



BIOMATH 2013

International Conference on Mathematical
Methods and Models in Biosciences
Sofia, Bulgaria, 16-21 June 2013



DISTRIBUTION AT CONTINGENCY OF ALIGNMENT OF TWO LITERAL SEQUENCES



LORENTZ JÄNTSCHI &
SORANA D. BOLBOACA



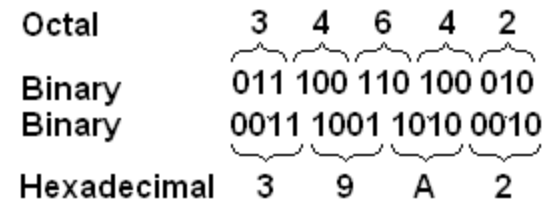
OUTLINE

2

- INTRODUCTION & AIM
- SYMILARITY
- OUR APPROACH
- LITERAL SEQUENCE CONTINGENCY
- ALIGNMENT PROBABILITY
- PDF & CDF PLOTS
- ALIGNMENT STATISTICS
- ALIGNMENT DISTRIBUTION

INTRODUCTION & AIM

- Literals sequences ← encode information
 - Written language: words, phrases
 - Computer's encoding systems:
 - binary (0,1)
 - hexadecimal (0-9, A-F)
- Genetic code:
 - DNA: A, G, C, T
 - RNA: A, G, C, U
- Proteins: amino acids



sorana.academicdirect.ro/pages/collagen/amino_acids/aa_web.htm

Amino Acid	Systemic name	Symbol		Structure	Molecular Formula	Molar weight
		3	1			
Alanine	(S)-2-Aminopropanoic Acid	Ala	A	$\begin{array}{c} \text{NH}_2 \quad \text{O} \\ \quad \\ \text{CH}_3\text{CH}-\text{C}-\text{OH} \end{array}$	C ₃ H ₇ NO ₂	89.09
Arginine	2-Amino-5-(Diaminomethylidene Amino) Pentanoic Acid	Arg	R	$\begin{array}{c} \text{NH} \quad \quad \quad \text{NH}_2 \quad \text{O} \\ \quad \quad \quad \quad \\ \text{H}_2\text{N}-\text{C}-\text{NHCH}_2\text{CH}_2\text{CH}_2\text{CH}-\text{C}-\text{OH} \end{array}$	C ₆ H ₁₄ N ₄ O ₂	174.20

AIM

4

- Hypothesis: distribution of alignments could provide useful information about the chance that a certain alignment occur or not by chance.
- Literals sequences alignment: analysis of similarity
 - Similarity of two texts (e.g. copy-paste issue)
 - Similarity of two computer encodings (e.g. copy rights)
 - Similarity of two genetic codes (biological compatibility; phylogeny)
 - Similarity of two proteins (biological functionality replacements)

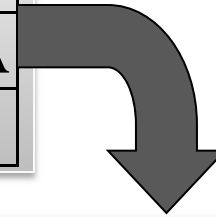
TYPES OF SYMILARITIES

5

- Exact matching
 - Gives probabilities of non-driven (or by chance) match (categorical data [1])
 - Statistical issues addressed
- Using insertions, deletions and shifts
 - Gives minimal set of insertions, deletions and shifts to match two (or more) sequences → costs (energetic for instance) for matching
 - Meta-heuristic issues addressed

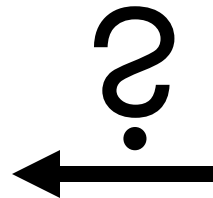
OUR APPROACH: EXACT MATCH

j	1	2	3	4	...	n-4	n-3	n-2	n-1	n
Seq1	A	C	C	G	...	T	A	G	A	C
Seq2	C	A	C	T	...	A	A	G	C	A
Seq1Seq2	AC	CA	CC	GT	...	TA	AA	GG	AC	CA
?	✗	✗	✓	✗	...	✗	✓	✓	✗	✗



Perfect alignment:
 $\Sigma AA, \Sigma CC, \Sigma GG, \Sigma TT$
 (main diagonal)

Alignment ratio:
 $(\Sigma AA + \Sigma CC + \Sigma GG + \Sigma TT) / n$



R\C	A	C	G	T	Σ
	ΣAA	ΣAC			
		ΣCC	ΣCG		
			ΣGG	ΣGT	
				ΣTT	
Σ					

LITERALS SEQUENCES CONTINGENCY (GENERAL CASE)

7

$R \setminus C$	c_1	\dots	c_q	Σ
r_1	a_{11}	\dots	a_{1q}	
\dots	\dots	\dots	\dots	
r_q	a_{q1}	\dots	a_{qq}	
Σ				

For two sequences of q literals with size n , the alignment is defined by

$$\text{SqA}(n, q) \stackrel{\text{def}}{=} \sum_{k=1}^q a_{k,k}$$

$\text{SqA}(n, q)$ is a distribution of integers ranging from 0 (no alignment) to n (perfect alignment):

$$\text{SqA}_F(0; n, q) \quad \dots \quad \text{SqA}_F(i; n, q) \quad \dots \quad \text{SqA}_F(n; n, q)$$

$$\{0, \dots, i, \dots, n\}$$

PROBABILITY OF AN CERTAIN ALIGNMENT

- Imposed conditions (alignment of U_i vs V_i , $1 \leq i \leq n$):
 - $\sum_j a_{j,k} \neq 0$ for $1 \leq k \leq q$ (all letters of the alphabet (or contingency) are present at least once in the second sequence of literals)
 - $\sum_k a_{j,k} \neq 0$ for $1 \leq j \leq q$ (all letters of the alphabet (or contingency) are present at least once in the first sequence of literals)
- A math calculation gives their frequencies:

$$\text{SqA}_F(i; n, q) = \begin{pmatrix} n - i + q^2 - q - 1 \\ q^2 - q - 1 \end{pmatrix} \begin{pmatrix} i + q - 1 \\ q - 1 \end{pmatrix}$$

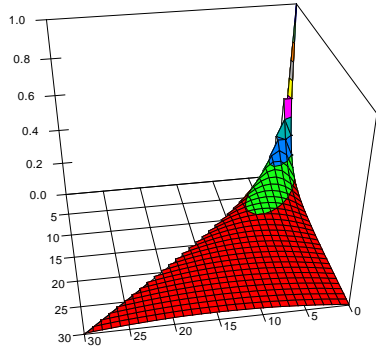
CDF & PDF

$$\text{SqA}_{\text{PDF}}(i; n, q) = \frac{\Gamma(q^2)}{\Gamma(q^2 - q)\Gamma(q)} \frac{\Gamma(n+1)}{\Gamma(n+q^2)} \cdot \frac{\Gamma(i+q)}{\Gamma(i+1)} \frac{\Gamma(n-i+q^2-q)}{\Gamma(n-i+1)}$$

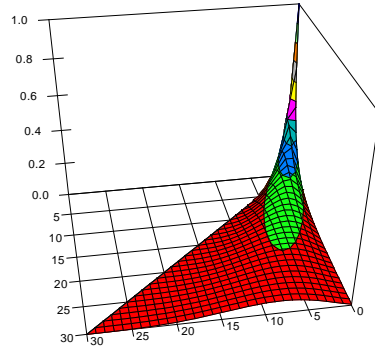
$$\text{SqA}_{\text{CDF}}(i; n, q) = \frac{\Gamma(q^2)}{\Gamma(q^2 - q)\Gamma(q)} \frac{\Gamma(n+1)}{\Gamma(n+q^2)} \sum_{j=0}^i \frac{\Gamma(j+q)}{\Gamma(j+1)} \frac{\Gamma(n-j+q^2-q)}{\Gamma(n-j+1)}$$

CDF: no close form!

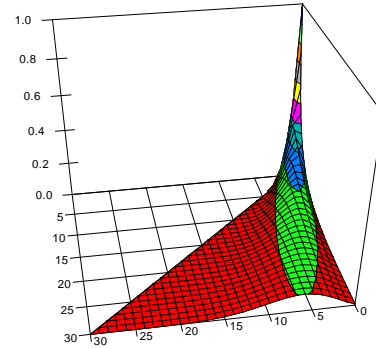
PDF & CDF PLOTS (DENSITY OF PROBABILITY VS. 'ALIGNED' (Σ MAIN DIAGONAL) AND 'TO BE ALIGNED' (N))



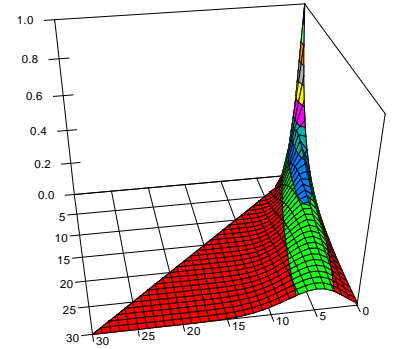
PDF2 (q = 2)



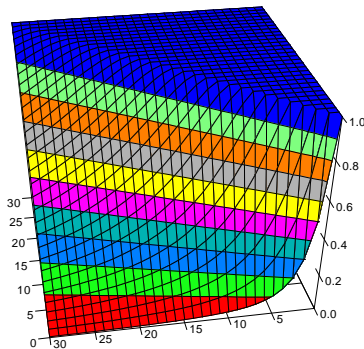
PDF3 (q = 3)



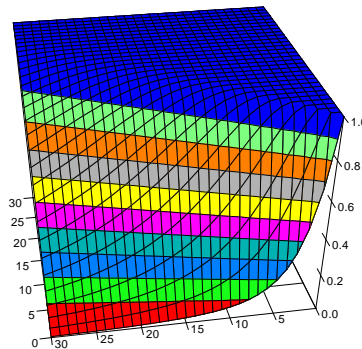
PDF4 (q = 4)



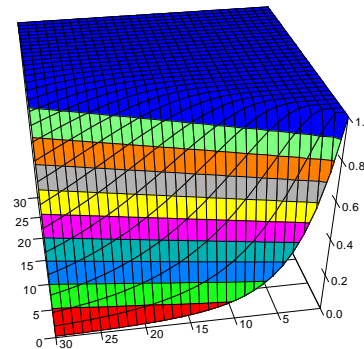
PDF5 (q = 5)



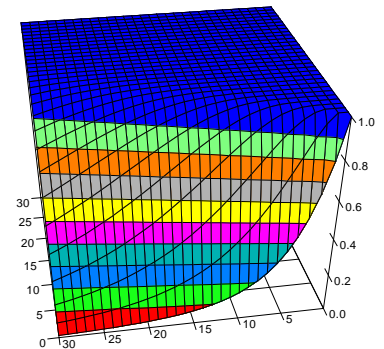
CDF2 (q = 2)



CDF3 (q = 3)



CDF4 (q = 4)



CDF5 (q = 5)

PERFECT ALIGNMENT STATISTICS

Mean

$$\mu(n, q) = \frac{n}{q}$$

Standard deviation

$$\sigma^2(n, q) = n(n + q^2) \frac{(q-1)}{q^2(q^2 + 1)}$$

Mode

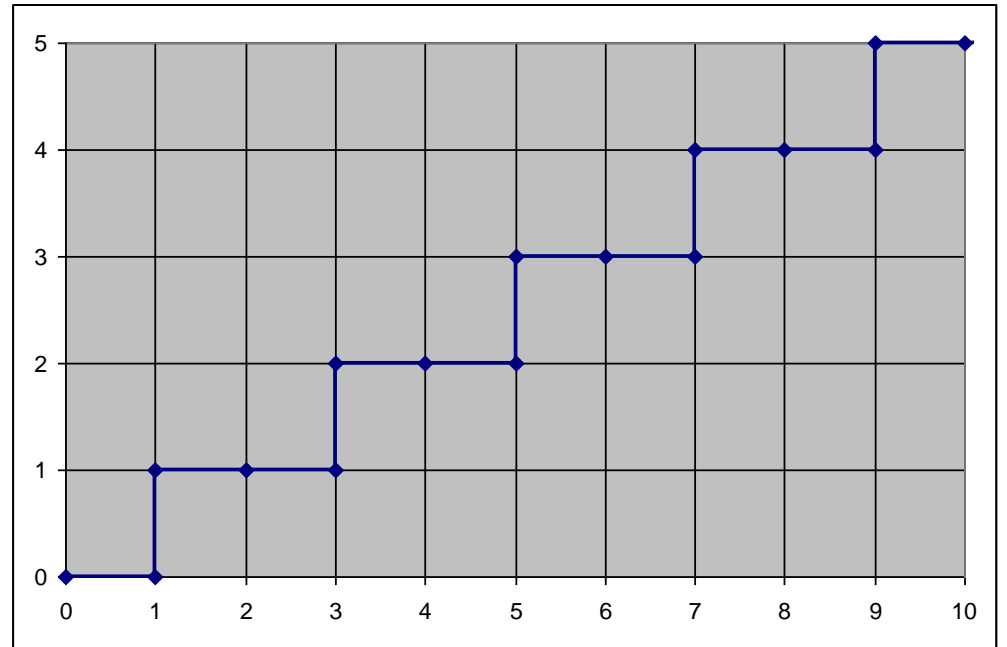
for q=2: When n is odd are two modes
(figure above)

for q=3: when $7 \cdot i = 2 \cdot n - 5$ (n = 20, 27, 41, 55, ...)

for q=4: when $14 \cdot i = 3 \cdot n - 11$ (n = 27, 83, ...)

for q=5: when $23 \cdot i = 4 \cdot n - 19$ (n = 91, ...)

in general: when $(q^2 - 2) \cdot i = (q - 1) \cdot n - (q^2 - q - 1)$



q = 2 (encoding with 0 and 1)

Mode: alignment with highest probability to be observed by chance

Figure: 'aligned' (by chance) vs. 'to be aligned' (sequences size)

OTHER DISTRIBUTION STATISTICS

$$\sigma^2(n, q) = n(n + q^2) \frac{(q - 1)}{q^2(q^2 + 1)}$$

$$\mu_3(n, q) = n(n + q^2)(2n + q^2) \frac{(q - 1)(q - 2)}{q^3(q^2 + 1)(q^2 + 2)}$$

$$\gamma_1(n, q) = 2 \frac{n + q^2/2}{\sqrt{n(n + q^2)}} \frac{q - 2}{\sqrt{q - 1}} \frac{\sqrt{q^2 + 1}}{q^2 + 2}$$

$$\mu_4(n, q) = \frac{n(n + q^2)(q - 1)(q^6 + 3(n - 2)q^5 + (3n + 5)q^4 + 3n(n - 6)q^3 + 3n(n + 6)q^2 - 18n^2(q - 1))}{q^4(q^2 + 1)(q^2 + 2)(q^2 + 3)}$$

$$\beta_2(n, q) = \frac{(q^2 + 1)(q^6 + 3(n - 2)q^5 + (3n + 5)q^4 + 3n(n - 6)q^3 + 3n(n + 6)q^2 - 18n^2(q - 1))}{n(n + q^2)(q^2 + 2)(q^2 + 3)(q - 1)}$$

$$\gamma_2(n, q) = \beta_2(n, q) - 3$$

ALIGNMENT DISTRIBUTION APPLICATIONS

13

"Alignment distribution":

- If Alignment ratio greater than the threshold of 95% (InvCDF_{95}) then with a risk smaller than 5% being in error the alignment of the sequences is not by chance.

ALIGNMENT DISTRIBUTION APPLICATIONS

14

Important points for $q=4$ (gene sequence alignment):

- Thresholds to reject the alignment by chance

n	Alignment %
8	70
13	60
21	55
39	50
282	45
$\rightarrow \infty$	44

Thank you for your attention

