# Macromolecular Docking - project seminar

## Peter Mühlbacher - a<br/>1253030 $\,$

## May 17, 2015

## Contents

1	Week 1		
	1.1	The Original Goal	2
<b>2</b>	Week 2		
	2.1	Benchmarks	3
	2.2	PDB	3
	2.3	Modelling the Potential	4
3	Week 3		
	3.1	Ordering Dihedral Angles by Variance	5
	3.2	Rotation About an Arbitrary Axis in 3D	
4	Week 4		
	4.1	Ordering Dihedral Angles by Variance (continued)	7
5	· · · · · · · · · · · · · · · · · · ·		
	5.1	The Node Class	7
	5.2	Absolute to Internal Coordinates	8
	5.3	The NeRF Algorithm	9
	5.4	Calculating the Children's Absolute Coordinates	10

#### Abstract

Weekly reports on advances, encountered difficulties and implementations for the seminar on artificial intelligence will be gathered in this document.

#### Goals

In this seminar I aim to examine the problem of (flexible) protein-protein docking.

#### 1 Week 1

#### 1.1 The Original Goal

Originally I wanted to approach the problem of protein folding by investigating the asymptotic  $(t \to \infty)$  behaviour of the probability density function of a given protein. As the changes over time of a given protein can be modelled by a Langevin equation  $\dot{x} = -\nabla U(x) + \beta \dot{w}$ , the corresponding probability density function p(t,x) is described by the forward Fokker-Planck equation  $\frac{\partial p}{\partial t} = \ddot{\beta} \Delta p + \text{div}(p \nabla U)$  whose asymptotic behaviour has been investigated in Nadler et al. (2008).

However, there are (at least) two reasons why this does not work:

- In order to get a meaningful low dimensional representation of the system, a considerable margin between two adjacent eigenvalues  $\lambda_{k+1} \gg \lambda_k$  of the Fokker-Planck operator is needed and it can be shown that this roughly corresponds to having a potential function U with k local minima<sup>2</sup>. Taking these properties into consideration it seems futile to employ this method of dimensionality reduction as the potential function in the area of protein folding has thousands of local minima.
- Even if solutions could be found explicitly in a reasonable amount of time the problem of finding and exact potential function U still remains unsolved<sup>3</sup>. As we are working with approximations it does not seem to make any sense trying to study its asymptotic long-term behaviour.

As a result of above considerations my topic of this seminar was restricted to docking. In the special case of rigid docking both molecules are assumed to be in a metastable state which drastically cuts the complexity of the state space.

<sup>&</sup>lt;sup>1</sup>understood as a single point in high dimensional space

<sup>&</sup>lt;sup>2</sup>c.f. (Nadler et al., 2008, pp.9), p.9

<sup>&</sup>lt;sup>3</sup>c.f. Neumaier (2006), p.14ff

#### 2 Week 2

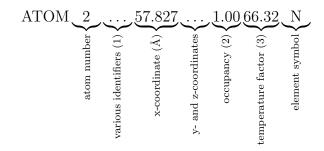
#### 2.1 Benchmarks

Visit http://zlab.umassmed.edu/benchmark/ for a benchmark on various docking problems that are ordered by difficulty. There are test cases for rigid docking as well as flexible docking, again ordered by amount of typically changing parameters like torsion angles.

#### 2.2 PDB

The *Protein Data Bank Format* is a plain text data file, storing (amongst other information) the coordinates of every single atom in a protein.

A typical *PDB*-file's row might look like this:



- 1. These identifiers are not of importance for this project.
- 2. In case an atom has been found in more than one locations all the different coordinates and according occupancies are listede.g. if it has been found in two different locations, both locations are represented as distinct rows with an occupancy factor of 0.5 each.
- 3. The temperature factor is an indication of how rigidly the atom is held in place in the given structure. Lower values imply that it is a more "stable" position.

For a list of PDB files for various proteins visit http://www.rcsb.org/.

#### 2.3 Modelling the Potential

Given two metastable molecules  $M_i = (\mathbf{x_i}, \alpha_i, \beta_i, \gamma_i, \{\mathbf{a_{ij}}\}_{j=1}^{n_i}), i \in \{1, 2\}$ . In the following  $\mathbf{x_i} \in \mathbb{R}^3$  will be called the molecule's (or the system's) point of reference,  $\alpha_i, \beta_i, \gamma_i \in [0, 2\pi)$  the system's rotation (where  $\alpha_i$  corresponds to the rotation around the x-axis relative to its point of reference; analogously for  $\beta_i, \gamma_i$ ) and  $\mathbf{a_{ij}} \in \mathbb{R}^3$  the system's atoms with (absolute) coordinates  $\mathbf{x_i} + \mathbf{a_{ij}}$ , where  $n_i$  is the number of atoms.

The Lennard-Jones potential of a system consisting of two points  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^3$  with distance  $r^2 = ||a - b||$  is defined as

$$V_{LJ}(r) = \varepsilon \left( \left( \frac{\sigma}{r} \right)^{12} - 2 \left( \frac{\sigma}{r} \right)^{6} \right), \tag{1}$$

where  $\varepsilon$  is the size of the potential well and  $\sigma$  the distance at which the potential reaches its minimum. Both are usually given by experiments.

In order to investigate the energy landscape for a docking problem it is natural to consider the Lennard-Jones potential for a system consisting of two proteins  $M_1, M_2$  given by

$$\tilde{V}_{LJ}(M_1, M_2) = \sum_{i=1}^{n_1+n_2} \sum_{j>i}^{n_1+n_2} V_{LJ}(r_{ij})$$
(2)

if we set  $r_{ij}$  to be the distance between the *i*th and *j*th atom<sup>4</sup>.

If we want to explicitly calculate its gradient with respect to the free parameters  $\mathbf{x}_1, \alpha_1, \beta_1, \gamma_1$  it is instructive to write down the entire formula as a function of those variables. First of all, to faciliate computation, we split up the summation into three parts  $\tilde{V}_{LJ}(M_1, M_2) = \tilde{V}_{LJ}^1(M_1) + \tilde{V}_{LJ}^2(M_2) + \tilde{V}_{LJ}^3(M_1, M_2)$ , where  $\tilde{V}_{LJ}^1$  refers to all combinations of atoms of  $M_1$ ,  $\tilde{V}_{LJ}^2$  to those of  $M_2$  and  $\tilde{V}_{LJ}^3$  to the summation over all pairs  $(a_{1i}, a_{2j})_{i \in \{1, \dots, n_1\}, j \in \{1, \dots, n_2\}}$  It is apparent that  $\tilde{V}_{LJ}^1$  and  $\tilde{V}_{LJ}^2$  do not change as a function of the arguments chosen above (however, they would in case we also introduced variable torsion angles and such). To calculate the gradient it thus suffices to inspect  $\tilde{V}_{LJ}^3$  as a function of the free parameters  $x, y, z, \alpha, \beta, \gamma$  describing the position of  $M_1$ :

$$\tilde{V}_{LJ}^{3}(\mathbf{x}_{1}, \alpha_{1}, \beta_{1}, \gamma_{1}) = \varepsilon \sum_{i=1}^{n_{1}} \sum_{j=1}^{n_{2}} \left(\frac{\sigma}{r_{ij}}\right)^{12} - 2\left(\frac{\sigma}{r_{ij}}\right)^{6}, \tag{3}$$

 $<sup>4</sup>i, j \in \{1, \dots, n_1 + n_2\}$ 

where  $r_{ij}^2 = \|\tilde{a}_{1i} - a_{2j}\|$  and  $\tilde{a}_{1i}$  are the first system's *i*th particle's absolute coordinates given by

$$R_x(\alpha_1)R_y(\beta_1)R_z(\gamma_1)\mathbf{a}_{1i} + \mathbf{x}_1,$$

where

$$R_x(\alpha_1) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \alpha_1 & -\sin \alpha_1 \\ 0 & \sin \alpha_1 & \cos \alpha_1 \end{pmatrix},$$

$$R_y(\beta_1) = \begin{pmatrix} \cos \beta_1 & 0 & -\sin \beta_1 \\ 0 & 1 & 0 \\ \sin \beta_1 & 0 & \cos \beta_1 \end{pmatrix}$$

and

$$R_z(\gamma_1) = \begin{pmatrix} \cos \gamma_1 & -\sin \gamma_1 & 0\\ \sin \gamma_1 & \cos \gamma_1 & 0\\ 0 & 0 & 1 \end{pmatrix}.$$

Now one can explicitly compute the partial derivatives  $\frac{\partial \tilde{V}_{LJ}}{\partial p} = \frac{\partial \tilde{V}_{LJ}^3}{\partial p}$ , for p being one of the parameters.

As of now we can compute the potential  $\tilde{V}_{LJ}(\mathbf{X})$  for a given point  $\mathbf{X} = (\mathbf{x}_1, \alpha_1, \beta_1, \gamma_1)$  in the state space and its gradient  $\nabla \tilde{V}_{LJ}(\mathbf{X})$  and that is all we need for algorithms like steepest descent (which would still require rescaling of the parameters as to avoid extensive zig-zagging) or diffusion maps.

Furthermore, note that the computation of the potential and the gradient can be combined into one double loop, making<sup>5</sup> it  $O((n_1 + n_2)^2)$ .

## 3 Week 3

### 3.1 Ordering Dihedral Angles by Variance

To be able to approximate the real docking process as closely as possible while only letting some chosen dihedral angles be variable, one first needs to order them accordingly. One such order might be given by regarding each such angle  $\Theta_i$  as a discrete random variable and ordering them by their variances.

 $<sup>^{5}</sup>$ although the computation of the gradient is likely to result in a big constant factor of operations

More precisely, a set of a single molecule's uniformly distributed<sup>6</sup> conformations  $\{M_k\}_{k=0}^n$  induces a discrete probability space  $(\{M_k\}_{k=0}^n, \mathscr{A}, P)$  where  $P(M_k) := \frac{V_{LJ}(M_k)}{Z}$ , where  $Z := \sum_{k=1}^n P(M_k)$  is a normalising factor and for all A in  $\mathscr{A} : P(A) = \sum_{M \in A} P(M)$ .

Another approach might be to compute sample trajectories and use them (regardless of their distribution) instead of uniformly distributed samples.

For a fixed  $\Theta_i$  we obtain the expected value  $E[\Theta_i] = \sum_{k=1}^n P(M_k)\Theta_{ik}$ , were  $\Theta_{ik}$  is the k-th conformation's i-th dihedral angle. However, this approach does not consider the fact that these values are circularly distributed; thus it may be better to assume that the angles' values are normally distributed and to estimate the parameters of the von Mises distribution.

The variance  $Var[\Theta_i]$  is given by  $E[(\Theta_i - E[\Theta_i])^2]$ . Other ideas included:

- Looking if such a list already existed for some molecules, but I couldn't find any.
- Taking different conformations from rcsb.com and using PCA on the phase space consisting of the dihedral angles, but given the number of dimensions compared to the amount of samples available<sup>7</sup> this is not a promising approach either.
- Take Ramachandran plots, interpret them as conditional probabilities and combine all of them. The data for this approach should be available, but combining all those single results may be hard.

### 3.2 Rotation About an Arbitrary Axis in 3D

To simulate flexible docking it is essential to be able to determine the single atoms' Cartesian coordinates, given one point of reference, the distance between two adjacent atoms, the angles of three consecutive atoms and the dihedral angles.

Given four atoms  $\{a_i\}_{i=0}^3$ , s.t.  $a_j$  and  $a_{j+1}$  are covalently bonded their dihedral angle is the smallest angle between the two planes  $\pi_0$  and  $\pi_1$  which are uniquely determined by  $\{a_i\}_{i=0}^2$  and  $\{a_i\}_{i=1}^3$ , respectively.

 $<sup>^6</sup>$ This requirement is necessary because we do not want to count "very similar" conformations twice.

<sup>&</sup>lt;sup>7</sup>For albumin there are 104 entries as of 14.04.2015.

Given  $\{a_i\}_{i=0}^2$  one sets  $a_3$  in  $_0$ , s.t. it has the prescribed distance as well as the angle between  $a_1, a_2$  and  $a_3$ . Now we want to rotate  $a_3$  about the line  $a_2 + \overrightarrow{a_1 - a_2}$ .

Clearly this is not a linear mapping unless  $a_2 = 0$ , but by introducing a forth dimension we can turn three-dimensional translations into linear mappings.

See http://inside.mines.edu/fs\_home/gmurray/ArbitraryAxisRotation/for a derivation of the according formula.

#### 4 Week 4

#### 4.1 Ordering Dihedral Angles by Variance (continued)

A solution that is better suited for this problem only considers the current configuration. Thus it suffices to consider the given configuration<sup>8</sup>  $C = (\Theta_i)_{i=1}^n$ , calculate its potential energy  $V_{LJ}(C)$  and then, for every dihedral angle  $\Theta_i$ , evaluate the difference  $\Delta_i^{\varepsilon} := V_{LJ}(C) - V_{LJ}(C_i^{\varepsilon})$ , where  $C_i^{\varepsilon}$  is the configuration  $(\Theta_1, \ldots, \Theta_{i-1}, \Theta_i + \varepsilon, \Theta_{i+1}, \ldots, \Theta_n)$ .

To get an ordering of the dihedral angles' variance we set  $k <_{\varepsilon} j$  if  $\Delta_k^{\varepsilon} + \Delta_k^{-\varepsilon} < \Delta_j^{\varepsilon} + \Delta_j^{-\varepsilon}$  which clearly defines an absolute order for every  $\varepsilon$ 

## 5 Weeks 5,6,7

#### 5.1 The Node Class

To store the data I introduced a *Node* class with the following fields:

- 1. atomname(String): atom name, e.g. \_C\_\_, \_CG1, where underscores are used in lieu of whitespaces to remove ambiguities; so far this field doesn't serve any particular purpose except making debugging easier
- 2. coord(double[]): coordinates of the node given as a list (e.g. [0,0,0]; usually that's what we're calculating from the other arguments, also referred to as "internal coordinates")
- 3. dist(double): distance to the parent node (bond length)

<sup>&</sup>lt;sup>8</sup>All other parameters are assumed to be fixed.

- 4. phi(double): bond angle between this, this.parent and this.parent.parent
- 5. theta(double): dihedral angle between this, this.parent, this.parent.parent and this.parent.parent
- 6. children(Node[]): list of other node elements

In a real world applications does not seem to be the way to go since languages like Python have a recursion limit which is easily reached. In Python one can kind of avoid this problem by manually increasing this limit by import sys and sys.setrecursionlimit(n), where n should be adjusted to the molecule one is currently dealing with. It seems to be common practice<sup>10</sup> to rewrite the code to use iterative instead of recursive methods in those cases, but as I couldn't find a source stating the recursive approach is slower I kept the current, recursive methods which are easier to implement.

#### 5.2 Absolute to Internal Coordinates

The task of extracting a tree-like structure as defined in the Node class from the PDB file was much more tedious than I originally thought. The main problem was that no parent fields were in use (see above why) and thus a separate structure containing information about the current item's predecessors had to be maintained. In a linear setting this would not have caused any problems, but if there are branching points it gets harder and harder to keep track of what's happening. I am really lucky to have chosen Python as my programming language as variables are "passed by assignment" which makes everything much more convenient as if variables were passed by value. Note that the "link" created by the "passing by pointer" method for mutable objects like lists can be broken by an assignment which was used a lot in the actual implementation.

Another issue was that there are actually a lot more implicit assumptions about protein structure in the PDB files than I had previously imagined and

<sup>&</sup>lt;sup>9</sup>this.parent is only pseudo syntax, a parent field is not implemented because its recursive nature would slow down the program considerably.

 $<sup>^{10}{\</sup>rm e.g.}$  c.f. http://stackoverflow.com/questions/8177073/python-maximum-recursion-depth-exceeded

<sup>11</sup>https://docs.python.org/3/faq/programming.html# how-do-i-write-a-function-with-output-parameters-call-by-reference

from author to author further undocumented "conventions" are introduced. 12

Obvious conventions are the structure of the backbone which is always assumed to be obvious. However, when it comes to hydrogen atoms it gets complicated. In some PDB files they are ignored completely, in others some are given and then there are some where every H atom's position got an entry. Due to the inconsistency with H atoms I choose to ignore them completely.

Concerning the "atomname" field (e.g. \_C\_\_, \_CG1, \_NE2) the official documentation<sup>14</sup> was no help; what really helped was this unofficial "guide": http://haldane.bu.edu/needle-doc/new/atom-format.html and taking a look at a real PDB file<sup>15</sup> while comparing each residue to its corresponding Wikipedia entry.

Another convention is that if branches were already in use and an atom doesn't have a branch, nor a distance identifier then it is assumed to be a "child" of the previous atom.

Also, it is clear that ring-like structures like in Tryptophan are not compatible with the tree-like structures. The simple solution is to just ignore one bond and not let the ring's dihedral angles vary.

The actual implementation is not particularly interesting from a mathematical point of view as it is basically a long list of if...else branches. In case you are interested, take a look at http://htmlpreview.github.io/?https://github.com/petermuehlbacher/reports/blob/master/docking/implementation/proteindocking.html.

## 5.3 The NeRF Algorithm

As proposed in http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.83.8235&rep=rep1&type=pdf we calculate an atom's absolute coordinates. This function, depending on the atom's internal and its three ancestors' absolute coordinates, will henceforth be referred to as nerf().

 $<sup>^{12}{</sup>m c.f.}$  http://chemistry.stackexchange.com/questions/30807/about-undocumented-conventions-in-pdb-files for an example with H.

<sup>&</sup>lt;sup>13</sup>However, it is not clear what they are connected to, as there are no additional CONECT records and the usual convention when it comes to numbering branches is ignored.

 $<sup>^{14} \</sup>verb|http://www.wwpdb.org/documentation/file-format-content/format33/sect9. \\ \verb|html#ATOM|.$ 

<sup>&</sup>lt;sup>15</sup>I chose http://www.rcsb.org/pdb/files/1E7I.pdb.

### 5.4 Calculating the Children's Absolute Coordinates

#### **Algorithm 1:** buildChildren()

Note that in the actual implementation the NeRF algorithm is called with parent.dist as additional parameter to get rid of one additional computation.

Also, the first three levels' atoms are assumed to be static due to the internal coordinates' nature. In practice these first three levels would be the first residue's N, CA, CB (in case there is one) and C. Note that, in order to bypass this problem, one could have defined "imaginary" auxiliary atoms at some coordinates, fixed by convention, in order to retrieve internal coordinates like a dihedral angle for the first residue's CA, etc.

### References

Nadler, Boaz, Ronald R. Coifman, Ioannis G. Kevrekidis, Stéphane Lafon, and Mauro Maggioni (2008), "Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems." Society for Industrial and Applied Mathematics, URL http://www.wisdom.weizmann.ac.il/~nadler/Publications/diffusion\_map\_MMS.pdf. Multiscale Model. Simul.

Neumaier, Arnold (2006), "Molecular modeling of proteins and mathematical prediction of protein structure." URL http://www.mat.univie.ac.at/~neum/ms/protein.pdf.