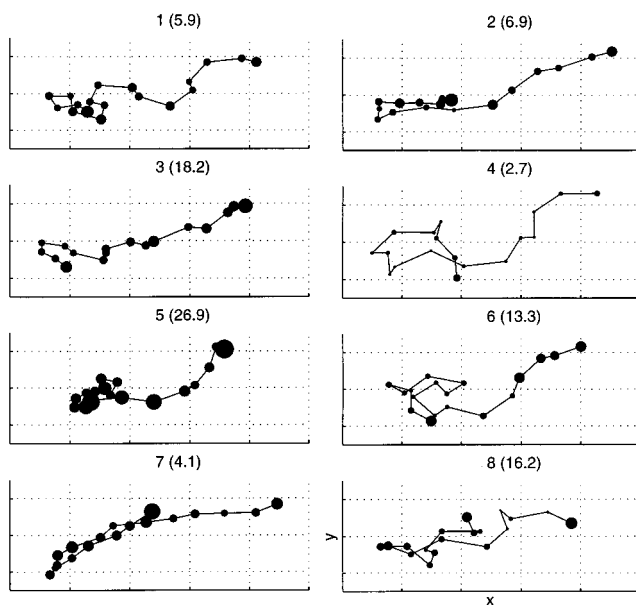


Figure 10. Mean structures and variances of the eight clusters from the distances-based dendrogram.

in the distance-based representations the relative importance of structure features between remote atoms is much bigger.

Overall, the mean structures and variances of Figures 7–10 confirm the conclusions made on the basis of Figure 4. For the torsion angle representation only two very similar structures are found (mean structures 1 and 4), while for the other representations besides these structures also other similar mean structures are found. Ultimately this is expressed by the larger B_k values in Figure 4. The appearance of similar structures to structures 1 and 4 for all different representations may give an indication of the robustness of these classes.

CONCLUSIONS

Many different ways exist of representing molecular structure information. All have a different effect on the

outcome of the clustering experiments. By means of Fowlkes' and Mallows' B_k values, the influence of these different representations on the clusterings can be analyzed. The results in this article show that the torsion angle representation dendrogram differs substantially from the Cartesian coordinates representation dendrogram. The pseudotorsion and distances-based representations give groupings of molecular structures of RNA trinucleotides that look more similar to the groupings using the Cartesian coordinates representation. When the objects under investigation contain elongated structures, the distances-based representation appears to approach the Cartesian coordinates clustering the best.

As shown in this paper, the representation employed to describe molecular structure may have a profound effect on the results of the clustering. To obtain meaningful results, one should carefully consider which representation to use; some representations stress global features (such as distance-based representations, where the overall structure of an object is described) or more local features (e.g. a torsion angle description, where the neighborhoods of several individual chemical bonds are described). A major problem is that it is not always possible to state a priori which representation is best for a particular application. Validation of the results therefore is of prime importance.

REFERENCES AND NOTES

- (1) Wehrens, R.; de Gelder, R.; Kemperman, G. J.; Zwanenburg, B.; Buydens, L. M. C. Molecular challenges in modern chemometrics. *Anal. Chim. Acta* **1999**, *400*, 413–424.
- (2) Duarte, C. M.; Pyle, A. M. Stepping Through an RNA Structure: A Novel Approach to Conformational Analysis. *J. Mol. Biol.* **1998**, *284*, 1465–1478.
- (3) Beckers, M. L. M.; Buydens, L. M. C. Multivariate analysis of a data matrix containing A-DNA and B-DNA dinucleotides. Multidimensional Ramachandran plots for nucleic acids. *J. Comput. Chem.* **1998**, *19*, 695–715.
- (4) Buydens, L. M. C.; Reijmers, T. H.; Beckers, M. L. M.; Wehrens, R. Molecular data-mining: a challenge for chemometrics. *Chemometrics. Intell. Lab. Sys.* **1999**, *49*, 121–133.
- (5) Donate, L. E.; Rufino, S. D.; Canard, L. H. J.; Blundell, T. L. Conformational analysis and clustering of short and medium size loops connecting regular secondary structures: A database for modeling and prediction. *Prot. Sci.* **1996**, *5*, 2600–2616.
- (6) Lin, T.; Lin, J.; Huang, Y.; Liu, J. Clustering Peptide Structures through Identification of Commonly Exposed Groups. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 622–629.
- (7) Barnard, J. M.; Downs, G. M. Clustering of Chemical Structures on the Basis of Two-Dimensional Similarity Measures. *J. Chem. Inf. Comput. Sci.* **1992**, *32*, 644–649.
- (8) Kaufman, L.; Rousseeuw, P. J. *Finding Groups in Data. An Introduction to Cluster Analysis*; Wiley: New York, 1989.
- (9) Jurs, P. C. Chemometrics and Multivariate Analysis in Analytical Chemistry. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH Publishers: New York, 1990; pp 169–212.
- (10) Berman, H. M.; Olson, W. K.; Beveridge, D. L.; Westbrook, J.; Gelbin, A.; Demeny, T.; Hsieh, S. H.; Srinivasan, A. R.; Schneider, B. The Nucleic Acid Database: a comprehensive relational database of three-dimensional structures of nucleic acids. *Biophys. J.* **1991**, *63*, 751–759.
- (11) Bernstein, F. C.; Koetzle, T. F.; Williams, G. J. B.; Meyer, E. F.; Brice, M. D.; Rogers, J. R.; Kennard, O.; Shimanouchi, T.; Tasumi, M. The Protein Data Bank: a computer-based archival file for macromolecular structures. *J. Mol. Biol.* **1977**, *112*, 535–542.
- (12) Fowlkes, E. B.; Mallows, C. L. A Method for Comparing Two Hierarchical Clusterings. *J. Am. Stat. Assoc.* **1983**, *78*, 553–569.
- (13) Gower, J. C. Generalised Procrustes analysis. *Psychometrika* **1975**, *40*, 33–50.
- (14) ten Berge, J. M. F. Orthogonal procrustes rotation for two or more matrixes. *Psychometrika* **1977**, *42*, 267–276.
- (15) Ward Jr., J. H. Hierarchical Grouping to Optimize an Objective Function. *J. Am. Stat. Assoc.* **1963**, *58*, 236–244.
- (16) Ihaka, R.; Gentleman, R. A language for data analysis and graphics. *J. Comput. Graph. Stat.* **1996**, *5*, 299–314.

CI0103626