Structural biomathematics: an overview of molecular simulations and protein structure prediction

Bernat Anton



Parc Recerca Biomèdica Barcelona



Bernat Anton Structural biomathematics



Figure: Parc de Recerca Biomèdica de Barcelona (PRBB).

<ロ> <四> <四> <四> <三</td>









Direct-Coupling Analysis for Prediction of Protein Folding

・ 回 ト ・ ヨ ト ・ ヨ ト

æ





3 Direct-Coupling Analysis for Prediction of Protein Folding

Bernat Anton Structural biomathematics

ヘロト 人間 ト ヘヨト ヘヨト

æ

All the biological information of the human body is encoded in

our DNA. Human Genome Project: Sequentiation of the whole human genome completed on 2001, by Francis Collins (*Public Project*) & Craig Venter (*Celera Genomics*).

- About 3 billion base pairs (A, C, T and G).
- Estimation of 30000 genes (around 3000bp per gene).
- Less than 2% of the genome codes for proteins.
- Unknown function for over the half of the discovered genes!

・ 同 ト ・ ヨ ト ・ ヨ ト

All the biological information of the human body is encoded in our DNA. **Human Genome Project**: Sequentiation of the whole human genome completed on 2001, by Francis Collins (*Public Project*) & Craig Venter (*Celera Genomics*).

- About 3 billion base pairs (A, C, T and G).
- Estimation of 30000 genes (around 3000bp per gene).
- Less than 2% of the genome codes for proteins.
- Unknown function for over the half of the discovered genes!

(4回) (日) (日)

All the biological information of the human body is encoded in our DNA. **Human Genome Project**: Sequentiation of the whole human genome completed on 2001, by Francis Collins (*Public Project*) & Craig Venter (*Celera Genomics*).

- About 3 billion base pairs (A, C, T and G).
- Estimation of 30000 genes (around 3000bp per gene).
- Less than 2% of the genome codes for proteins.
- Unknown function for over the half of the discovered genes!

・聞き ・ ほき・ ・ ほき

All the biological information of the human body is encoded in our DNA. **Human Genome Project**: Sequentiation of the whole human genome completed on 2001, by Francis Collins (*Public Project*) & Craig Venter (*Celera Genomics*).

- About 3 billion base pairs (A, C, T and G).
- Estimation of 30000 genes (around 3000bp per gene).
- Less than 2% of the genome codes for proteins.
- Unknown function for over the half of the discovered genes!

CAGGCGGCCTCTGAGGGAAACAGTGACTGCTACTTTGGGAATGGGTCAGCCTACCG TGGCACGCACAGCCTCACCGAGTCGGGTGCCTCCTGCCTCCCGTGGAATTCCATGAT CCTGATAGGCAAGGTTTACACAGCACAGAACCCCAGTGCCCAGGCACTGGGCCTGG GCAAACATAATTACTGCCGGAATCCTGATGGGGATGCCAAGCCCTGGTGCCACGTG CTGAAGAACCGCAGGCTGACGTGGGAGTACTGTGATGTGCCCTCCTGCTCCACCTGC GGCCTGAGACAGTACAGCCAGCCTCAGTTTCGCATCAAAGGAGGGCTCTTCGCCGA CATCGCCTCCCACCCCTGGCAGGCTGCCATCTTTGCCAAGCACAGGAGGTCGCCCGG AGAGCGGTTCCTGTGCGGGGGGCATACTCATCAGCTCCTGCTGGATTCTCTCTGCCGC CCACTGCTTCCAGGAGAGGGTTTCCGCCCCACCACCTGACGGTGATCTTGGGCAGAAC ATACCGGGTGGTCCCTGGCGAGGAGGAGGAGCAGAAATTTGAAGTCGAAAAATACATTG TCCATAAGGAATTCGATGATGACACTTACGACAATGACATTGCGCTGCTGCAGCTGA AATCGGATTCGTCCCGCTGTGCCCAGGAGAGAGCAGCGTGGTCCGCACTGTGTGCCTTC CCCCGGCGGACCTGCAGCTGCCGGACTGGACGGAGTGTGAGCTCTCCGGCTACGGC AAGCATGAGGCCTTGTCTCCTTTCTATTCGGAGCGGCTGAAGGAGGCTCATGTCAGA CTGTACCCATCCAGCCGCTGCACATCACAACATTTACTTAACAGAACAGTCACCGAC AACATGCTGTGTGCTGGAGACACTCGGAGCGGCGGGCCCCAGGCAAACTTGCACGA CGCCTGCCAGGGCGATTCGGGAGGCCCCCTGGTGTGTCTGAACGATGGCCGCATGA GTGTACACAAAGGTTACCAACTACCTAGACTGGATTCGTGACAACATGCGACCG (SEO ID NO:2)

ヘロト ヘアト ヘビト ヘビト

ъ



Bernat Anton Structural biomathematics

$\mathsf{DNA} \stackrel{\mathsf{Transcription}}{\longrightarrow} \mathsf{RNA} \stackrel{\mathsf{Translation}}{\longrightarrow} \mathsf{Protein}$

Standard genetic code									
1st	2nd base								3rd
base		U		С		Α		G	
U	UUU	(Phe/F) Phenylalanine	UCU	(Ser/S) Serine	UAU	(Tyr/Y) Tyrosine	UGU	(0	U
	UUC		UCC		UAC		UGC	(Cys/C) Cysteine	С
	UUA	(Leu/L) Leucine	UCA		UAA	Stop (Ochre)	UGA	Stop (Opal)	Α
	UUG		UCG		UAG	Stop (Amber)	UGG	(Trp/W) Tryptophan	G
с	CUU		CCU	(Pro/P) Proline	CAU	(His/H) Histidine	CGU	(Arg/R) Arginine	U
	CUC		CCC		CAC		CGC		С
	CUA		CCA		CAA	(Gln/Q) Glutamine	CGA		Α
	CUG		CCG		CAG		CGG		G
А	AUU	(IIe/I) Isoleucine	ACU	(Thr/T) Threonine	AAU	(Asn/N) Asparagine	AGU	(Ser/S) Serine	U
	AUC		ACC		AAC		AGC		С
	AUA		ACA		AAA	(Lys/K) Lysine	AGA	(Arg/R) Arginine	Α
	AUG ^[A]	(Met/M) Methionine	ACG		AAG		AGG		G
G	GUU	(Val/V) Valine	GCU	(Ala/A) Alanine	GAU	(Asp/D) Aspartic acid	GGU	(Gly/G) Glycine	U
	GUC		GCC		GAC		GGC		С
	GUA		GCA		GAA	(Glu/E) Glutamic acid	GGA		Α
	GUG		GCG		GAG		GGG		G

Bernat Anton

¹ Table taken from *Wikipedia*.

1

æ

ヘロン 人間 とくほど くほとう



Protein structure $\xrightarrow{???}$ Protein function \longrightarrow Gene function

イロト イポト イヨト イヨト

ъ



- Primary: Amino acid linear sequence.
- Secondary: α -helices and β -strands.
- Tertiary / Domains: Functionally independent part of the sequence.
- Quaternary: Multi-subunit complex of domains or proteins.

² Figure taken from C.Branden & J.Tooze, Introduction to Protein Structure. D > < 🗇 > < 🖹 > < 🖹 > 📃 🔊 🤉 🖓

Main question:

How can we find the structure of a given protein?

- Crystallography.
- Nuclear magnetic resonance spectroscopy.
- Molecular simulation.
- Prediction of structure (structural biology).

NOT AN EASY TASK!

・ 同 ト ・ ヨ ト ・ ヨ ト …

ъ

Main question:

How can we find the structure of a given protein?

- Crystallography.
- Nuclear magnetic resonance spectroscopy.
- Molecular simulation.
- Prediction of structure (structural biology).

NOT AN EASY TASK!

・ 同 ト ・ ヨ ト ・ ヨ ト …

1





3 Direct-Coupling Analysis for Prediction of Protein Folding

Bernat Anton Structural biomathematics

ヘロト 人間 ト ヘヨト ヘヨト

æ



3

Structural biomathematics

<ロト <回 > < 注 > < 注 > 、

æ



And these are not the only forces and energies implied in a molecular simulation!

PLC- β 2 simulation

This simulation lasts around 20*ns*, with timesteps of $4fs^4$, using the ACEMD software with the AMBER forcefield. The simulation has been visualized using VMD software.

The protein has 708 amino acids, for a total of around 150000 atoms in the simulation (counting water and lipid molecules).

In the simulation can be observed the folding of the X/Y linker in order to cover the hydrophobic active site of the protein.

Afinsen's Dogma

The native structure of a protein is unique and is determined only by it's amino acid sequence. The folding to its native state is almost spontaneous.

Levinthal's Paradox

Due to the huge number of degrees of freedom in an unfolded protein, the number of possible conformations is astronomically large.

Then... how can proteins fold?

- Partially folded transition states.
- Funnel-like energy landscapes.
- ...?

ヘロト ヘ戸ト ヘヨト ヘヨト

Afinsen's Dogma

The native structure of a protein is unique and is determined only by it's amino acid sequence. The folding to its native state is almost spontaneous.

Levinthal's Paradox

Due to the huge number of degrees of freedom in an unfolded protein, the number of possible conformations is astronomically large.

Then... how can proteins fold?

- Partially folded transition states.
- Funnel-like energy landscapes.
- •...?

・ 同 ト ・ ヨ ト ・ ヨ ト





Oirect-Coupling Analysis for Prediction of Protein Folding

ヘロト 人間 ト ヘヨト ヘヨト

ъ

Let X, Y be two (discrete) random variables.

- The (self-)information of X is I(X) = -log(P(X)).
- The **entropy** of *X* is the measure of uncertainty associated with *X*: S(X) = E(I(X)).
- The mutual information of X and Y (also called **Kullback-Leibler divergence**) is

$$MI(X; Y) = \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \left(\frac{P(x, y)}{P(x)P(y)} \right)$$

Maximum Entropy Principle

Given a proposition that expresses testable information, the probability distribution that best represents the current state of knowledge is the one with largest entropy.

・ロト ・ 理 ト ・ ヨ ト ・

ъ



Figure: Multiple Sequence Alignment (MSA) for aaTHEP1.

Bernat Anton Structural biomathematics

くロト (過) (目) (日)

æ

From the previous MSA let's define:

 $f_i(A) =$ 'frequency of apparitions of aa A in the position i of the MSA' $f_{i,j}(A, B) =$ 'frequency of simultaneous apparitions of aa A and B in respective positions i and j of the MSA'

$$MI_{i,j} = \sum_{A,B} f_{i,j}(A,B) ln\left(\frac{f_{i,j}(A,B)}{f_i(A)f_j(B)}\right)$$

Be careful!

By definition, this mutual information of these frequencies is local in the amino acid chain, thus is noised by transitivity of correlations.

イロト 不得 とくほ とくほとう

From the previous MSA let's define:

 $f_i(A) =$ 'frequency of apparitions of aa A in the position i of the MSA' $f_{i,j}(A, B) =$ 'frequency of simultaneous apparitions of aa A and B in respective positions i and j of the MSA'

$$MI_{i,j} = \sum_{A,B} f_{i,j}(A,B) ln\left(\frac{f_{i,j}(A,B)}{f_i(A)f_j(B)}\right)$$

Be careful!

By definition, this mutual information of these frequencies is local in the amino acid chain, thus is noised by transitivity of correlations.

æ

We want

$$P(A_1,\ldots,A_L)$$

a general model for the probability of a particular amino acid sequence $A_1 \dots A_L$ to be member of the iso-structural family under consideration, and such that

$$P_i(A) \approx f_i(A),$$

 $P_{i,j}(A,B) \approx f_{i,j}(A,B),$

where

$$\begin{array}{lll} P_i(A) & = & \sum\limits_{A_k \neq A} P(A_1, \ldots, A_L), \\ P_{i,j}(A,B) & := & \sum\limits_{A_k \neq A, B} P(A_1, \ldots, A_L). \end{array}$$

Many distributions satisfying this: Maximum Entropy Principle!!!!

イロト イポト イヨト イヨト

We want

$$P(A_1,\ldots,A_L)$$

a general model for the probability of a particular amino acid sequence $A_1 \dots A_L$ to be member of the iso-structural family under consideration, and such that

$$P_i(A) \approx f_i(A),$$

 $P_{i,j}(A,B) \approx f_{i,j}(A,B),$

where

$$\begin{array}{lll} P_i(A) & = & \sum\limits_{A_k \neq A} P(A_1, \ldots, A_L), \\ P_{i,j}(A,B) & := & \sum\limits_{A_k \neq A, B} P(A_1, \ldots, A_L). \end{array}$$

Many distributions satisfying this: Maximum Entropy Principle!!!!

くロト (過) (目) (日)

Optimization problem:

maximize
$$S = -\sum_{A_i \mid i=1,...,L} P(A_1,...,A_L) ln P(A_1,...,P_L)$$

subject to $P_{i,j}(A,B) = f_{i,j}(A,B)$
 $P_i(A) = f_i(A)$

Solution: *disordered Q*-state *Potts model*

$$P(A_1,\ldots,A_L) = \frac{1}{\mathcal{Z}} exp\left\{\sum_{1 \leq i < j \leq L} e_{i,j}(A_i,A_j) + \sum_{1 \leq i \leq L} h_i(A_i)\right\}$$

where:

- the parameters e_{i,j}(A_i, A_j), h_i(A_i) are the Lagrange multipliers of the system,
- \mathcal{Z} is the normalization constant (*partition function*).

ヘロト 人間 とくほとく ほとう

3

Optimization problem:

maximize
$$S = -\sum_{A_i \mid i=1,...,L} P(A_1, ..., A_L) ln P(A_1, ..., P_L)$$

subject to $P_{i,j}(A, B) = f_{i,j}(A, B)$
 $P_i(A) = f_i(A)$

Solution: disordered Q-state Potts model

$$P(A_1,\ldots,A_L) = \frac{1}{\mathcal{Z}} exp\left\{\sum_{1 \leq i < j \leq L} e_{i,j}(A_i,A_j) + \sum_{1 \leq i \leq L} h_i(A_i)\right\}$$

where:

- the parameters e_{i,j}(A_i, A_j), h_i(A_i) are the Lagrange multipliers of the system,
- \mathcal{Z} is the normalization constant (*partition function*).

イロト 不得 とくほと くほとう

э.

Geometrically, this probability distribution is given by the Boltzmann-Gibbs distribution:

$$P(A_1,\ldots,A_L)=\frac{1}{\mathcal{Z}}e^{-\mathcal{H}(A_1,\ldots,A_L)}$$

Formally, the marginals of this distribution are obtained from

$$\frac{\frac{\partial \ln \mathcal{Z}}{\partial h_i(A)}}{\frac{\partial^2 \ln \mathcal{Z}}{\partial h_i(A)\partial h_j(B)}} = -P_{i,j}(A, B) + P_i(A)P_j(B)$$

but the direct computation is computationally prohibitive.

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Geometrically, this probability distribution is given by the Boltzmann-Gibbs distribution:

$$P(A_1,\ldots,A_L)=\frac{1}{\mathcal{Z}}e^{-\mathcal{H}(A_1,\ldots,A_L)}$$

Formally, the marginals of this distribution are obtained from

$$\frac{\frac{\partial \ln \mathcal{Z}}{\partial h_i(A)}}{\frac{\partial^2 \ln \mathcal{Z}}{\partial h_i(A)\partial h_j(B)}} = -P_{i,j}(A, B) + P_i(A)P_j(B)$$

but the direct computation is computationally prohibitive.

ヘロン 人間 とくほ とくほ とう

3

The Lagrange multipliers can be obtained using **Mean Field Aproximation** technique⁵:

Introduce a new parameter *α* in the partition function (via the disturbed Hamiltonian):

$$\mathcal{H}(\alpha) = \sum_{i=1,\dots,L} exp\left\{ \alpha \sum_{1 \le i < j \le L} e_{i,j}(A_i, A_j) + \sum_{1 \le i \le L} h_i(A_i) \right\}$$

• Consider the Legendre transform of the *Gibbs free energy* $\mathcal{F} = -\ln \mathcal{Z}(\alpha)$ (*Gibbs potential*):

$$\mathcal{G}(\alpha) = ln \mathcal{Z}(\alpha) - \sum_{i=1,\dots,L-A} \sum_{A} h_i(A) P_i(A).$$

⁵Plefka, T., Convergence condition of the TAP equation for the infinite-ranged Ising spin glass model (1982) 📃 🗠 9.9.0

The Lagrange multipliers can be obtained using **Mean Field Aproximation** technique⁵:

Introduce a new parameter *α* in the partition function (via the disturbed Hamiltonian):

$$\mathcal{H}(\alpha) = \sum_{i=1,\dots,L} exp\left\{ \alpha \sum_{1 \le i < j \le L} e_{i,j}(A_i, A_j) + \sum_{1 \le i \le L} h_i(A_i) \right\}$$

• Consider the Legendre transform of the *Gibbs free energy* $\mathcal{F} = -\ln \mathcal{Z}(\alpha)$ (*Gibbs potential*):

$$\mathcal{G}(\alpha) = ln \mathcal{Z}(\alpha) - \sum_{i=1,...,L} \sum_{A} h_i(A) P_i(A).$$

⁵Plefka, T., Convergence condition of the TAP equation for the infinite-ranged Ising spin glass_model (1982) 🚊 🗠 🔍 🖓

• Considering the *empirical connected correlation matrix*:

$$C_{i,j}(\boldsymbol{A},\boldsymbol{B}) = f_{i,j}(\boldsymbol{A},\boldsymbol{B}) - f_i(\boldsymbol{A})f_j(\boldsymbol{B}).$$

As a consecuence of the functional form of the Legendre transform

$$\begin{aligned} h_i(A) &= \frac{\partial \mathcal{G}(\alpha)}{\partial P_i(A)} \\ (C^{-1})_{i,i}(A,B) &= \frac{\partial h_i(A)}{\partial P_j(B)} = \frac{\partial^2 \mathcal{G}(\alpha)}{\partial P_i(A) \partial P_j(B)} \end{aligned}$$

 Expand the Gibbs potential up to first order Taylor expansion around α = 0:

$$\mathcal{G}(\alpha) \approx \mathcal{G}(\mathbf{0}) + \alpha \frac{\partial \mathcal{G}(\alpha)}{\partial \alpha}|_{\alpha = \mathbf{0}}$$

イロト イポト イヨト イヨト 三日

A computation over the two terms of the Taylor expansion of \mathcal{G} leads us to an expression which is easily derivable. First and second derivatives with respect the marginal distributions $P_i(A)$ provide self-consistent equations for the local fields, from which we obtain

$$(C^{-1})_{i,j}(A,B)_{|\alpha=0} = -e_{i,j}(A,B), \text{ for } i \neq j.$$

Finally, the parameters h_i can be estimated imposing empirical single-site frequency counts as marginal distributions and considering gauge conditions:

$$f_i(A) = \sum_B P_{i,j}(A, B).$$

くロト (過) (目) (日)

A computation over the two terms of the Taylor expansion of \mathcal{G} leads us to an expression which is easily derivable. First and second derivatives with respect the marginal distributions $P_i(A)$ provide self-consistent equations for the local fields, from which we obtain

$$(C^{-1})_{i,j}(A,B)_{|\alpha=0} = -e_{i,j}(A,B), \text{ for } i \neq j.$$

Finally, the parameters h_i can be estimated imposing empirical single-site frequency counts as marginal distributions and considering gauge conditions:

$$f_i(A) = \sum_B P_{i,j}(A, B).$$

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶

This leads us to the effective pair probabilities

$$\mathcal{P}_{i,j}^{Dir}(A,B) = rac{1}{\mathcal{Z}_{i,j}} exp\left\{e_{i,j}(A,B) + \tilde{h}_i(A) + \tilde{h}_j(B)
ight\}.$$

From which we can define its Kullback-Leibler divergence, that will be called **Direct Information**:

$$DI_{i,j} := \sum_{A,B} P_{i,j}^{Dir}(A,B) ln\left(rac{P_{i,j}^{Dir}(A,B)}{f_i(A)f_j(B)}
ight).$$

イロト イポト イヨト イヨト

3

And now... what?

- Depending on the previous, not an unique folding for the protein is possible. We must remove knotted structures (*Alexander polynomial* or *Heegaard Floer homology*).
- A scoring method over the resulting foldings must be defined, in order to decide which one of the structures is better.
- A short simulation of the system may be run in order to optimize the energies of the folding.

< 回 > < 回 > < 回

And now... what?

- Depending on the previous, not an unique folding for the protein is possible. We must remove knotted structures (*Alexander polynomial* or *Heegaard Floer homology*).
- A scoring method over the resulting foldings must be defined, in order to decide which one of the structures is better.
- A short simulation of the system may be run in order to optimize the energies of the folding.

< 回 > < 三 > < 三

And now... what?

- Depending on the previous, not an unique folding for the protein is possible. We must remove knotted structures (*Alexander polynomial* or *Heegaard Floer homology*).
- A scoring method over the resulting foldings must be defined, in order to decide which one of the structures is better.
- A short simulation of the system may be run in order to optimize the energies of the folding.

A (1) > A (2) > A

And now... what?

- Depending on the previous, not an unique folding for the protein is possible. We must remove knotted structures (*Alexander polynomial* or *Heegaard Floer homology*).
- A scoring method over the resulting foldings must be defined, in order to decide which one of the structures is better.
- A short simulation of the system may be run in order to optimize the energies of the folding.

・ 戸 ・ ・ ヨ ・ ・



Figure: Chewbacca mounted on a squirrel wants to thank you for your assistance!

・ロト ・ ア・ ・ ヨト ・ ヨト

ъ