

Note

Rapid Calculation of Coordinates from Distance Matrices

1. INTRODUCTION

One of the major difficulties in the calculation of conformation by the "distance geometry" approach [1] has been that of producing a set of three dimensional Cartesian coordinates v_i , $i = 1, \dots, n$ for n points, given an $n \times n$ matrix of upper bounds U and lower bounds L on the interpoint distances:

$$l_{ij} \leq d_{ij} = \|v_i - v_j\| \leq u_{ij}, \quad \text{for } i, j = 1, \dots, n. \quad (1)$$

We are not concerned here with the determination of U and L , but rather with the calculation of any set of v 's satisfying (1), assuming that a solution exists. A proposed distance matrix D chosen at random such that $l_{ij} \leq d_{ij} \leq u_{ij}$ for all i and j , does not correspond in general to any set of v 's in the three dimensional Euclidean space because the corresponding $(n + 1) \times (n + 1)$ bordered matrix of squared distances, C , is in general not of rank 5 (see [1] for proof). Here $c_{ij} = d_{ij}^2$ for $i, j = 1, \dots, n$ and $c_{i, n+1} = c_{n+1, i} = 1$ for $i = 1, \dots, n$ and $c_{n+1, n+1} = 0$. Furthermore, generating the v 's analytically according to (6) of [1] directly from D (or C) results in increasingly larger violations of (1) for larger i and j . Direct alteration of C (and hence D) in order to reduce its rank to 5, permits the straightforward calculation of v 's from the modified D , but the modification algorithm in [1] (involving Gaussian elimination and direct alteration of individual matrix elements) required computer time proportional to n^3 and memory space proportional to n^2 .

One can devise a more efficient method, time and memory requirements depending only linearly on n , by resolving the real, symmetric matrix C into its eigenvalues and eigenvectors:

$$C = E\Lambda E^T \quad (2)$$

where $E = (e_{ij})$ is the column matrix of eigenvectors and $\Lambda = (\lambda_i)$ is the diagonal matrix of corresponding eigenvalues. The λ 's are real and may be ordered so that $|\lambda_1| > |\lambda_2| > \dots$; the eigenvectors are real and may be taken as orthonormal [2]. The condition that C have rank 5 is equivalent to requiring $|\lambda_5| > \lambda_6 = \lambda_7 = \dots = \lambda_{n+1} = 0$. Thus we may calculate a modified C matrix, denoted by C' , from (2) after setting the sixth and subsequent eigenvalues to zero:

$$c'_{ij} = \sum_{k=1}^5 e_{ik} \lambda_k e_{jk} \quad \text{for } i, j = 1, \dots, n + 1 \quad (3)$$

Then C' is the rank 5 matrix closest to C in the spectral sense, but some elements may violate (1), some diagonal elements may be non-zero, and some border elements may not be equal to unity. We have tried a number of iterative schemes to correct these discrepancies, but in our hands they either converged very slowly or not at all. In the examples we have tried, C is well-conditioned with respect to its eigenvalues but ill-conditioned with respect to its eigenvectors. By numerically minimizing the above deviations using as variables the 5 eigenvalues of largest modulus and the components of their corresponding eigenvectors, one can avoid the instability difficulties, but the constraints are awkward to handle. There are [5 variable eigenvalues] + [5($n + 1$) variable components of the five eigenvectors] - [($n + 1$) constraints of zero diagonal elements] - [n constraints of unity border elements] - [15 constraints from the orthonormality of the five eigenvectors] = $3n - 6$ net degrees of freedom, which is what one would expect, since we are seeking the x , y and z coordinates of n points, up to a rigid translation and rotation (and mirror inversion). C' is automatically symmetric, as can be seen from (3).

2. ALGORITHM

Our most successful method for large n consists of: (i) choosing D at random, bounded element-wise by U and L as described; (ii) calculate the first 5 eigenvalues and eigenvectors of the corresponding C matrix; (iii) calculate C' from (3); (iv) calculate the v 's from the first four rows of C' according to (6) of Ref. [1]; and (v) minimize $f(v_1, \dots, v_n)$ with respect to the $3n$ Cartesian coordinates

$$f(v_1, \dots, v_n) = \sum_{j>i} \begin{cases} (d_{ij}^2 - u_{ij}^2)^2, & d_{ij} > u_{ij} \\ 0, & l_{ij} \leq d_{ij} \leq u_{ij} \\ (d_{ij}^2 - l_{ij}^2)^2, & d_{ij} < l_{ij} \end{cases} \quad (4)$$

where $d_{ij} = \|v_i - v_j\|$, the ordinary Euclidean norm. Step (ii) is conveniently done by the method of "exhaustion," where the eigenvalue of largest absolute value, λ_1 , and its corresponding eigenvector e_1 are found by iterating for $k = 1, 2, \dots$

$$\frac{(C'^k y) \cdot (C'^k y)}{(C'^{k-1} y) \cdot (C'^k y)} \rightarrow \lambda_1, \quad \text{as } k \rightarrow \infty \quad (5)$$

and $C'^k y \rightarrow e_1$, for arbitrary y .

Then λ_2 is found similarly by using $C' - \lambda_1 e_1 e_1^T$ in the place of C' in (5), and so on [3]. Convergence of the λ 's to six significant figures is regularly achieved well before $k = 100$, even when two eigenvalues are close in magnitude. The computational effort involved in the matrix-vector multiplication goes up as n^2 , but since only the first 5 eigenvalues are sought, the cost otherwise depends linearly on n . Step (iv) is remarkably stable. For instance, in our calculations of conformations of the small protein, bovine pancreatic trypsin inhibitor (BPTI), $n = 58$, $d_{i,i+1} = 3.8 \text{ \AA}$, and $6.0 \text{ \AA} \leq d_{5,55}, d_{14,38}, d_{30,51} \leq 6.5 \text{ \AA}$. In other words, U and L enforce a chain of 57 fixed length steps with three short crosslinks spanning nearly the whole chain.

Other distances are much less constrained. Calculation of the v 's from C would result in $d_{5,55} \sim 100$ and so on, whereas using C' we regularly achieve crosslink distances between 5 and 10. Thus step (iv) initially positions the points in an overall correct fashion, although with still some substantial errors. Step (v) is a nonlinear least squares minimization of the function f , which has zero value at any solution, strictly positive values elsewhere, and continuous first derivatives everywhere. The gradient can easily be calculated analytically, so we use initially steepest descent minimization until $\|\nabla f\|^2 < 10^5 \text{ \AA}^6$ and then Fletcher-Reeves conjugate gradient minimization [4] until $\|\nabla f\|^2 < 10^3 \text{ \AA}^6$, resulting in $f < 1 \text{ \AA}^4$ and a largest term in (4) of 0.01 \AA^4 or less. Typical initial values of f are about 10^{10} \AA^4 . Both minimization algorithms have storage requirements proportional to n , so large numbers of variables are easily handled. As a test of the method, we chose $u_{ij}^2 = 1.05g_{ij}^2$ and $l_{ij}^2 = 0.95g_{ij}^2$, where the g_{ij} are the distances between C^α atoms of BPTI taken from the x -ray crystallographic coordinates [5]. On the Lawrence Berkeley Laboratory's CDC 7600, each conformation generated according to the algorithm above took an average of 2.7 seconds. By way of comparison, minimizing from arbitrary starting points with the same U and L takes about 5 times as long [6]. The performance with looser upper and lower bounds is similar. The starting configuration in step (iv) is close enough to the minimal region of f that convergence to spurious local minima rarely occurs, in our experience. This is in sharp contrast to the difficulties experienced when using arbitrary starting points [6].

The method also performs well when the upper and lower bound matrices are much less restrictive. When U and L correspond to only fixed i to $i + 1$ distances and the three crosslinks of BPTI (plus their logical consequences regarding all other distances, as deduced from the triangle inequality [1]) then only 1 second of computer time per structure is required. Convergence is about as rapid and precise as with tight bounds, and we have encountered no spurious local minima. Of course there is a great deal of conformational variability among the resultant structures since the bounds are so loose.

3. SUMMARY

In conclusion, we have devised and tested a practical algorithm for rapidly solving (1) in the difficult case of large n and numerous strong constraints. Computer time increases only quadratically and memory requirements increase only linearly with n , and there is little difficulty with multiple minima. We hope this will make the distance geometry approach to conformational calculation more feasible for large, highly constrained systems.

ACKNOWLEDGMENTS

This work was supported by a grant from the Academic Senate of the University of California. The author thanks Drs. I. D. Kuntz, K. Miller, and B. N. Parlett for their stimulating and helpful discussions.

REFERENCES

1. G. M. CRIPPEN, *J. Computational Phys.* **24** (1977), 96.
2. J. H. WILKINSON, "The Algebraic Eigenvalue Problem," Oxford Univ. Press, Oxford, 1965.
3. D. K. FADDEEV AND V. N. FADDEVA, "Computational Methods of Linear Algebra," pp. 307, 328-330, Freeman, San Francisco, 1963.
4. R. FLETCHER AND C. M. REEVES, *Comput. J.* **7** (1964), 149.
5. R. HUBER *et al.*, *Cold Spring Harbor Symp. Quant. Biol.* **36** (1971), 141.
6. I. D. KUNTZ, G. M. CRIPPEN, AND P. A. KOLLMAN, submitted for publication (1978).

RECEIVED: May 3, 1977

G. M. CRIPPEN

*University of California
School of Pharmacy,
San Francisco, California 94143*