

MASTER I INFORMATIQUE
TRAVAIL ENCADRÉ DE RECHERCHE
JUIN 2018

Rapport QuChemPedia
Sous titre

Auteur

Jules LEGUY

Encadrants

Benoit DA MOTA

Thomas CAUCHY

Table des matières

1	Introduction	2
2	Description des problèmes	3
3	Représentations géométriques moléculaires	4
3.1	Matrice des coordonnées	4
3.2	Matrice réduite des distances inter-atomiques	4
3.2.1	Motivation	4
3.2.2	Formalisation	4
3.2.3	Reconstruction des molécules	5
3.3	Matrice réduite des distances à des points fixes	10
3.4	10
4	Modèles prédictifs	11
5	Conclusion	12
	Appendices	12

Chapitre 1

Introduction

Chapitre 2

Description des problèmes

Chapitre 3

Représentations géométriques moléculaires

3.1 Matrice des coordonnées

3.2 Matrice réduite des distances inter-atomiques

3.2.1 Motivation

Cette représentation est issue du travail qui a été fait précédemment sur ce projet, et consiste à représenter une molécule par ses distances inter-atomiques. L'intérêt de cette représentation est que les réseaux de neurones qui l'utilisent travaillent dans des repères relatifs. Lorsqu'ils effectuent des prédictions, ils n'ont pas besoin de *comprendre* les notions mathématiques de géométrie permettant de déduire la position d'un point dans un repère à partir de ses distances à d'autres points, contrairement à la représentation décrite en (REF REPR ABS). Cette représentation est donc très commode pour les modèles prédictifs dont l'objectif est de corriger les distances entre deux atomes, puisqu'elle est basée sur les distances entre les paires d'atomes.

De plus, l'utilisation d'une représentation basée sur les distances relatives permet d'offrir une représentation unique pour les molécules ayant des ensembles d'atomes pouvant effectuer des rotations, contrairement aux représentations basées sur les coordonnées (REF REPR COORDS) ou sur des distances à des points fixes (REF REPR DIST ABS).

Lorsque les modèles utilisent cette représentation en sortie, ou plus précisément que l'on déduit la matrice réduite des distances inter-atomiques de la sortie du modèle (voir REF SORTIE DELTA_DIST+H), nous devons toutefois trouver une méthode (voir REF PRINC RECONSTRUCT) pour reconstruire les molécules sous la forme d'une matrice de coordonnées (REF REPR COORDS).

3.2.2 Formalisation

Pour ne pas surcharger les modèles d'information, nous ne travaillons pas sur la matrice de distances inter-atomiques complète, mais sur un sous-ensemble de cardinalité minimale de cette matrice telle que nous pouvons reconstruire sans ambiguïté un ensemble de coordonnées représentant les positions des atomes de la molécule. La matrice des distances étant symétrique et la diagonale étant nulle, toute l'information est contenue dans chaque demi-matrice triangulaire privée de la diagonale.

De plus, nous n'avons besoin que des distances à quatre points pour retrouver la position de chaque atome (voir REF RECONSTRUCT), nous nous contentons donc de garder les quatre premières distances de chaque ligne de la matrice triangulaire supérieure privée de la diagonale.

	a_0	a_1	a_2	a_3	a_4	\dots	a_{n-4}	a_{n-3}	a_{n-2}	a_{n-1}	a_n
a_0	$d_{0,0}$	$d_{0,1}$	$d_{0,2}$	$d_{0,3}$	$d_{0,4}$	\dots	$d_{0,n-4}$	$d_{0,n-3}$	$d_{0,n-2}$	$d_{0,n-1}$	$d_{0,n}$
a_1	$d_{1,0}$	$d_{1,1}$	$d_{1,2}$	$d_{1,3}$	$d_{1,4}$	\dots	$d_{1,n-4}$	$d_{1,n-3}$	$d_{1,n-2}$	$d_{1,n-1}$	$d_{1,n}$
a_2	$d_{2,0}$	$d_{2,1}$	$d_{2,2}$	$d_{2,3}$	$d_{2,4}$	\dots	$d_{2,n-4}$	$d_{2,n-3}$	$d_{2,n-2}$	$d_{2,n-1}$	$d_{2,n}$
a_3	$d_{3,0}$	$d_{3,1}$	$d_{3,2}$	$d_{3,3}$	$d_{3,4}$	\dots	$d_{3,n-4}$	$d_{3,n-3}$	$d_{3,n-2}$	$d_{3,n-1}$	$d_{3,n}$
a_4	$d_{4,0}$	$d_{4,1}$	$d_{4,2}$	$d_{4,3}$	$d_{4,4}$	\dots	$d_{4,n-4}$	$d_{4,n-3}$	$d_{4,n-2}$	$d_{4,n-1}$	$d_{4,n}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\vdots	\vdots
a_{n-4}	$d_{n-4,0}$	$d_{n-4,1}$	$d_{n-4,2}$	$d_{n-4,3}$	$d_{n-4,4}$	\dots	$d_{n-4,n-4}$	$d_{n-4,n-3}$	$d_{n-4,n-2}$	$d_{n-4,n-1}$	$d_{n-3,n}$
a_{n-3}	$d_{n-3,0}$	$d_{n-3,1}$	$d_{n-3,2}$	$d_{n-3,3}$	$d_{n-3,4}$	\dots	$d_{n-3,n-4}$	$d_{n-3,n-3}$	$d_{n-3,n-2}$	$d_{n-3,n-1}$	$d_{n-3,n}$
a_{n-2}	$d_{n-2,0}$	$d_{n-2,1}$	$d_{n-2,2}$	$d_{n-2,3}$	$d_{n-2,4}$	\dots	$d_{n-2,n-4}$	$d_{n-2,n-3}$	$d_{n-2,n-2}$	$d_{n-2,n-1}$	$d_{n-2,n}$
a_{n-1}	$d_{n-1,0}$	$d_{n-1,1}$	$d_{n-1,2}$	$d_{n-1,3}$	$d_{n-1,4}$	\dots	$d_{n-1,n-4}$	$d_{n-1,n-3}$	$d_{n-1,n-2}$	$d_{n-1,n-1}$	$d_{n-1,n}$
a_n	$d_{n,0}$	$d_{n,1}$	$d_{n,2}$	$d_{n,3}$	$d_{n,4}$	\dots	$d_{n,n-4}$	$d_{n,n-3}$	$d_{n,n-2}$	$d_{n,n-1}$	$d_{n,n}$

FIGURE 3.1 – Matrice complète des distances inter-atomiques d’une molécule

	a_0	a_1	a_2	a_3	a_4	\dots	a_{n-4}	a_{n-3}	a_{n-2}	a_{n-1}	a_n
a_0	$d_{0,0}$	$d_{0,1}$	$d_{0,2}$	$d_{0,3}$	$d_{0,4}$	\dots	$d_{0,n-4}$	$d_{0,n-3}$	$d_{0,n-2}$	$d_{0,n-1}$	$d_{0,n}$
a_1	$d_{1,0}$	$d_{1,1}$	$d_{1,2}$	$d_{1,3}$	$d_{1,4}$	\dots	$d_{1,n-4}$	$d_{1,n-3}$	$d_{1,n-2}$	$d_{1,n-1}$	$d_{1,n}$
a_2	$d_{2,0}$	$d_{2,1}$	$d_{2,2}$	$d_{2,3}$	$d_{2,4}$	\dots	$d_{2,n-4}$	$d_{2,n-3}$	$d_{2,n-2}$	$d_{2,n-1}$	$d_{2,n}$
a_3	$d_{3,0}$	$d_{3,1}$	$d_{3,2}$	$d_{3,3}$	$d_{3,4}$	\dots	$d_{3,n-4}$	$d_{3,n-3}$	$d_{3,n-2}$	$d_{3,n-1}$	$d_{3,n}$
a_4	$d_{4,0}$	$d_{4,1}$	$d_{4,2}$	$d_{4,3}$	$d_{4,4}$	\dots	$d_{4,n-4}$	$d_{4,n-3}$	$d_{4,n-2}$	$d_{4,n-1}$	$d_{4,n}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots	\vdots	\vdots
a_{n-4}	$d_{n-4,0}$	$d_{n-4,1}$	$d_{n-4,2}$	$d_{n-4,3}$	$d_{n-4,4}$	\dots	$d_{n-4,n-4}$	$d_{n-4,n-3}$	$d_{n-4,n-2}$	$d_{n-4,n-1}$	$d_{n-4,n}$
a_{n-3}	$d_{n-3,0}$	$d_{n-3,1}$	$d_{n-3,2}$	$d_{n-3,3}$	$d_{n-3,4}$	\dots	$d_{n-3,n-4}$	$d_{n-3,n-3}$	$d_{n-3,n-2}$	$d_{n-3,n-1}$	$d_{n-3,n}$
a_{n-2}	$d_{n-2,0}$	$d_{n-2,1}$	$d_{n-2,2}$	$d_{n-2,3}$	$d_{n-2,4}$	\dots	$d_{n-2,n-4}$	$d_{n-2,n-3}$	$d_{n-2,n-2}$	$d_{n-2,n-1}$	$d_{n-2,n}$
a_{n-1}	$d_{n-1,0}$	$d_{n-1,1}$	$d_{n-1,2}$	$d_{n-1,3}$	$d_{n-1,4}$	\dots	$d_{n-1,n-4}$	$d_{n-1,n-3}$	$d_{n-1,n-2}$	$d_{n-1,n-1}$	$d_{n-1,n}$
a_n	$d_{n,0}$	$d_{n,1}$	$d_{n,2}$	$d_{n,3}$	$d_{n,4}$	\dots	$d_{n,n-4}$	$d_{n,n-3}$	$d_{n,n-2}$	$d_{n,n-1}$	$d_{n,n}$

FIGURE 3.2 – Matrice réduite des distances inter-atomiques d’une molécule (en gras)

$d_{0,1}$	$d_{0,2}$	$d_{0,3}$	$d_{0,4}$
$d_{1,2}$	$d_{1,3}$	$d_{1,4}$	$d_{1,5}$
$d_{2,3}$	$d_{2,4}$	$d_{2,5}$	$d_{2,6}$
$d_{3,4}$	$d_{3,5}$	$d_{3,6}$	$d_{3,7}$
\vdots	\vdots	\vdots	\vdots
$d_{n-4,n-3}$	$d_{n-4,n-2}$	$d_{n-4,n-1}$	$d_{n-4,n}$
$d_{n-3,n-2}$	$d_{n-3,n-1}$	$d_{n-3,n}$	0
$d_{n-2,n-1}$	$d_{n-2,n}$	0	0
$d_{n-1,n}$	0	0	0

FIGURE 3.3 – Matrice réduite des distances inter-atomiques d’une molécule

3.2.3 Reconstruction des molécules

Lorsqu’un modèle a pour sortie une matrice réduite des distances inter-atomiques lorsqu’il effectue des prédictions, il faut définir une méthode pour reconstruire une matrice des coordonnées (ref REPR MAT COORDS) de façon automatique à partir de cette sortie, la seule contrainte étant que la distance relative entre chaque paire d’atomes soit respectée. Il ne s’agit pour autant pas d’une tâche triviale, elle s’est en effet avérée impossible en pratique pour les grosses molécules à cause de la propagation des erreurs qu’elle induit (voir REF PROPAG ERREURS).

3.2.3.1 Formalisation de la méthode de reconstruction

Nécessité et limite de l'introduction d'un atome fictif Notre méthode de reconstruction des atomes doit permettre de respecter la chiralité¹ des molécules. Or, en déduisant uniquement la position d'un atome de ses distances aux quatre atomes précédents, il existe des cas pour lesquels il existe plusieurs solutions pour la position de l'atome (deux si les quatre atomes précédents sur un même plan ou une infinité si les quatre atomes précédents appartiennent à une droite). Pour pallier ce problème, la méthode retenue précédemment a été d'introduire un nouveau point (que l'on nomme atome fictif) et que l'on place arbitrairement dans la molécule, à une position telle qu'il n'appartient pas au plan formé par les trois premiers atomes, ou à la droite formée par les trois premiers atomes s'ils sont alignés. De cette façon, les atomes suivants seront placés sans ambiguïté.

Cependant, on peut imaginer des cas pour lesquels la technique de l'introduction d'un atome fictif ne permet pas de lever l'ambiguïté, notamment pour les molécules possédant une chaîne de carbones liés par des doubles liaisons (et formant donc une droite). La méthode ne permettra pas dans ce cas de déterminer les positions des atomes en bout de chaîne tel que leurs distances relatives soient respectées, cette information étant perdue lors de la création de la matrice réduite des distances inter-atomiques.

Cette représentation n'est donc pas viable en pratique. Cela fait partie des raisons (voir également REF PB SQRT) pour lesquelles nous sommes passés à la représentation par matrice réduite des distances à des points fixes (REF MATR RED FIXES).

Placement de l'atome fictif Puisque l'on définit la position de chaque atome en fonction de ses distances aux quatre atomes précédents, on doit d'abord placer les quatre premiers atomes de façon partiellement arbitraire. Le premier atome de la molécule dans notre représentation étant l'atome fictif a_0 , nous commençons par le placer à la position qui lui a été attribuée.

Placement de l'atome a_1 Une fois l'atome a_0 placé, il existe une infinité de solutions pour la position de l'atome a_1 . On peut en effet le placer à tout point appartenant à la surface de la sphère de centre a_0 et de rayon $d_{0,1}$.

Placement de l'atome a_2 L'atome a_2 appartient au cercle solution de l'intersection entre les sphères de centres a_0 et a_1 et de rayons $d_{0,2}$ et $d_{1,2}$. On choisit donc arbitrairement une position appartenant à ce cercle.

Placement de l'atome a_3 Dans le cas général, il existe deux solutions pour le placement de l'atome a_3 , l'intersection non nulle de trois sphères étant deux points si tous les points ne sont pas sur un même plan ou une même droite. On choisit arbitrairement un point parmi ces deux solutions, car il n'y a pas à ce stade d'ambiguïté de chiralité de la molécule, une molécule composée de trois atomes ne possédant pas de chiralité (l'atome fictif a_0 ne fait pas partie de la molécule).

Placement de l'atome a_n Pour placer l'atome a_n (n étant strictement inférieur à la taille de la molécule), nous généralisons la méthode de placement de l'atome a_3 . Plutôt que de travailler sur l'intersection de quatre sphères, nous travaillons toujours sur l'intersection de trois sphères et nous utilisons la dernière distance pour discriminer les deux solutions obtenues. Cela facilite grandement la résolution des équations mathématiques associées et permet d'obtenir des solutions sensiblement équivalentes.

Formellement, nous calculons les positions des deux points solutions de l'intersection des trois sphères de centres a_{n-4} , a_{n-3} , et a_{n-2} et de rayons $d_{n-4,n}$, $d_{n-3,n}$ et $d_{n-2,n}$, et nous discriminons les deux solutions selon la distance $d_{n-1,n}$.

3.2.3.2 Reconstruction automatique des positions en utilisant un solveur

Nous développons ici une méthode permettant de déterminer les coordonnées d'un atome quelconque en utilisant un solveur d'équations non linéaires². Nous utilisons la bibliothèque Sympy³.

1. Un composé chimique est dit chiral s'il n'est pas superposable à son image dans un miroir. ([https://fr.wikipedia.org/wiki/Chiralité_\(chimie\)](https://fr.wikipedia.org/wiki/Chiralité_(chimie)))

2. https://en.wikipedia.org/wiki/Nonlinear_system

3. <http://www.sympy.org/fr/>

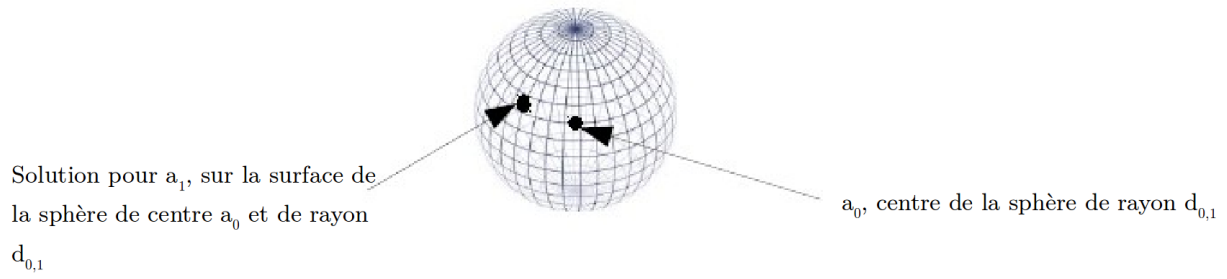


FIGURE 3.4 – Placement de l'atome a_1 (image extraite du rapport de N.Roux)

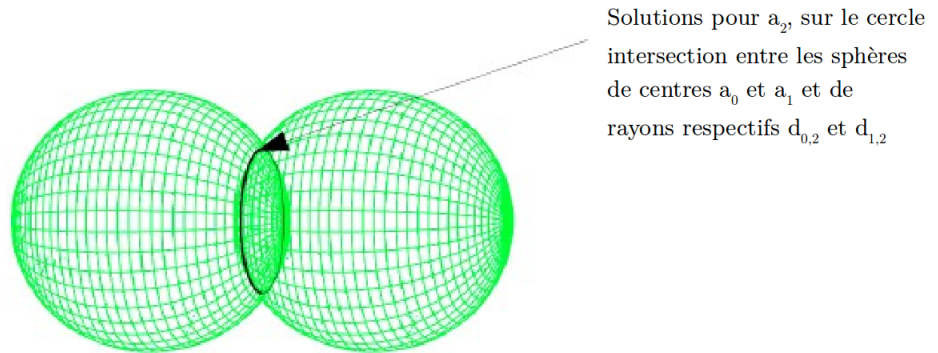


FIGURE 3.5 – Placement de l'atome a_2 (image extraite du rapport de N.Roux)

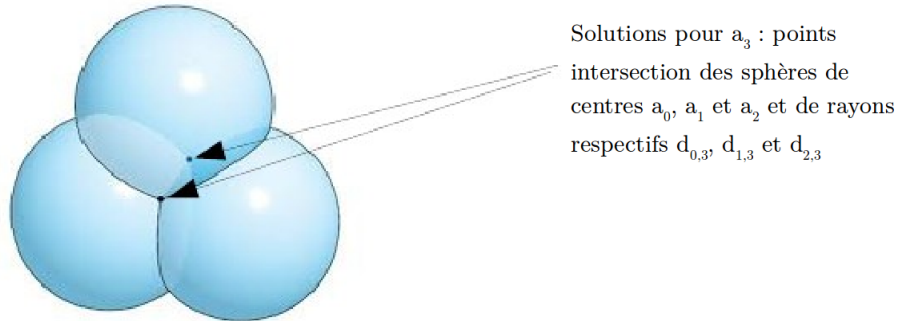


FIGURE 3.6 – Placement de l'atome a_3 (image extraite du rapport de N.Roux)

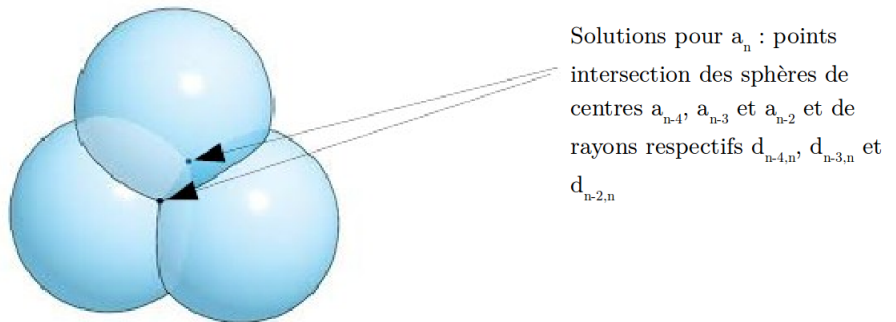


FIGURE 3.7 – Placement de l'atome a_n (image extraite du rapport de N.Roux)

Tout d'abord, l'atome fictif a_0 doit être placé à la position qui lui a été attribuée (REF PLACEMENT AT FICTIF). Nous plaçons ensuite arbitrairement les trois atomes suivants, de sorte que leurs distances relatives soient respectées. Pour simplifier le problème, nous effectuons une translation temporaire telle que a_0' est à l'origine du repère. Nous plaçons alors a_1' sur l'axe x , à une distance $d_{0,1}$ de l'origine, et a_2' sur le plan tel que $z = 0$, à une position telle que les distances $d_{0,2}$ et $d_{1,2}$ sont respectées. Pour finir, nous plaçons a_3' à l'une des deux solutions de l'intersection des sphères associées au problème (REF PLACEMENT A3). Le choix de la solution est arbitraire car la reconstruction de la bonne chiralité de la molécule ne dépend pas du placement des trois premiers atomes non fictifs.

$$a_0' \begin{cases} x_0' = 0 \\ y_0' = 0 \\ z_0' = 0 \end{cases} \quad a_1' \begin{cases} x_1' = d_{0,1} \\ y_1' = 0 \\ z_1' = 0 \end{cases} \quad a_2' \begin{cases} x_2' = \frac{d_{0,2}^2 - d_{1,2}^2 + x_1'^2}{2x_1'} \\ y_2' = \sqrt{d_{2,0}^2 - x_2'^2} \\ z_2' = 0 \end{cases} \quad a_3' \begin{cases} x_3' = \frac{d_{0,3}^2 + x_1'^2 - d_{1,3}^2}{2x_1'} \\ y_3' = \frac{-2x_3'x_2' + d_{0,2}^2 - d_{2,3}^2 + d_{0,3}^2}{2y_2'} \\ z_3' = \sqrt{-x_3'^2 - y_3'^2 + d_{0,3}^2} \end{cases}$$

FIGURE 3.8 – Placement des atomes a_0' , a_1' , a_2' et a_3'

Une fois que les quatre premiers atomes sont placés, nous leur appliquons une translation selon le vecteur \vec{a}_0 , de sorte que l'atome fictif soit à sa position originale, et que les distances relatives des atomes a_0 , a_1 , a_2 et a_3 soient toujours consistantes. Nous faisons alors appel au solveur pour résoudre les équations associées au placement des autres atomes de la molécule. Pour chaque atome, nous sélectionnons la solution respectant au mieux la distance $d_{n-1,n}$ (REF PLACEMENT AN).

$$\begin{cases} d_{n-4,n}^2 = (x_n - x_{n-4})^2 + (y_n - y_{n-4})^2 + (z_n - z_{n-4})^2 \\ d_{n-3,n}^2 = (x_n - x_{n-3})^2 + (y_n - y_{n-3})^2 + (z_n - z_{n-3})^2 \\ d_{n-2,n}^2 = (x_n - x_{n-2})^2 + (y_n - y_{n-2})^2 + (z_n - z_{n-2})^2 \end{cases}$$

FIGURE 3.9 – Équations de sphères permettant d'obtenir la position d'un atome quelconque de la molécule

Limites de l'approche par solveur L'utilisation d'un solveur calculant les solutions au cas par cas pose deux problèmes importants. Le premier concerne les performances de la solution. En effet, la résolution des systèmes d'équations consomme beaucoup de ressources et prend donc un temps non négligeable si on souhaite appliquer la méthode à un grand nombre de molécules.

Le second problème est lié à la propagation des erreurs lors de la reconstruction (REF RECONSTRUCT TRI-LAT). À cause du manque de précision de certaines valeurs, certaines intersections de sphères sont vides. Le solveur renvoie alors des solutions imaginaires que nous ne pouvons pas interpréter. Ce problème se manifeste avant tout sur les molécules de taille importante, mais il est impossible de déterminer une taille limite au delà de laquelle nous ne pouvons pas reconstruire les molécules. Cela implique qu'il existe des molécules que nous ne pouvons pas reconstruire, et que nous ne pouvons pas déterminer à l'avance si une molécule donnée peut être reconstruite.

3.2.3.3 Reconstruction automatique des positions en utilisant des équations de trilatération

Afin de pallier les problèmes liés à l'utilisation d'un solveur pour construire l'ensemble des positions des atomes d'une molécule à partir de la matrice réduite des distances inter-atomiques, nous utilisons une méthode permettant de calculer les positions de chaque point à partir d'un ensemble d'équations. Cette méthode est décrite sur Wikipédia⁴. Il s'agit d'une méthode de trilatération de points, c'est à dire que l'on cherche à déterminer la position d'un point en fonction de ses distances à trois points dont les positions sont connues, par opposition à la triangulation⁵ pour laquelle on détermine la position d'un point en fonction de ses angles à des points dont

4. <https://en.wikipedia.org/wiki/Trilateration>

5. <https://fr.wikipedia.org/wiki/Triangulation>

les positions sont connues.

De même que pour la méthode utilisant un solveur (REF SOLV), nous commençons par placer l'atome fictif a_0 à la position qui lui a été attribuée, puis les atomes a_1 , a_2 et a_3 de façon arbitraire telle que les distances relatives des atomes a_i , $i \in \{0, \dots, 3\}$ soient respectées. Nous utilisons pour cela les équations décrites en (REF FIG PLACEMENT).

Une fois les quatre premiers atomes placés, nous cherchons à placer l'atome a_n de la molécule en fonction de ses distances aux quatre atomes précédents. Nous calculons les solutions en considérant que a_{n-4}' est à l'origine du repère, que a_{n-3}' est sur l'axe x , et que a_{n-2}' est sur le plan tel que $z = 0$, puis nous effectuons une translation des solutions dans le système de coordonnées original. Pour cela, nous définissons les quantités et vecteurs suivants.

La notation \hat{u} indique un vecteur u de norme 1, et nous considérons que $\overline{a_i}$ représente le vecteur allant de l'origine au point a_i , dans le but de simplifier l'écriture des équations.

Vecteur unitaire dans la direction de a_{n-4} à a_{n-3} :

$$\hat{e}_x = \frac{\overline{a_{n-3}} - \overline{a_{n-4}}}{d_{n-4,n-3}}$$

Ordre de grandeur signé de la composante x dans le nouveau système de coordonnées du vecteur $\overline{a_{n-4}a_{n-2}}$:

$$i = \hat{e}_x \cdot (\overline{a_{n-4}} - \overline{a_{n-2}})$$

Vecteur unitaire dans la direction y par rapport à \hat{e}_x :

$$\hat{e}_y = \frac{\overline{a_{n-2}} - \overline{a_{n-4}} - i\hat{e}_x}{\|\overline{a_{n-2}} - \overline{a_{n-4}} - i\hat{e}_x\|}$$

Vecteur unitaire dans la direction z par rapport à \hat{e}_x et \hat{e}_y :

$$\hat{e}_z = \hat{e}_x \times \hat{e}_y$$

Ordre de grandeur signé de la composante y dans le nouveau système de coordonnées du vecteur $\overline{a_{n-4}a_{n-2}}$:

$$j = \hat{e}_y \cdot (\overline{a_{n-4}} - \overline{a_{n-2}})$$

On calcule alors les deux solutions pour a_n' selon les équations suivantes.

$$a_n' \begin{cases} x_n' = \frac{d_{n-4,n}^2 - d_{n-3,n}^2 + d_{n-4,n-3}^2}{2d_{n-4,n-2}} \\ y_n' = \frac{d_{n-4,n}^2 - d_{n-2,n}^2 + i^2 + j^2}{2j} - \frac{i}{j}x_n' \\ z_n' = \pm \sqrt{d_{n-4,n}^2 - x_n'^2 - y_n'^2} \end{cases}$$

Enfin, nous translatons les deux solutions a_n' dans le système de coordonnées original selon le vecteur suivant.

$$\bar{p} = \overline{a_{n-4}} + x_n'\hat{e}_x + y_n'\hat{e}_y + z_n'\hat{e}_z.$$

Nous obtenons alors deux solutions a_n , et nous sélectionnons celle telle que la distance $d_{n-1,n}$ est la plus cohérente.

3.3 Matrice réduite des distances à des points fixes

3.4

Chapitre 4

Modèles prédictifs

Chapitre 5

Conclusion

Annexes