# From genetic coding

to

protein structure prediction

May 2014

Jean-Luc Jestin

Département de Biologie Structurale et Chimie Institut Pasteur Biological chemistry Molecular biology

Quantitative approaches

**Mathematics** 

Arithmetic Basic number theory

Genetic code

Error-correcting codes Coding theory

A method to identify rules between protein sequences and structures

#### OUTLINE

1. Basic number theory, the p-adic distance and applications

#### 2. Symmetries in the genetic code

3. Applications in protein structure prediction

Conclusion

Representation of numbers by quadratic forms

X<sup>2</sup> represents 1, 4, 9 and 25 over Q

but not 2 and 3

X<sup>2</sup> and Y<sup>2</sup> represent the same numbers over Q

but not  $X^2$  and  $2Y^2$ 

#### Hasse-Minkowski theorem

For two quadratic forms to represent 0 over Q,

it is necessary and sufficient that

they represent 0 over  $Q_{\rm p}$  for all p and over R

#### **Metrics**

# $norm(a+b) \le norm(a) + norm(b)$

For ultrametric distances such as the p-adic distance, there is a stronger inequality:

 $norm(a+b) \le sup(norm(a); norm(b))$ 

#### The 3-adic distance and a 3-adic tree



# Phylogeny:

The smaller the distance, the higher the number of ancestors in common.

In: « DNA replication and mechanism » Nova, 2011

#### The codon substitutions matrix M



The « frequentist interpretation » of probabilities in statistical physics *Dill, Pressé et al., Rev. Mod. Phys. 2013* 

The same mutation, a factor in common, p-adically close numbers: the p-adic distance

Symmetries by base substitutions of degeneracy in the genetic code are symmetries of the classes of quadratic forms over Q<sub>p</sub> *Biosystems 2010*  1. Basic number theory, the p-adic distance and applications

2. Symmetries in the genetic code

# A representation of the genetic code

00.	TAC	i yi	TGC	Cys
	TAA	stop	TGA	stop
	TAG		TGG	Trp
Pro	CAT	His	CGT	Arg
	CAC		CGC	
	CAA	Gln	CGA	
	CAG		CGG	
Thr	AAT	Asn	AGT	Ser
		Lve	AGC	٨ra
	AAA AAG	LyS	AGA	Aig
Ala	GAT GAC	Asp	GGT GGC	Gly
	GAA	Glu	GGA	
	GAG		GGG	
	Pro Thr Ala	TAC TAA TAG Pro CAT CAC CAA CAG Thr AAT AAC AAA AAG Ala GAT GAC GAA GAG	TAC TAA stop TAG Pro CAT His CAC CAA GIn CAG Thr AAT Asn AAC AAA Lys AAG Ala GAT Asp GAC GAA GIu GAG	TAC TAA TAGTGC TGA TGGProCAT CAC CAC CAA CAGHis CGT CGC CGA CGGThrAAT AAC AAC AAGAsn AGT AGC AAG AAGAlaGAT GAC GAC GAGAsp GGT GGC GGA GGG

# Symmetries of degeneracy by base subtitutions

	TTT TTC	Phe	TCT TCC	Ser	TAT TAC	Tyr	TGT TGC	Cys	
(A)	TTA	Leu	TCA		TAA	stop	TGA	Trp	
	TTG		TCG		TAG		TGG		
	CTT	Leu	ССТ	Pro	CAT	His	CGT	Arg	
$(\mathbf{R})$	CTC		CCC		CAC		CGC	6	
$(\mathbf{D})^{-}$	CTA		CCA		CAA	Gln	CGA		
	CTG		CCG	0	CAG		CGG		
				U					
	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser	
$(\mathbf{C})$ –	ATC		ACC		AAC		AGC		-R
$(\mathbf{C})$	ATA	Met	ACA		AAA	Lys	AGA	Arg	P
	ATG		ACG		AAG		AGG		
	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly	
$(\mathbf{D})$	GTC		GCC		GAC		GGC		
	GTA		GCA		GAA	Glu	GGA		
	GTG		GCG		GAG		GGG		

Comp. Rend. Biol. 2006 Rumer, Proc. Natl. Acad. USSR 1966

# Genetic code codon assignments preserve information minimize the deleterious effects of mutations

#### **Mutations**

*Transitions at the 3d base* 

*Transversions at the 3d base* 

*Transversions at the 1st base* 

Single-base substitution leaving T as the 2d base

Single-base deletion within stop codons

Single-base substitution for side chains with square mass

*Three-bases substitutions* (GT/AC, AG/CT, GT/AC)

Conserved property
Amino acid
Metabolic pathway
Degeneracy
Hydrophobicity

Genes: full-length ORF

Kinetic energy over Q (not R)

2' or 3' –OH tRNA-aminoacylation

#### Reference

Goldberg & Wittes Science 1965 Wong PNAS 1975

C. R. Biol. 2006

Woese PNAS 1965

with Kempf FEBS Letters 1997

with Guilloux Biosystems 2012

with Soulé J Theor Biol 2007

#### Kinetic energy conservation by mutation

protein evolution by mutation Amino acid *1* with side-chain mass m

$$\mathsf{E}_1 = \frac{\mathsf{m} \mathsf{X}^2}{2} \qquad (1)$$

1/0

Amino acid 2 with side-chain mass M

$$\mathsf{E}_2 = \frac{\mathsf{M}\mathsf{Y}^2}{2} \qquad (2)$$

if 
$$M = mA^2$$
 then  $E_2 = \frac{m(AY)^2}{2}$ 

and same numbers represented by (1) and (2) over Q, i.e. kinetic energy disturbances are minimized during evolution consistently with the codon arrangement

with A.Guilloux, Biosystems 2012

1. Basic number theory, the p-adic distance and applications

2. Symmetries in the genetic code

3. Applications in protein structure prediction for beta-sheets

# Elementary protein folding step



Total energy conservation law during the elementary folding step



There is an infinite number of solutions (X,Y) over Q,  $\forall E$  if m/M is a square

### Sequences optimized for folding



SOF length (red)  $\leq$  20 Protein length (x-axis)  $\leq$  250

using pdb2 and pdb23 on mobyle.pasteur.fr

### Edge strand prediction

Open beta-sheets with 4 strands or more more than 3 amino acids-long

Two (training / test) sets

Rule: edge strands predicted from the sequence by extreme mean number of SOFs (Sequences Optimized for Folding) except if 2D-loops



Human transthyretin structure

96 predictions *p*-value prediction accuracy

< 9x10<sup>-7</sup> 74-78% improved using multiple sequence alignments

#### Edge strand prediction methods

<u>Structural bioinformatics & physical biochemistry methods:</u> strand length, hydrophobicity Sternberg & Thornton J. Mol. Biol. 1977

Von Heijne & Blomberg J. Mol. Biol. 1977

protein folding, total energy conservation edge strand prediction accuracy: 74-78% with A.Guilloux & B.Caudron, Comp. Struct. Biotech. J. 2013

Machine learning methods:

support vector machines

decision tree algorithms

edge strand prediction accuracy: 70-83%

Westhead et al. Prot. Sci. 2003 Faulon et al. J. Mol. Model. 2006

# Mathematics | Physical biochemistry

Symmetry Christophe Soulé Quadratic forms Fields Q, Q<sub>p</sub>, R Antonin Guilloux Ultrametricity p-adic distance From protein sequences to structures *Bernard Caudron* 

Protein folding Structural bioinformatics Genetic coding Molecular evolution Classical mechanics Achim

Kempf

Approximations as in physics on lengths & probabilities

# Acknowledgements

Achim Kempf

Department of Applied Mathematics Waterloo University, Canada

Christophe Soulé

Institut des Hautes Etudes Scientifiques, Bures Centre National de la Recherche Scientifique, France

Antonin Guilloux

Institut de Mathématiques de Jussieu, Paris Ecole Normale Supérieure de Lyon

**Bernard Caudron** 

Centre d'Informatique pour la Biologie Institut Pasteur, Paris