

Algorithms in Bioinformatics: Molecular Distance Geometry

Antonio Mucherino

`www.antoniomucherino.it`

IRISA, University of Rennes 1, Rennes, France

last update: October 5th 2016

Proteins

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

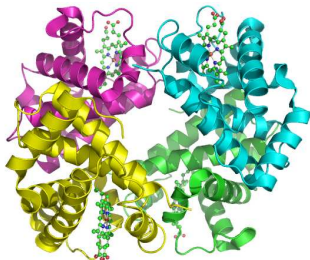
More research

Proteins are biochemical molecules consisting of **one or more polypeptides**, typically **folded** into a globular or fibrous form, which perform a certain **biological function**.

They are **chains** of smaller molecules called **amino acids**.

Their **three-dimensional conformations** can give clues about their biological function.

Google finds about 315,000,000 documents containing the word "protein".



Wikipedia: <http://en.wikipedia.org/wiki/Protein>
YouTube: <http://www.youtube.com/watch?v=Q7dxi4ob2O4>

The Protein Data Bank (PDB)

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP
the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP
the BP algorithm

Vertex orders

Making order

Consecutivity

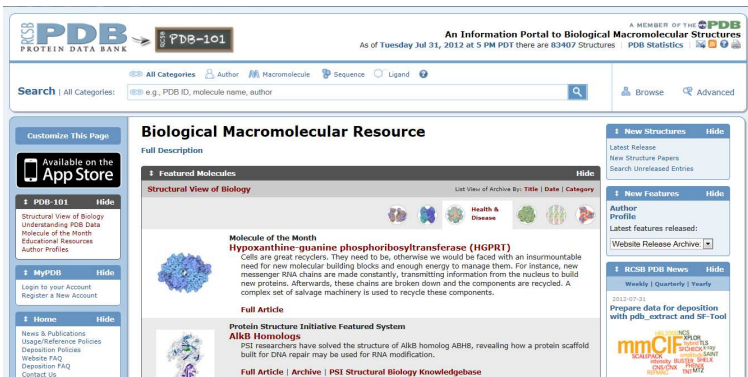
de Bruijn
Optimization

Ending

Challenge

More research

It's a database containing several protein three-dimensional conformations.



The screenshot shows the PDB website interface. At the top, there's a header with the PDB logo, a search bar, and navigation links. The main content area is titled "Biological Macromolecular Resource" and features a "Molecule of the Month" section for Hypoxanthine-guanine phosphoribosyltransferase (HGPRT). The sidebar on the left contains links to various resources like "PDB-101", "MyPDB", and "Home". The right sidebar shows "New Structures" and "New Features".

The database is experiencing a great expansion: this snapshot was taken a few years ago, meanwhile the total number of conformations in the database reached the 110,000 threshold!

<http://www.rcsb.org/pdb/>

Identifying protein conformations

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

How to identify the three-dimensional conformation of a protein?

- Experimental methods
 - X-ray crystallography
 - Nuclear Magnetic Resonance (NMR)
 - ...
- Computational methods
 - Homology modeling
 - Ab-initio approaches
 - ...

This is a non-exhaustive list.

Identifying protein conformations

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

How to identify the three-dimensional conformation of a protein?

- **Experimental methods**

- X-ray crystallography
- Nuclear Magnetic Resonance (NMR)
- ...

- **Computational methods**

- Homology modeling
- Ab-initio approaches
- ...

This is a non-exhaustive list.

X-ray crystallography

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

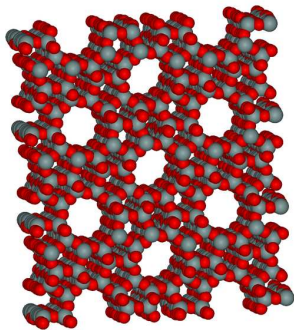
More research

X-ray crystallography is an experimental method for determining the arrangement of atoms within a crystal.

Crystals of proteins are generated in order to discover their conformation.

The crystal must have a certain size in order to be used.

The process of generating the crystal can be very difficult and expensive.



Wikipedia: http://en.wikipedia.org/wiki/X-ray_crystallography

YouTube: http://www.youtube.com/watch?v=j4HgLf_eJoc

The **Nuclear Magnetic Resonance** (NMR) studies the behavior of the magnetic moments of spin nuclei.

The protein sample is submitted to an external intense magnetic field, which induces the alignment of the magnetic moment of nuclei.

The analysis of this phenomenon allows to estimate the distance between pairs of nuclei (i.e., between pairs of atoms).

NMR do not directly provide information about the coordinates of the atoms.



Wikipedia: http://en.wikipedia.org/wiki/Nuclear_magnetic_resonance_spectroscopy

YouTube: <http://www.youtube.com/watch?v=IGk3NAziVWs>

Identifying protein conformations

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

How to identify the three-dimensional conformation of a protein?

- **Experimental methods**

- X-ray crystallography
- Nuclear Magnetic Resonance (NMR)
- ...

- **Computational methods**

- Homology modeling
- Ab-initio approaches
- ...

Identifying protein conformations

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

How to identify the three-dimensional conformation of a protein?

- **Experimental methods**

- X-ray crystallography
- **Nuclear Magnetic Resonance (NMR)**
- ...

- **Computational methods**

- Homology modeling
- Ab-initio approaches
- ...

We will study in details the problem of identifying protein conformations from the data obtained through NMR experiments.

the **Molecular Distance Geometry Problem**

MDGP

The Molecular Distance Geometry Problem

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Connectivity
de Bruijn
Optimization

Ending

Challenge

More research

Let $G = (V, E, d)$ be a **simple weighted undirected graph**, where

V **the set of vertices of G** – it is the set of atoms;

E **the set of edges of G** – it is the set of known distances;

$E' \subset E$ the subset of E where distances are exact;

d **the weights associated to the edges of G**
the numerical value of each weight corresponds to the
known distance; it can be an interval.

Definition

The **DGP** is the problem of finding an **embedding** $x : V \rightarrow \mathbb{R}^K$
such that:

$$\begin{aligned} \forall (u, v) \in E' \quad & \|x_u - x_v\| = d(u, v), \\ \forall (u, v) \in E \setminus E' \quad & \underline{d}(u, v) \leq \|x_u - x_v\| \leq \overline{d}(u, v). \end{aligned}$$

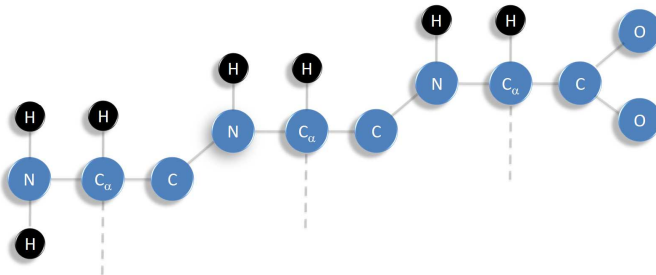
Equality constraints represent (hyper) spheres;

Inequality constraints represent (hyper) spherical shells.

The **MDGP** is NP-hard.

Where to find the necessary information about the distances?

- when working with molecules, a set of distances can be **derived** from their **chemical structure**:



- additional distances can be obtained by experimental techniques, such as **NMR**.

By definition, the MDGP is a **constraint satisfaction** problem.

However, it is generally reformulated as a **global optimization** problem, where the objective is to minimize a **penalty function** capable of measuring the **violation** of the constraints:

$$\frac{1}{|E|} \sum_{(u,v) \in E} \left[\frac{\max(\underline{d}(u,v) - \|x_u - x_v\|, 0)}{\underline{d}(u,v)} + \frac{\max(\|x_u - x_v\| - \bar{d}(u,v), 0)}{\bar{d}(u,v)} \right]$$

When all distances are correct, the **value** of the penalty function in the solution is **zero**.

The penalty function

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

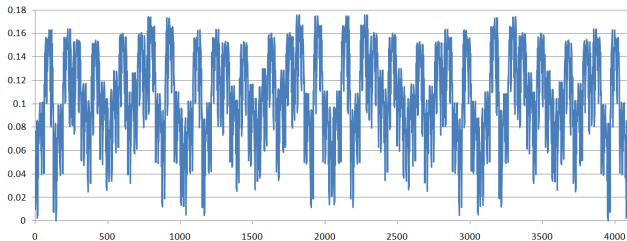
Optimization

Ending

Challenge

More research

The penalty function of the optimization problem is **strongly non-smooth**:



- this search space is, a priori, **continuous**,
- optimization methods risk to **get stuck** at local minima with objective value *very close* to the optimal one.

Function graphic from: C. Lavor, A. Mucherino, L. Liberti, N. Maculan, *On the Computation of Protein Backbones by using Artificial Backbones of Hydrogens*, **Journal of Global Optimization** 50(2), 329–344, 2011.

The Simulated Annealing (SA)

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

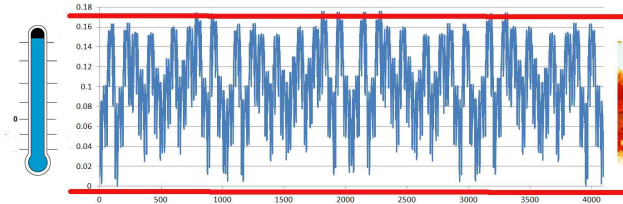
Optimization

Ending

Challenge

More research

The **Simulated Annealing** (**SA**) is based on the idea of simulating the physical annealing process for the solution of a global optimization problem.



SA is a **meta-heuristic** search:

- it can be applied to any optimization problem
- it can give no guarantees of optimality

SA and the MDGP

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

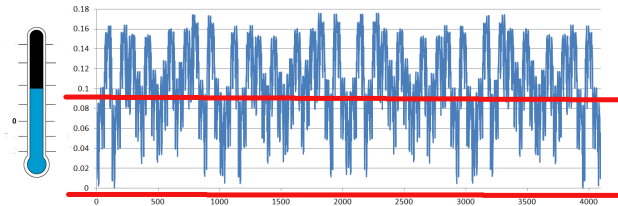
de Bruijn

Optimization

Ending

Challenge

More research



About SA and the MDGP:

- There are currently more than 110,000 molecular conformations on the [PDB](#)
- about the 10% of such conformations were obtained through [NMR](#) experiments
- in the detailed description of about 5% of such conformations, the name “*Simulated Annealing*” appears

SA and the MDGP

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

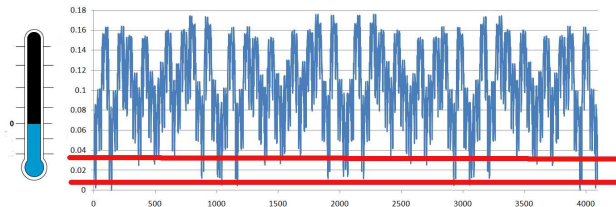
de Bruijn

Optimization

Ending

Challenge

More research



Disadvantages in using SA:

- there exist **other meta-heuristic searches** that are able to provide better quality results
- when the found solution has a penalty function **value larger than 0**, we cannot distinguish between
 - the given set of distances is not compatible
 - SA was not able to converge
- there is no hope to identify **all optimal solutions**.

the **Discretizable DGP**

DDGP

Intersecting spheres and spherical shells

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

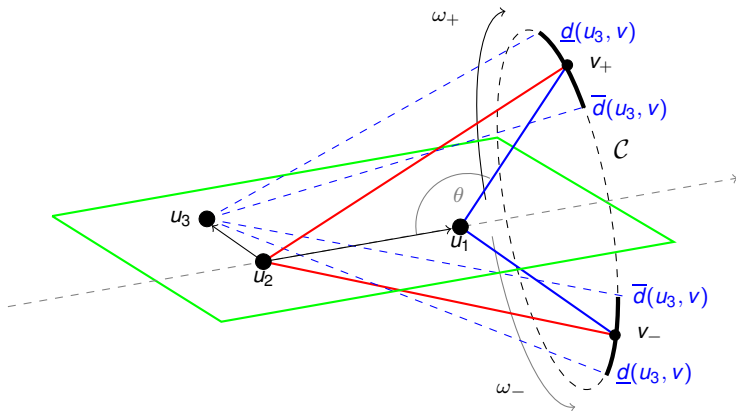
de Bruijn

Optimization

Ending

Challenge

More research



This drawing was made by Douglas Gonçalves, ancient postdoc student at University of Rennes 1.

The Discretizable DGP (DDGP in dimension $K = 3$)

Definition

A graph $G = (V, E, d)$ represents a **DDGP** instance if there exists a **vertex order** such that

A1 $G[\{1, 2, 3\}]$ is a clique consisting of exact distances;

A2 $\forall v \in V : v > 3, \quad \exists u_1, u_2, u_3 :$

$$\begin{cases} u_1 < v, u_2 < v, u_3 < v, \\ \{(u_1, v), (u_2, v)\} \subset E', (u_3, v) \in E, \\ A(u_1, u_2, u_3) > 0, \end{cases}$$

where A is the area of the triangle with vertices u_1, u_2, u_3 .

Notice that

- for all vertices $v > 3$, the atomic positions can be found by intersecting **2 spheres** with **1 spherical shell**
- the computation of A can be performed by using the distances (when available); this is a **probability 1** constraint
- this definition can be extended to any dimension $K > 0$

The new search space

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

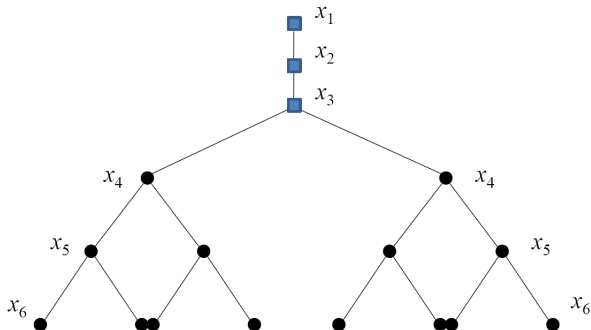
Optimization

Ending

Challenge

More research

When the discretization assumptions are satisfied, the domain of the penalty function can be reduced to a tree.



Notice that

- the tree is **binary** if only exact distances are available
- otherwise, D **sample positions** are selected from each arc for generating D new branches

Complexity of the DDGP

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance

Geometry

the MDGP

the Simulated
Annealing

Discrete

Distance

Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Definition

SUBSET-SUM. Given nonnegative integers a_1, \dots, a_n , is there a partition into two sets, encoded by $s \in \{-1, +1\}^n$, such that each subset has the same sum, i.e. $\sum_{i=1}^n s(i)a_i = 0$?

By reduction from the Subset-sum problem (which is known to be NP-hard), we can prove the following:

Theorem

The DDGP is NP-hard.

C. Lavor, L. Liberti, N. Maculan, A. Mucherino, *The Discretizable Molecular Distance Geometry Problem, Computational Optimization and Applications* 52, 115–146, 2012.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

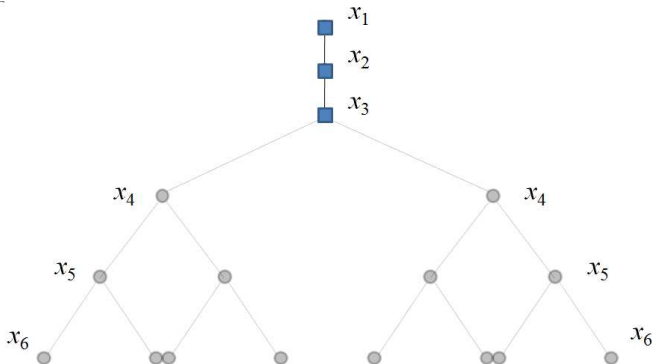
Optimization

Ending

Challenge

More research

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

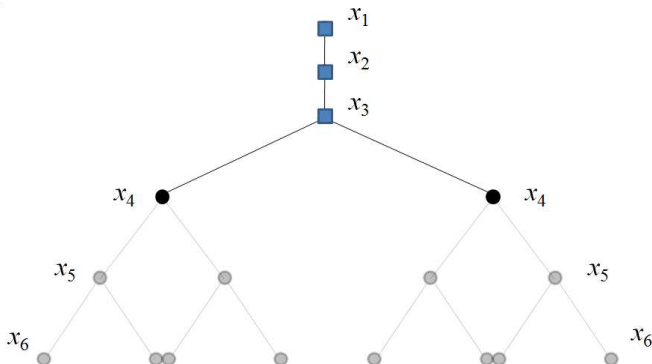
Optimization

Ending

Challenge

More research

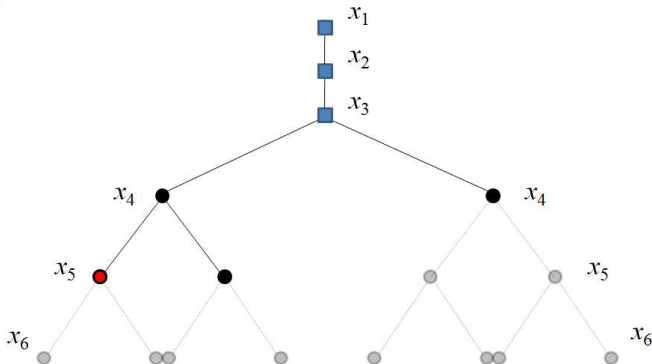
The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

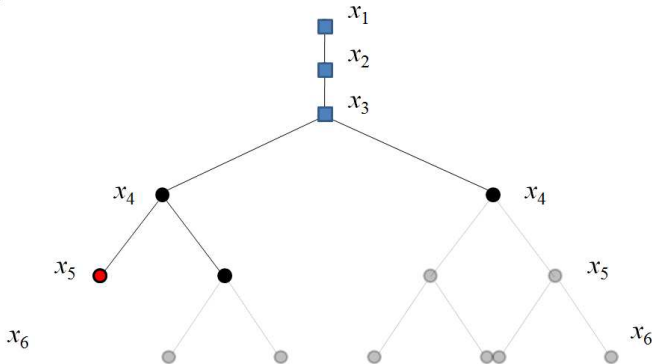
Optimization

Ending

Challenge

More research

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

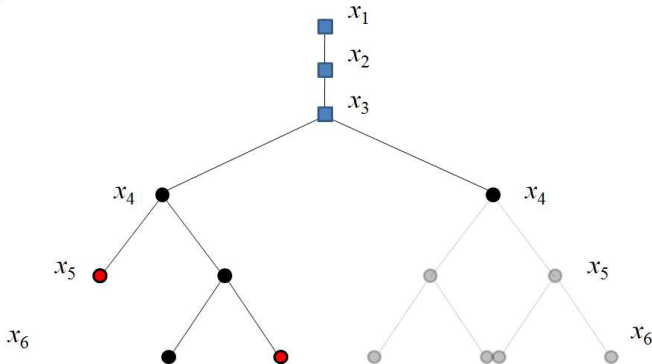
Optimization

Ending

Challenge

More research

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

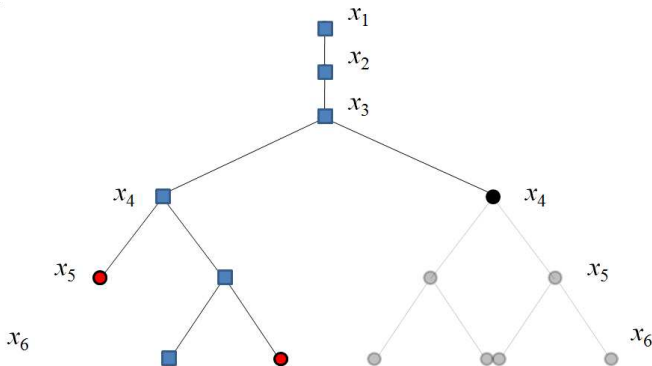
Optimization

Ending

Challenge

More research

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

The Branch & Prune (BP) algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

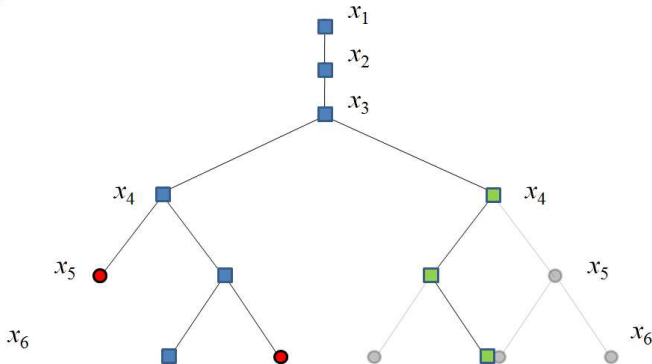
Optimization

Ending

Challenge

More research

The **Branch & Prune** (BP) algorithm is based on the idea of **branching** over all possible positions for each atom, and of **pruning** tree branches by using *additional distances* that are not used in the discretization process.



In this animation, it is supposed that all available distances are exact.

How to compute the **coordinates** of atoms during the execution of the BP algorithm?

- intersection of three spheres
 - solution of a quadratic system
 - solution of two linear systems
- method based on matrix multiplication
- method based on change of basis

For more details, see references in the last slides . . .

How to compute the **coordinates** of atoms during the execution of the BP algorithm?

- intersection of three spheres
 - solution of a quadratic system ... numerically unstable
 - solution of two linear systems ... numerically unstable
- method based on matrix multiplication ... stable
- method based on change of basis ... stable, fast

For more details, see references in the last slides ...

Pruning devices

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

The simplest and probably most efficient **pruning device** to be used with the BP algorithm is:

- **DDF** – **Direct Distance Feasibility**

if for some $v > 3$, $\exists u_j : u_j \notin \{u_1, u_2, u_3\}$, $u_j < v$, $d(u_j, v)$ is known, then verify whether:

$$\|x_v - x_{u_j}\| \in [\underline{d}(v, u_j), \bar{d}(v, u_j)].$$

Other pruning devices can be based on:

- information about torsion angles
- information about secondary structures
- potential energy for the molecule

A. Mucherino, C. Lavor, T. Malliavin, L. Liberti, M. Nilges, N. Maculan, *Influence of Pruning Devices on the Solution of Molecular Distance Geometry Problems*, **Lecture Notes in Computer Science 6630**, P.M. Pardalos, S. Rebennack (Eds.), Proceedings of the 10th International Symposium on Experimental Algorithms (SEA11), Crete, Greece, 206–217, 2011.

Many advantages but ...

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

The advantages of **BP**:

- the search space is built step by step;
- thanks to pruning devices, parts of the search space can be removed and never explored;
- the complete enumeration of the solution set may be performed.

Disadvantages:

- in order to apply BP, the discretization assumptions need to be satisfied!

Many advantages but ...

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

The advantages of **BP**:

- the search space is built step by step;
- thanks to pruning devices, parts of the search space can be removed and never explored;
- the complete enumeration of the solution set may be performed.

Disadvantages:

- in order to apply BP, the discretization assumptions need to be satisfied!

the **Ordering Problem**

Making order among the atoms

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance

Geometry

the MDGP

the Simulated
Annealing

Discrete

Distance

Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Let S be the set of all subsets $s \subseteq V$.

A **sequence of subsets** of V can be represented by a function $r : \mathbb{N} \longrightarrow S$ with length $|r| \in \mathbb{N}$ (for which $r_i = \emptyset$ for all $i > |r|$) such that, for each $v \in V$, there exist

- a non-empty subset $s \in S$ containing v
- an index $i \in \mathbb{N}$

such that $r_i = s$.

A *sequence of subsets* naturally implies a **partial order** on V .

A. Mucherino, *Optimal Discretization Orders for Distance Geometry: a Theoretical Standpoint*, **Lecture Notes in Computer Science 9374**, Proceedings of the 10th International Conference on Large-Scale Scientific Computations (LSSC15), Sozopol, Bulgaria, June 2015.

Total or partial orders?

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance

Geometry

the MDGP

the Simulated
Annealing

Discrete

Distance

Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

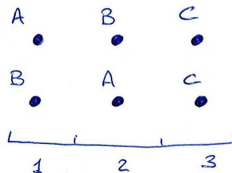
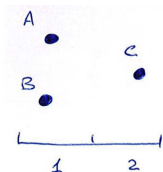
Challenge

More research

Definition

An **order** r is **total** if and only if, for each $i = 1, 2, \dots, |r|$, $|r_i| = 1$.

Notice that, if r is not total, different atoms may take the “same place” in the order. This kind of order is named **partial order**:



From every partial order, a set of total orders can be defined.

Definition

An **order** without **repetitions** is order where, for each pair r_i and r_j , with $i \neq j$, the intersection $r_i \cap r_j$ is empty.

Repetitions in atomic orders are necessary to satisfy some particular conditions.

Theorem

Every order without repetitions has finite length $|r|$.

Reference atoms

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Given an order r and an atom $v \in V$ such that $v \in r_i$,
how many **references** v has?

We define two sets of edges:

$$\Lambda_{\alpha}(r_i, v) = \{(u, v) \in E \mid \exists j < i : u \in r_j\}$$

$$\Lambda_{\beta}(r_i, v) = \{(v, u) \in E \mid \exists j \geq i : u \in r_j\}$$

We introduce four counters:

$$\alpha(r_i) = \min_{v \in r_i} |\Lambda_{\alpha}(r_i, v)|$$

$$\beta(r_i) = \max_{v \in r_i} |\Lambda_{\beta}(r_i, v)|$$

$$\alpha_{\text{ex}}(r_i) = \min_{v \in r_i} |\Lambda_{\alpha}(r_i, v) \cap E'|$$

$$\beta_{\text{ex}}(r_i) = \max_{v \in r_i} |\Lambda_{\beta}(r_i, v) \cap E'|$$

Discretization orders

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP
the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP
the BP algorithm

Vertex orders

Making order

Consecutivity
de Bruijn
Optimization

Ending

Challenge
More research

Let $G = (V, E, d)$ be a simple weighted undirected graph.
Let K be a positive integer.

Definition

A **discretization order** in dimension K is an order $r : \mathbb{N} \rightarrow S$ having finite length such that:

- (a) $r_1 = V_C$ where $G[V_C] = (V_C, E_C)$ is a clique with $V_C \subset V$, $|V_C| = K$ and $E_C \subset E'$;
- (b) $\forall i \in \{2, \dots, |r|\}, \alpha(r_i) \geq K$ and $\alpha_{\text{ex}}(r_i) \geq K - 1$.

Theorem

Necessary condition for G to admit a **discretization order** in dimension K is that, for any order r on V without repetitions,

$$\forall i \in \{1, 2, \dots, |r|\}, \quad \alpha(r_i) + \beta(r_i) \geq K, \quad \alpha_{\text{ex}}(r_i) + \beta_{\text{ex}}(r_i) \geq K - 1.$$

The ordering problem

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Definition

Given a simple weighted undirected graph $G = (V, E, d)$ and a positive integer K , establish whether there exists an order r in dimension K such that:

- (a) $r_1 = V_C$ where $G[V_C] = (V_C, E_C)$ is a clique with $V_C \subset V$, $|V_C| = K$ and $E_C \subset E'$;
- (b) $\forall i \in \{2, \dots, |r|\}$, $\alpha(r_i) \geq K$ and $\alpha_{\text{ex}}(r_i) \geq K - 1$.
- (c) a set of objectives f_ℓ ($\ell = 1, \dots, M$) is optimized for every r_i (with priority order)

The consecutivity assumption

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Definition

An order satisfies the **consecutivity assumption** if, for every subset r_i such that $i > 1$,

- $|r_i| = 1$
- if $P = \{(u, v) \in E \mid \exists j \in \{i - K, \dots, i\} : u \in r_j\}$,
then $|P| \geq K$

Why?

- it allows us to represent the order as a sequence of *overlapping cliques*
- in order to satisfy this additional assumption, atoms generally need to be *repeated* in the order
- the feasibility of each clique can be *a priori* verified

The consecutivity assumption

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Definition

An order satisfies the **consecutivity assumption** if, for every subset r_i such that $i > 1$,

- $|r_i| = 1$
- if $P = \{(u, v) \in E \mid \exists j \in \{i - K, \dots, i\} : u \in r_j\}$,
then $|P| \geq K$

However:

- Finding an order *with* the consecutivity assumption is NP-hard
- Finding an order *without* consecutivity assumption has polynomial complexity when K is fixed

A handcrafted order

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

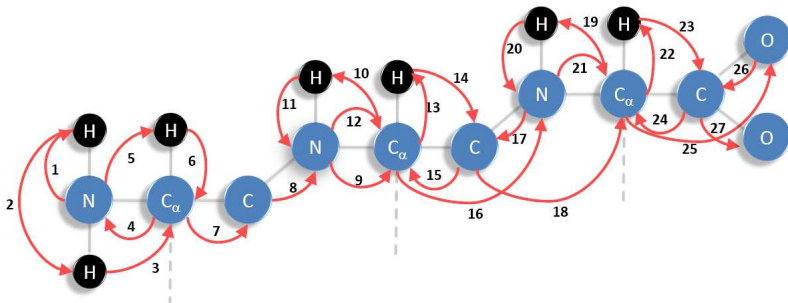
Optimization

Ending

Challenge

More research

Only information about bond length, bond angles and torsion angles are here considered (NMR data not included).



Remarks:

- no objectives are here optimized
- the consecutivity assumption is satisfied

Orders satisfying the consecutivity assumption can be obtained by exploring *pseudo de Bruijn graphs* B .

Let $B = (V_B, E_B)$ be a directed graph, defined as follows:

- 1 $c \in V_B$ is a $(K + 1)$ -clique of G
- 2 $(b, c) \in E_B$ if the cliques b and c admit a K -overlap

K -overlap: the K -suffix of b coincides with the K -prefix of c , in a possible *internal ordering* for the atoms in b and the atoms in c .

Discretization order with consecutivity assumption:

A *path* on B such that:

- the internal order of cliques c is constant on the path
- the set of vertices deduced from the set of cliques covers V

Finding this path has exponential complexity.

Expected: finding an order with consecutivity assumption is
NP-hard!

This table contains all 4-cliques in a 3-amino acid backbone:

name	atoms	edge $\{r_{i-3}, r_i\}$	name	atoms	edge $\{r_{i-3}, r_i\}$
c_1	$N^1 C_\alpha^1 H_\alpha^1 C^1$	exact	c_7	$N^2 C_\alpha^2 H_\alpha^2 C^2$	exact
c_2	$H_\alpha^1 C_\alpha^1 C^1 N^2$	interval	c_8	$H_\alpha^2 C_\alpha^2 C^2 N^3$	interval
c_3	$C_\alpha^1 C^1 N^2 H^2$	exact	c_9	$C_\alpha^2 C^2 N^3 H^3$	exact
c_4	$C_\alpha^1 C^1 N^2 C_\alpha^2$	exact	c_{10}	$C_\alpha^2 C^2 N^3 C_\alpha^3$	exact
c_5	$C^1 N^2 H^2 C_\alpha^2$	exact	c_{11}	$C^2 N^3 H^3 C_\alpha^3$	exact
c_6	$H^2 N^2 C_\alpha^2 H_\alpha^2$	interval	c_{12}	$H^3 N^3 C_\alpha^3 H_\alpha^3$	interval
			c_{13}	$N^3 C_\alpha^3 H_\alpha^3 C^3$	exact

Auxiliary cliques can be added to B by duplicating one atom in the 3-cliques of G : allowed internal orders in auxiliary cliques are the ones where the repeated atom takes the first and the last position.

This introduces atomic repetitions in the orders, which are often necessary for finding orders satisfying the consecutivity assumption.

de Bruijn order

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

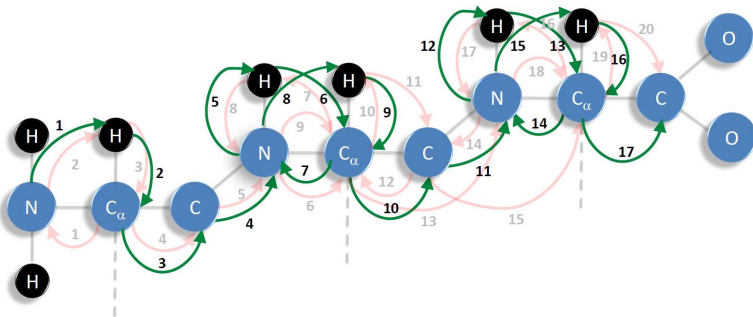
Optimization

Ending

Challenge

More research

A comparison between the handcrafted order and one possible *de Bruijn* order.



A. Mucherino, *A Pseudo de Bruijn Graph Representation for Discretization Orders for Distance Geometry*, **Lecture Notes in Computer Science 9043**, Lecture Notes in Bioinformatics series, F. Ortuño, I. Rojas (Eds.), Proceedings of the 3rd International Work-Conference on Bioinformatics and Biomedical Engineering (IWBBIO15), Part I, Granada, Spain, 514–523, 2015.

A greedy algorithm

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance

Geometry

the MDGP

the Simulated
Annealing

Discrete

Distance

Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Greedy algorithm *in: G out: r*

// initial clique

choose a K -clique $G_C = (V_C, E_C)$ in V with edges in E'

set $r_1 = V_C$

let $A = V \setminus V_C$

set $i = 2$

// constructing the rest of the order

while ($A \neq \emptyset$) **do**

let $A^0 = \{v \in A : \alpha(v) \geq K, \alpha_{ex}(v) \geq K - 1\}$

if ($A^0 = \emptyset$) **then**

break: no possible orders; **choose** another initial clique

else

for each objective f_ℓ ($\ell = 1, \dots, M$) **do**

$A^\ell = \{v \in A^{\ell-1} : f_\ell(v) \text{ is optimized}\}$

end for

set $r_i = A^M$

let $A = A \setminus \{r_i\}$

let $i = i + 1$

end if

end while

A. Mucherino, *Optimal Discretization Orders for Distance Geometry: a Theoretical Standpoint*, **Lecture Notes in Computer Science 9374**, Proceedings of the 10th International Conference on Large-Scale Scientific Computations (LSSC15), Sozopol, Bulgaria, June 2015.

The objectives

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance

Geometry

the MDGP

the Simulated
Annealing

Discrete

Distance

Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

Assumption **(c)** of the ordering problem allows for **constructing** discretization orders having **some additional properties**.

We may try to generate orders that **make** the BP algorithm **more efficient**.

One objective can correspond to the counter α :

$$f_1(v) = \alpha(v)$$

- maximizing f_1 means selecting the vertices having the maximal number of reference atoms
- since A^0 contains vertices having at least K references (necessary for the discretization), f_1 enforces the use of vertices where pruning distances are also available
- early pruning on the discrete search domain allows for a more efficient execution of BP

Another order without repetitions

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP
the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP
the BP algorithm

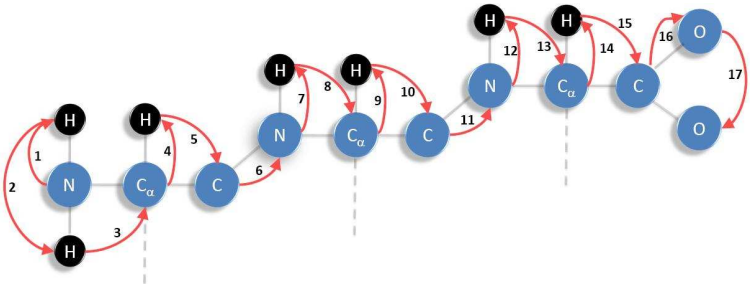
Vertex orders

Making order
Consecutivity
de Bruijn
Optimization

Ending

Challenge
More research

This order was automatically obtained by the greedy algorithm. It is a total order.



A. Mucherino, *On the Identification of Discretization Orders for Distance Geometry with Intervals*, **Lecture Notes in Computer Science 8085**, F. Nielsen and F. Barbaresco (Eds.), Proceedings of GSI13, Paris, France, 231–238, 2013.

More objectives?

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

What about including some other objectives f_ℓ ???

- anticipate the use of exact distances
- maximize the use of distances between hydrogens
- minimize the rank-difference of distances used for pruning
- maximize the number of cliques used for computing atomic coordinates
- ...

When **discretizing interval distances**, the predefined number of samples D taken from each arc plays a **critical role**:

- D too small \longrightarrow a few chances to catch the “true distance”
- D too large \longrightarrow high increase of computational cost

Possible solutions for overcoming this issue:

- 1 choose the best D value layer by layer
- 2 try to discover *in advance* whether sample points will lead to infeasibilities in deeper layers
- 3 avoid the discretization of the intervals: make the search locally continuous

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

L. Liberti, C. Lavor, N. Maculan, A. Mucherino,
Euclidean Distance Geometry and Applications,
SIAM Review 56(1), 3–69, 2014.

More references about discretization orders

AlgBioInfo

A. Mucherino

Introduction

Proteins

Methods for protein
determination

Distance
Geometry

the MDGP

the Simulated
Annealing

Discrete
Distance
Geometry

The DDGP

the BP algorithm

Vertex orders

Making order

Consecutivity

de Bruijn

Optimization

Ending

Challenge

More research

- D.S. Gonçalves, A. Mucherino, *Optimal Partial Discretization Orders for Discretizable Distance Geometry*, International Transactions in Operational Research **23**(5), 947–967, 2016.
- C. Lavor, J. Lee, A. Lee-St.John, L. Liberti, A. Mucherino, M. Sviridenko, *Discretization Orders for Distance Geometry Problems*, Optimization Letters **6**(4), 783–796, 2012.

Symmetry properties of discretizable instances

- A. Mucherino, C. Lavor, L. Liberti, *Exploiting Symmetry Properties of the Discretizable Molecular Distance Geometry Problem*, Journal of Bioinformatics and Computational Biology **10**(3), 1242009(1–15), 2012.
- A. Mucherino, C. Lavor, L. Liberti, *A Symmetry-Driven BP Algorithm for the Discretizable Molecular Distance Geometry Problem*, IEEE Conference Proceedings, Computational Structural Bioinformatics Workshop (CSBW11), International Conference on Bioinformatics & Biomedicine (BIBM11), Atlanta, GA, USA, 390–395, 2011.
- ...

Efficient generation of atomic coordinates

- D.S. Gonçalves, A. Mucherino, *Discretization Orders and Efficient Computation of Cartesian Coordinates for Distance Geometry*, Optimization Letters **8**(7), 2111–2125, 2014.
- A. Mucherino, C. Lavor, L. Liberti, *The Discretizable Distance Geometry Problem*, Optimization Letters **6**(8), 1671–1686, 2012.
- ...

Parallel and distributed versions of the BP algorithm

- W. Gramacho, A. Mucherino, C. Lavor, N. Maculan, *A Parallel BP Algorithm for the Discretizable Distance Geometry Problem*, IEEE Conference Proceedings, Workshop on Parallel Computing and Optimization (PCO12), 26th IEEE International Parallel & Distributed Processing Symposium (IPDPS12), Shanghai, China, 1756–1762, 2012.
- A. Mucherino, C. Lavor, L. Liberti, E-G. Talbi, *A Parallel Version of the Branch & Prune Algorithm for the Molecular Distance Geometry Problem*, IEEE Conference Proceedings, ACS/IEEE International Conference on Computer Systems and Applications (AICCSA10), Hammamet, Tunisia, 1–6, 2010.
- ...

Management of *real* NMR instances

- A. Cassioli, B. Bardiaux, G. Bouvier, A. Mucherino, R. Alves, L. Liberti, M. Nilges, C. Lavor and T.E Malliavin, *An Algorithm to Enumerate all Possible Protein Conformations Verifying a Set of Distance Restraints*, BMC Bioinformatics **16**:23, 15 pages, 2015.
- A. Mucherino, C. Lavor, T. Malliavin, L. Liberti, M. Nilges, N. Maculan, *Influence of Pruning Devices on the Solution of Molecular Distance Geometry Problems*, Lecture Notes in Computer Science **6630**, P.M. Pardalos and S. Rebennack (Eds.), Proceedings of the 10th International Symposium on Experimental Algorithms (SEA11), Crete, Greece, 206–217, 2011.
- ...

The End