12208

Awakening to reality Available online at www.elixirpublishers.com (Elixir International Journal)

# **Discrete Mathematics**





# Mathematical representation of high-resolution protein structures using graph

theory

M.El-Ghoul<sup>1</sup>, A.El-Guoshy<sup>2</sup>, F.El-Fiki<sup>2</sup> and A.El-Refy<sup>2</sup> <sup>1</sup>Department of Mathematics, Faculty of Science, Tanta University. <sup>2</sup>Biotechnology Department, Faculty of Agriculture, Al-Azhar University.

## ARTICLE INFO

Article history: Received: 24 November 2011; Received in revised form: 22 November 2012; Accepted: 10 December 2012;

# ABSTRACT

In this paper we developed new matrices (angle matrix ,connection matrix and connection angle matrix) in addition to a set of matrices (adjacent matrix, incidence matrix) that previously used to introduce the mathematical representation of high-resolution protein structures an accurate representatin that facilitate analysis of its structures.

© 2012 Elixir All rights reserved.

# Keywords

Protein, Amino acid, Graph theory, Matrices.

### Introduction

#### Definition and background:-

- Abstract graphs: An abstract graphs G is a diagram consisting of a finite non empty set of the elements, called "vertices" denoted by V (G) together with a set of unordered pairs of these elements, called "edges" denoted by E (G). The set of vertices of the graph G is called "the vertex set of G" and the list of edges is called "the edge – list of G" [Gibbons, 1995; Giblin, 1977].

- **Adjacency and incidence:** let v and w be vertices of a graph. If v and w are joined by an edge e. then v and w are said to be adjacent. Moreover, v and w are said to be incident with e, and e is said to be incident with v and w [Wilson, 1972].

- The adjacency matrix: let G be a graph without loops, with n-vertices labeled 1, 2, 3, . . . , n. The adjacency matrix A(G) is the nxn matrix in which the entry in row i and column j is the length of edge in angstroms if the vertices i and j are joining and 0 otherwise. [Modification: Wilson, 1972].

- The incidence matrix: let G be a graph without loops, with n-vertices labeled 1, 2, 3,  $\ldots$  n and m edges labeled 1, 2, 3,  $\ldots$ , m. the incidence matrix I (G) is the n x m matrix in which the entry in row i and column j is l if vertex i is incident with edge j and 0 otherwise [Gross, 1987; Wilson, 1990].

- **Angle matrix:** let G be a graph without loops, with n-vertices labeled 1, 2, 3,  $\ldots$  n and m edges labeled 1, 2, 3,  $\ldots$ , m. the angle matrix g(G) is the m x m matrix in which the entry in row i and column j is the angle in degrees if edge i and j are incident with the same vertex and 0 otherwise.(New matrix)

- **Connection matrix:** let g1 be a directed subgraph without loops, with n1-vertices labeled 1, 2, 3. . . n1 and g2 be another directed subgraph without loops, with n2-vertices labeled 1, 2, 3, . . . ., n2. The connection matrix C (g1, g2) is the n1xn2 matrix in which the entry in row i and column j is the length of edge in angstroms if the vertices i and j are joining and 0 otherwise. [New matrix].

- Connection Angle matrix: let g1 be a directed weighted subgraph without loops, with n1-vertices labeled  $1, 2, 3, \ldots, n1$ ,

Tele: E-mail addresses: hendelmorsy@yahoo.com m1- edges labeled 1, 2, 3...m1, g2 be another directed weighted subgraph without loops, with n2-vertices labeled 1, 2, 3...n2, m2- edges labeled 1, 2, 3...m2 and E is a directed connector between g1, g2. the connection angle matrix CL(G) is the m1 x E matrix or E x m2 in which the entry in row i and column j is the angle in degrees if edge i and j are incident with the same vertex and 0 otherwise.(New matrix)

## Main results

Now we will represent the Mean values of Main-chain bond lengths and bond angles extracted from PROCHECK Operating Manual Appendix A Stereochemical parameters TABLE A.1

(http://www.ebi.ac.uk/thornton

srv/software/PROCHECK/manual/manappa.html) as observed in small molecules (Engh & Huber, 1991) that derived from high-resolution protein structures using graph matrices (Figure 1).

## Definition

Representation of high-resolution protein structure backbone using a directed weighted graph (G) in consideration atoms as vertices and chemical bonds between atoms as edges in 2 steps. **The First step:** 

### The First step:-

Representation of each amino acid residue backbone (except glycine and proline) as undirected weighted sub graph (g) all in protein primary structure as a directed weighted graph (G) using 2 matrices.

A	(g)	all	=
	$\langle O \rangle$	an	

Ν Cα C O	N	Сα	С	0
	0	1.458	0	0
	1.458	0	1.525	0
	0	1.525	0	1.231
	0	0	1.231	0

Where N, C $\alpha$ , C, O represent the vertices of subgraph (g) <sub>all</sub> and values in cells represent the length of edges in angstrom. I (g) <sub>all</sub> =

$N\ C\alpha$	Ν Cα	CaC 0	CO
CO	1	110	0
	1		0
			1
	0		1
	0		
	0		

Where N,  $C\alpha$ , C,O represent the vertices and N C $\alpha$ , C $\alpha$ C, CO represent the edges of subgraph (g) <sub>all</sub> and values in cells represent the existence or absence of connection.

 $L(g)_{all} =$ 

	Ν Cα	CαC	CO
$N \ C \alpha$	0	111.2	0
CαC	111.2	0	120.8
CO	0	120.8	0

Where N C $\alpha$ , C $\alpha$ C, CO represent the edges of aa backbone graph (all aa except glycine and proline) and values in cells represent the angle values in degrees between two edges. **Representation of the amino acid residue glycine backbone** 

as subgraph  $(g)_g = A(g)_g =$ 

5)	g —				
		Ν	Сα	С	0
		0	1.451	0	0
		1.451	0	1.516	0
	Ν Cα	0	1.516	0	1.231
	СО	0	0	1.231	0

Where N, C $\alpha$ , C, O represent the vertices of subgraph (g)<sub>g</sub> and values in cells represent the length of edges in angstrom I (g) <sub>g</sub> =



Where N, C $\alpha$ , C,O represent the vertices and N C $\alpha$ , C $\alpha$ C, CO represent the edges of subgraph (g)<sub>g</sub> and values in cells represent the existence or absence of connection.

 $L(g)_g =$ 

	Ν Cα	CaC	CO
Ν Cα	0	112.5	0
CαC	112.5	0	120.8
CO	0	120.8	0

Where N C $\alpha$ , C $\alpha$ C, CO represent the edges of glycine backbone graph and values in cells represent the angle values in degrees between two edges.

Representation of the amino acid residue proline backbone as subgraph  $(g)_{p}. \label{eq:generalized_prod}$ 

A	(g)	p=
	10/	P

Ν Cα	N	Сα	С	0
CO	0	1.466	0	0
	1.466	0	1.525	0
	0	1.525	0	1.231
	0	0	1.231	0

Where N,  $\overline{C\alpha}$ , C, O represent the vertices of subgraph (g)  $_{p}$  and values in cells represent the length of edges in angstrom I (g)  $_{p}$  =



Where N, C $\alpha$ , C,O represent the vertices and NC $\alpha$ , C $\alpha$ C, CO represent the edges of subgraph (g)  $_p$  and values in cells represent the existence or absence of connection.

 $L(g)_{p}=$ 

	Ν Cα	CaC	CO
Ν Cα	0	111.8	0
CaC	111.8	0	120.8
CO	0	120.8	0

Where N Ca, CaC, CO represent the edges of proline backbone graph and values in cells represent the angle values in degrees between two *edges*.

### The second step: -

Representation of each connection among (g)  $_{all}$ , (g)  $_{g}$  and (g)  $_{p}$  subgraphs using a connection matrix and angle matrix. C (g  $_{anv}$ , g  $_{all}$ ) =

Ν Cα	Ν	Сα	С	0
0	0	0	0	0
	0	0	0	0
	1.329	0	0	0
	0	0	0	0

#### Figure 1 Main-chain bond lengths and bond angles extracted from PROCHECK Operating Manual Appendix A Stereochemical parameters TABLE A.1

### a. Bond lengths

Bond	X-FLOR label1	ing	Value	signa
C-3	C-8E1	(except Pro)	1.329	0.014
	C-3	(Pro)	1.341	0.016
C-0	C-0		1.291	0.020
Calpha-C	CRIE-C	(except Gly)	1.525	0.021
	CE2G*-C	(G1y)	1.516	0.018
Calpha-Coeta	CELE-CESE	(Ale)	1.521	0.033
	CE1E-CE1E	(Ile,Thr,Val)	1,540	0.027
	CHIE-CHIE	(the rest)	1,530	0.020
N-Calpha	NE1-CE1E	(except Gl7, Pro)	1.458	0.019
	881-C82G*	(GL7)	1,451	0.016
	SI-CELE	(Pro)	1.466	0.015

#### b. Bond Angles

Rngle	X-PLOS labellin	g	Velue	signa
C-N-Calpha	C-MH1-CH1E	except Cly, Fro)	121.7	1.8
	C-MH1-CH2G*	(Gly)	120.6	1.7
	C-N-CRIE	(Pro)	122.6	5.0
Calpha-C-N	CHIE-C-RHI	    except Gly, Froj	116.2	2.0
	CH2G*-C-NH1	(G1y)	116.4	2.1
	CH1E-C-R	(Pro)	116.9	1.5
Calpha-C-O	CHIE-C-D	except Gly)	120.8	1.7
	CH2G*-C-0	(Gly)	120.8	2.1
Cbeta-Calpha-C	CH3E-CH1E-C	(Ala)	110.5	1.5
	CHIE-CHIE-C	(Ile,Thr,Val)	109.1	2.2
	CH2E-CH1E-C	(the rest)	110.1	1.9
	1	1	1 .	L.
H-Calpha-C	1011-CE1E-C	except Gly, Fra)	111.2	2.8
	M81-CE28*-C	(G1¥)	112.5	Z.3
	N-CHIE-C	(Fro)	111.8	2.5
N-Calpha-Cheta	MH1-CB1E-CH3E	(11a)	110.4	1.5
	MH1-CR1E-CH1E	(Ile,Thr,Val)	111.5	1.1.7
	N-CHIE-CHIE	(Pro)	103.0	1.1
	NH1-CE1E-CH2E	( (the rest)	110.5	1.7
0-0-5	0-C-NE1	except Frc)	123.0	1.6
	1 0-C-N	(Pro)	1 122.0	1.4

Where the vertical axe (g1) represents (g)  $_{all or}$  (g)  $_{g or}$  (g)  $_{p}$  subgraph vertices and the horizontal one (g2) represents (g)  $_{all}$ 

subgraph vertices. And values in cells represent the length of edges in angstrom.

Figure 2 All amino acid residue backbone (except glycine and proline)



Figure 3





Figure 4 proline amino acid residue backbone



 $CL (g_{any}, CN \& CN, g_{all}) =$ 

g any	CN		g all
ΝCα	0	121.7	NCα
CaC	116.2	0	CaC
OC	123.0	0	OC

Where N C $\alpha$ , C $\alpha$ C, OC represents the edges of g1, g2 directed subgraph, g2 represents (g) <sub>all</sub>, CN represent the directed connection between g1,g2 and values in cells represent the angle values in degrees between two edge.

	Ν	Сα	С	0
	0	0	0	0
	0	0	0	0
Ν Cα	1.341	0	0	0
СО	0	0	0	0

Where the vertical axe (g1) represents (g)  $_{all or}$  (g)  $_{g or}$  (g)  $_{p}$  subgraph vertices and the horizontal one (g2) represents (g)  $_{g}$  subgraph vertices. And values in cells represent the length of edges in angstrom.

 $CL (g_{any}, CN \& CN, g_g) =$ 

cc c1, 8 g/					
g	any	CN		g g	
N	JCα	0	120.6	ΝCα	
C	CαC	116.2	0	CaC	
C	C	123.0	0	OC	

Where N Ca, CaC, OC represent the edges of g1, g2 directed subgraph, g2 represents (g)  $_{\rm g}$ , CN represent the directed connection between g1, g2 and values in cells represent the angle values in degrees between two edges.

 $C(g_{any},g_p) =$ 

Ν Cα	N	Сα	С	0
00	0	0	0	0
	0	0	0	0
	1.341	0	0	0
	0	0	0	0

Where the vertical axe (g1) represents (g)  $_{all or}$  (g)  $_{g or}$  (g)  $_{p}$  subgraph vertices and the horizontal one (g2) represents (g)  $_{p}$  subgraph vertices. And values in cells represent the length of edges in angstrom

CL(g any,	CN	&	CN,	g	g)	=
				_	· • •	

, 0	g/		
g any	CN		gg
ΝCα	0	122.6	NCα
CaC	116.9	0	CaC
OC	122.0	0	OC

Where N Ca, CaC, O C represent the edges of g1,g2 directed subgraph , g2 represents  $(g)_p$ , CN represent the directed connection between g1,g2 and values in cells represent the angle values in degrees between two edges.

#### Figure 5

Any amino acid residue backbone and All amino acid residue backbone (except glycine and proline) dipeptide.











#### References

1. Gibbons, A. (1995). Algorithmic graph theory. Cambridge University Press, Cambridge, UK.

2. Giblin, P.J. (1977). Graphs, surfaces and homology, an introduction to algebraic topology. Chapman and Hall. Ltd, London 1977.

3. Gross, J.L. and Tucker, T.W. (1987). Topological graph theory. Jon Wiley & Sons, Inc, Canada 1987.

4. Wilson, R.J. (1972). Introduction to graph theory. Olivar& Boyed, Edinburgh.

5. Wilson, R.J. and Watkins, J.J. (1990). Graphs, an introductory approach, a first course in discrete mathematics. Jon Wiley & Sons Inc, Canada.