

Alignment-free prediction of mycobacterial DNA promoters based on pseudo-folding lattice network or star-graph topological indices

Alcides Perez-Bello, Cristian Robert Munteanu, Florencio M. Ubeira, Alexandre Lopes de Magalhães, Eugenio Uriarte, Humberto González-Díaz

▶ To cite this version:

Alcides Perez-Bello, Cristian Robert Munteanu, Florencio M. Ubeira, Alexandre Lopes de Magalhães, Eugenio Uriarte, et al.. Alignment-free prediction of mycobacterial DNA promoters based on pseudo-folding lattice network or star-graph topological indices. Journal of Theoretical Biology, 2009, 256 (3), pp.458. 10.1016/j.jtbi.2008.09.035. hal-00554510

HAL Id: hal-00554510 https://hal.science/hal-00554510

Submitted on 11 Jan 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Author's Accepted Manuscript

Alignment-free prediction of mycobacterial DNA promoters based on pseudo-folding lattice network or star-graph topological indices

Alcides Perez-Bello, Cristian Robert Munteanu, Florencio M. Ubeira, Alexandre Lopes De Magalhães, Eugenio Uriarte, Humberto González-Díaz

PII: DOI: Reference: S0022-5193(08)00505-5 doi:10.1016/j.jtbi.2008.09.035 YJTBI 5311

To appear in:

Journal of Theoretical Biology

Received date:6 August 2008Revised date:23 September 2008Accepted date:25 September 2008

Cite this article as: Alcides Perez-Bello, Cristian Robert Munteanu, Florencio M. Ubeira, Alexandre Lopes De Magalhães, Eugenio Uriarte and Humberto González-Díaz, Alignment-free prediction of mycobacterial DNA promoters based on pseudo-folding lattice network or star-graph topological indices, *Journal of Theoretical Biology* (2008), doi:10.1016/j.jtbi.2008.09.035

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



www.elsevier.com/locate/yjtbi

Alignment-free prediction of Mycobacterial DNA promoters based on Pseudo-folding Lattice Network or Star-Graph Topological Indices

ALCIDES PEREZ-BELLO,^{1,2,3} CRISTIAN ROBERT MUNTEANU,⁴ FLORENCIO M. UBEIRA,¹ ALEXANDRE LOPES DE MAGALHÃES,⁴ EUGENIO URIARTE³ AND HUMBERTO GONZÁLEZ-DÍAZ^{1,*}

¹ Department of Microbiology and Parasitology, University of Santiago de Compostela, Santiago de Compostela 15782, Spain, <u>mpubeira@usc.es</u>, <u>humberto.gonzalez@usc.es</u>

 ² Department of Veterinary Medicine, UCLV, Santa Clara 54830, Cuba, <u>alcidopb@gmail.com</u>
 ³ Department of Organic Chemistry, Faculty of Pharmacy, University of Santiago de Compostela, Santiago de Compostela 15782, Spain, <u>eugenio.uriarte@usc.es</u>

⁴ REQUIMTE/University of Porto, Faculty of Science, Chemistry Department, Porto 4169-007, Portugal, <u>muntisa@gmail.com</u>, <u>almagalh@fc.up.pt</u>

Abstract. The importance of the promoter sequences in the function regulation of several important mycobacterial pathogens creates the necessity to design simple and fast theoretical models that can predict them. This work is proposing two DNA promoter QSAR models based on pseudo-folding Lattice Network (LN) and Star-Graphs (SG) topological indices. In addition, a comparative study with the previous RNA electrostatic-driven secondary structure folding representations have been carried out. The best model of this work was obtained with only two LN stochastic electrostatic potentials and is characterised by accuracy, selectivity and specificity of 90.87%, 82.96% and 92.95%, respectively. In addition, we pointed out the SG results dependence on the DNA sequence codification and we proposed a QSAR model based on codons and only three SG spectral moments.

Keywords: QSAR, Markov model, Mycobacterial promoters, Star graph, Lattice Network, Topological indices

**Correspondence to*: GONZÁLEZ-DÍAZ, H. Faculty of Pharmacy, University of Santiago de Compostela 15782, Spain. Email: <u>humberto.gonzalez@usc.es</u>, Tel: +34-981-563100, Fax: +34-981 594912.

1. Introduction

Protein synthesis promoter sequences play an important role in the function regulation of several important mycobacterial pathogens (Levine and Tjian, 2003; Wyrick and Young, 2002). In this sense, the prediction of the mycobacterial promoter sequences (Mps) could be interesting for the future discovery of new anti-mycobacterial drugs targets or in the study of proteins metabolism. Mycobacteria have a low transcription rate and a low RNA content per unit DNA. Thus, the transcription and translation signals in Mycobactaria may be different from those in other bacteria such as Esccherichia coli. The large variations among the characterized mycobacterial promoters suggest that the consensus sequences are not representative of these promoters. Consequently, a number of conflicting opinions regarding the presence and characteristics of consensus promoter sequences in the Mycobacteria have been presented in the literature (Mulder et al., 1997). Therefore, understanding the factors that are responsible for the low level of transcription and the possible mechanisms of regulation of gene expression in Mycobacteria, involve the examination of the mycobacterial promoter structure and the promoter transcription machinery, including chemical information about the involved RNA molecules (Arnvig et al., 2005; Harshey and Ramakrishnan, 1977). Efforts have been made to develop statistical algorithms for the sequence analysis and motif prediction by searching for homologous regions or by

2

comparing the sequence information with a consensus sequence (O'Neill and Chiafari, 1989). Wide variations existing within individual promoter sequences are primarily responsible for the unsatisfactory results yielded by the promoter-site-searching algorithms that in essence perform statistical analysis (Mulligan and McClure, 1986; Mulligan et al., 1984). Therefore, it can be inferred that the recognition of mycobacterial promoter sequences require a powerful technique that is capable of unravelling those hidden patterns in the promoter regions, which are difficult to identify directly by sequence alignment.

The Bioinformatics methods based on sequence alignment may fail in general for cases of low sequence homology between the databases query and the template sequences. The lack of function annotation (defined biological function) of the sequences used as template for function prediction constitutes another weakness of alignment approaches (Dobson and Doig, 2005; Dobson et al., 2004; Dobson et al., 2005). In addition, Chou demonstrated that the 3-dimensional structures developed based on homology modelling are very sensitive to the sequence alignment of the query protein with the structure-known protein (Chou, 2004). A group of researchers shows the growing importance of machine learning methods for predicting protein functional class independently of sequence similarity (Han et al., 2006). These methods often use as the input the 1D sequence numerical parameters, specifically defined to seek sequencefunction relationships. For instance, the so-called pseudo amino acid composition approach (Chou, 2001a; Chou, 2005) based on 1D sequence coupling numbers has been widely used to predict sub-cellular localization, enzyme family class, structural class, as well as other attributes of proteins based on their sequence similarity (Caballero et al., 2006; Chou and Shen, 2006; Du et al., 2008) Alternatively, the molecular indices that are classically used for small molecules (Aguero-Chapin et al., 2006; Liao and Wang,

2004a; Liao and Wang, 2004b; Liao and Wang, 2004c; Liao and Ding, 2005; Liao et al., 2005; Liao et al., 2006; Liu et al., 2002; Nandy, 1994; Nandy, 1996; Nandy and Basak, 2000; Randic and Vracko, 2000; Randic and Balaban, 2003; Randic and Zupan, 2004; Randic et al., 2000; Song and Tang, 2005; Woodcock et al., 1992; Zupan and Randic, 2005) have been adapted to describe the protein sequences. On the other hand, many authors have introduced 2D or higher dimension representations of sequences prior to the calculation of numerical parameters. This constitutes an important step in order to uncover useful higher-order information not encoded by 1D sequence parameters (Randic, 2004). One example of the 2D representations is the graphs used for proteins and DNA sequences. For example, the spectral-like and zigzag representations have been used suggesting an algorithm for encoding long strings of building blocks (like four DNA bases, twenty natural amino acids, or all 64 possible base triplets) (Aguero-Chapin et al., 2006). The use of the graphic approaches to study biological systems can provide useful insights, as indicated by many previous studies on a series of important biological topics, such as enzyme-catalyzed reactions (Andraos, 2008; Chou, 1981; Chou, 1989; Chou and Forsen, 1980; Chou and Liu, 1981; Chou et al., 1979; Cornish-Bowden, 1979; King and Altman, 1956; Kuzmic et al., 1992; Myers and Palmer, 1985; Zhou and Deng, 1984), protein folding kinetics (Chou, 1990), inhibition kinetics of processive nucleic acid polymerases and nucleases (Althaus et al., 1993a; Althaus et al., 1993b; Althaus et al., 1993c; Chou et al., 1994), analysis of codon usage (Chou and Zhang, 1992; Zhang and Chou, 1993; Zhang and Chou, 1994), analysis of DNA sequence (Qi et al., 2007). Moreover, graphical methods have been introduced for OSAR study (Gonzalez-Diaz et al., 2006c; Gonzalez-Diaz et al., 2007b; Prado-Prado et al., 2008) as well as utilized to deal with complicated network systems (Diao et al., 2007; Gonzalez-Diaz et al., 2007a; Gonzalez-Díaz et al., 2008). Recently, the "cellular

automaton image" (Wolfram, 1984; Wolfram, 2002) has also been applied to study hepatitis B viral infections (Xiao et al., 2006a), HBV virus gene missense mutation (Xiao et al., 2005b), and visual analysis of SARS-CoV (Gao et al., 2006; Wang et al., 2005), as well as representing complicated biological sequences (Xiao et al., 2005a) and helping to identify protein attributes (Xiao and Chou, 2007; Xiao et al., 2006b).

In this work, we are proposing a comparative study of the Mycobacterial DNA promoter prediction using pseudo-folding Lattice Network (LN) and Star-Graph (SG) topological indices. The first group of indices contains the mean stochastic electrostatic potential $({}^{LN}\xi_k)$, Markov spectral moments $({}^{LN}\pi_k)$ and Markov entropies $({}^{LN}\theta_k)$ of a Markov Model (MM) associated to a 2D network that numerically characterize DNA sequences and build a Quantitative Structure-Activity Relationships (QSAR) model to predict mycobacterial promoters sequence (Mps). The lattice-like representations (also called maps or graphs) for Mps and control group sequences (Cgs) were derived (González-Díaz et al., 2003; González-Díaz et al., 2006a; González-Díaz et al., 2005c; González-Díaz, 2007d). The ξ_k , π_k and θ_k values of several types of graphs/networks have been the base for different QSAR studies of DNA/RNA and protein sequences (Du et al., 2007a; Du et al., 2007b; Garcia-Garcia et al., 2004; Marrero-Ponce et al., 2004a; Marrero-Ponce et al., 2005b; Marrero-Ponce et al., 2004b; Meneses-Marcel et al., 2005; Santana et al., 2006). The second group of TIs is derived from the Star-Graph representations (Harary, 1969). We subsequently developed a classifier to connect Mps information (represented by the ξ_k , π_k , θ_k and star-graph TIs values) with the prediction of Cgs as Mps. The Linear Discriminant Analysis (LDA) was selected as a simple but powerful technique (González-Díaz et al., 2006b; González-Díaz, 2003a).

2. Materials and methods

2.1 Pseudo-folding Lattice Network

The first Markov Model (MM), also called MARCH-INSIDE, was used to codify the information of 135 Mps (González-Díaz et al., 2005a; González-Díaz et al., 2006a; González-Díaz et al., 2007d) and 511 random Cgs (see **Table S.1** in the supplementary material). Our methodology considers as states of the Markov Chain (MC) any atom, nucleotide or amino acid depending on the class of molecule to be described (González-Díaz et al., 2005e; González-Díaz, 2003b). Therefore, MM deals with the calculation of the probabilities $({}^{k}p_{ii})$ where the charge distribution of nucleotide moves from any nucleotide in the vicinity i at time t_0 to another nucleotide j along the protein backbone in discrete time periods until a stationary state is achieved (Yuan, 1999). As can be seen from the discussion above, we selected ${}^{LN}\xi_k$, ${}^{LN}\pi_k$ and ${}^{LN}\theta_k$ based on the utility of nonstochastic (González-Díaz and Uriarte, 2005; González-Díaz et al., 2005d; Ramos de Armas et al., 2004) and stochastic parameters (Randic and Vracko, 2000). Many researchers have demonstrated the possibility of predicting RNA from sequences (Aguero-Chapin et al., 2006) and we used 2D graphs to encode information about Mps sequences (Estrada, 2000; Estrada, 2002; Estrada and González-Díaz, 2003; González-Díaz et al., 2005b; Gonzalez and Moldes del Carmen Teran, 2004; Vilar et al., 2005; Vilar et al., 2006). This RNA 2D graphical representation is similar to those previously reported for DNA (Jacchieri, 2000; Nandy, 1994; Nandy, 1996) using four different nucleotides. The construction of the 2D lattice graph corresponding to the Mps of the gene Alpha in Mycobacterum bovis (BCG) is shown in Table 1 and Figure 1. Each nucleotide in the sequence is placed in a Cartesian 2D space starting with the first monomer at the (0, 0) coordinates. The coordinates of the successive nucleotide are calculated with the following rules:

a) Increase by +1 the abscissa axis coordinate for thymine (rightwards-step) or:

6

- b) Decrease by -1 the abscissa axis coordinate for cytosine (leftwards-step) or:
- c) Increase by +1 the ordinate axis coordinate for adenine (upwards-step) or:
- d) Decrease by -1 the ordinate axis coordinate for guanine (downwards-step).

Table 1 comes about here

Figure 1 comes about here

In the next step, we assigned to each graph a stochastic matrix ${}^{1}\Pi$. The elements of ${}^{1}\Pi$ are the probabilities ${}^{1}p_{ij}$ of reaching a node n_i with the charge Q_i moving through a walk of length of k = 1 from another node n_j with charge Q_j (Aguero-Chapin et al., 2006):

$$p_{ij} = \frac{\frac{Q_j}{d_{j0}}}{\sum_{m=l}^n \alpha_{il} \cdot \frac{Q_j}{d_{l0}}} = \frac{\varphi_j}{\sum_{m=l}^n \alpha_{il} \cdot \varphi_l}$$

$$p_j = \frac{\frac{Q_j}{d_{j0}}}{\sum_{m=l}^n \frac{Q_j}{d_{l0}}} = \frac{\varphi_j}{\sum_{m=l}^n \varphi_l}$$

$$(1)$$

where α_{ij} equals to 1 if the nodes n_i and n_j are adjacent in the graph or equal to 0 otherwise; Q_j is equal to the sum of the electrostatic charges of all nucleotide placed at this node. Note that the number of nodes (n) in the graph is equal to the number of rows and columns in ¹**H** but may be equal or even smaller than the number of DNA bases in the sequence. It then becomes straightforward calculating different types of invariant parameters for ¹**H** in order to numerically characterize the DNA sequence. In this work we calculated the following invariants:

$$^{LN}\pi_{k} = \sum_{i=j}^{n} {}^{k} p_{ij}$$
(3)

$${}^{LN}\xi_k = \sum_{i=j}^n {}^k p_j \cdot \varphi_j \tag{4}$$

$${}^{LN}\boldsymbol{\theta}_{k} = -\sum_{i=j}^{n} {}^{k}\boldsymbol{p}_{j} \cdot \log({}^{k}\boldsymbol{p}_{j})$$

$$\tag{5}$$

where ${}^{LN}\pi_k$ are the Markov spectral moments and indicate that we sum all the values in the main diagonal of the matrices ${}^{LN}\pi_k = \text{Tr}({}^k\Pi) = \text{Tr}[({}^{1}\Pi)^k]$ (Tr is the trace operator), ${}^{LN}\xi_k$ are the mean values of electrostatic potentials and ${}^{LN}\theta_k$ are the Markov entropies (González-Díaz et al., 2007a). All calculations of the ${}^{LN}\xi_k$, ${}^{LN}\pi_k$ and ${}^{LN}\theta_k$ values for the DNA sequences of both groups (Mps and Cgs) were carried out with our in-house software MARCH-INSIDE, *version 2.0* (González-Díaz et al., 2007a), including sequence representation.

2.2 Star-Graph topological indices

Each DNA sequence is a real network where the nucleotides are the vertices/nodes, connected in a specific sequence by the phosphodiester bonds. SG is an abstract representation of the real network having a dummy non-nucleotide centre and a number of "rays" equal with the nucleotide types. In the case of DNA, we can consider two codifications: the nucleotide code (as in the case of the amino acid protein sequences) and the DNA codons (the final incomplete codons are ignored). In the first codification, there are only four branches ("rays") of the star corresponding to the four types of nucleotides: adenosine (a), thymidine (t), cytidine (c) and guanosine (g). Using the codons, the DNA sequences are virtually translated into amino acid sequences that generate 21 branches, 20 standard amino acids and an extra X non-amino acid corresponding to the STOP DNA codons (Griffiths et al., 1999). Even if the promoters are not naturally translated in proteins, the second codification is useful for a comparison with the protein SG calculations. The same DNA/protein can be represented by different forms which are associated to distinct distance matrices (Randic et al., 2007). Standard star-graphs were constructed for each DNA promoter: each

nucleotide/vertex holds the position in the original sequence and the branches are labelled by the standard letters of the nucleotides (a, t, c and g). If the initial connectivity in the DNA sequence is included, the graph is embedded. In order to qualitatively evaluate the graphs, it is necessary to transform the graphical representation into correspondent connectivity matrix, distance matrix and degree matrix. In the case of the embedded graph, the matrices of the connectivity in the sequence and in the star graph are combined. These matrices and the normalized ones are the base for the calculation of the topological indices.

For a visual comparison of the lattice and star-graph representations, the same promoter sequence from **Table 1** was used to generate a standard SG based on codons that are virtually translated to amino acids (see **Table 2** and **Figure 2**).

Table 2 comes about here

Figure 2 comes about here

The star-graph topological indices are obtained with the in-house Sequence to Star Networks (S2SNet) python application. This tool can transform any character string in SG topological indices. Our recent works (Munteanu et al., 2008a; Munteanu et al., 2008b) proved the potential of S2SNet in protein QSAR models. The calculations presented in this work are characterized by embedded (E) and non-embedded (nE) TIs, no weights, Markov normalization and power of matrices/indices (n) up to 5. The result file contains the following embedded (super index "e") or non-embedded TIs (Todeschini and Consonni, 2002):

Shannon Entropy of the *n* powered Markov Matrices (${}^{SG}\theta_n$):

$${}^{\text{SG}}\theta_n^{(e)} = -\sum_i p_i * \log(p_i)$$
(6)

where p_i are the n_i elements of the *p* vector, resulted from the matrix multiplication of the powered Markov normalized matrix ($n_i \ge n_i$) and a vector ($n_i \ge 1$) with each

PTED MANUSCRIP 이 어 그

element equal to $1/n_i$;

The trace of the *n* connectivity matrices (${}^{SG}\pi_n$):

$$^{\mathrm{SG}}\pi_n^{(e)} = \sum_i (M^n)_{\mathrm{ii}} \tag{7}$$

where n = 0 – power limit, ${}^{SG}M = SG$ connectivity matrix (*i***i* dimension); *ii* = *i*th diagonal element;

Harary number (H):

$$H^{(e)} = \sum_{i < j} (m_{ij}/d_{ij})$$
(8)

where d_{ij} are the elements of the distance matrix and m_{ij} are the elements of the M nuscri connectivity matrix;

Wiener index (W):

$$W^{(e)} = \sum_{i < j} d_{ij}$$

Gutman topological index (S_6) :

$$S_6^{(e)} = \sum_{ij} \deg_i * \deg_j / d_{ij} \tag{10}$$

where deg_i are the elements of the degree matrix;

Schultz topological index (non-trivial part) (S):

$$S^{(e)} = \sum_{i < j} \left(deg_i + deg_j \right) * d \tag{11}$$

Balaban distance connectivity index (*J*):

$$J^{(e)} = (edges - nodes + 2) * \sum_{i < j} m_{ij} * \operatorname{sqrt}(\sum_{k} d_{ik} * \sum_{k} d_{kj})$$
(12)

where *nodes*+1 = AA numbers/node number in the Star Graph + origin, $\sum_k d_{ik}$ is the node distance degree;

Kier-Hall connectivity indices (^{n}X) :

$${}^{0}X^{(e)} = \sum_{i} 1 / \operatorname{sqrt}(deg_{i})$$
⁽¹³⁾

$${}^{2}X^{(e)} = \sum_{i < j < k} m_{ij} * m_{jk} / \operatorname{sqrt}(deg_{i} * deg_{j} * deg_{k})$$

$$\tag{14}$$

$${}^{3}X^{(e)} = \sum_{i < j < k < m} m_{ij} * m_{jk} * m_{km} / \operatorname{sqrt}(deg_{i} * deg_{j} * deg_{k} * deg_{m})$$
(15)

$${}^{4}X^{(e)} = \sum_{i < j < k < m < o} m_{ij} * m_{jk} * m_{km} * m_{mo} / \operatorname{sqrt}(\operatorname{deg}_{i} * \operatorname{deg}_{j} * \operatorname{deg}_{k} * \operatorname{deg}_{m} * \operatorname{deg}_{o})$$
(16)

(9)

 ${}^{5}X^{(e)} = \sum_{i < j < k < m < o < q} m_{ij} * m_{jk} * m_{km} * m_{mo} * m_{oq} / \operatorname{sqrt}(\operatorname{deg}_{i} * \operatorname{deg}_{j} * \operatorname{deg}_{k} * \operatorname{deg}_{m} * \operatorname{deg}_{o} * \operatorname{deg}_{q})$ (17) Randic connectivity index (¹XR):

$${}^{I}XR^{(e)} = \sum_{i < j} m_{ij} / \operatorname{sqrt}(deg_{i} * deg_{j})$$
(18)

The embedded and non-embedded SG TIs will be used to construct a DNA promoter classification model using the LDA statistical methods.

2.3 Linear Discriminant Analysis

LDA forward stepwise analysis from STATISTICA (StatSoft.Inc., 2002) was carried out for variable selection to build up the model (Garcia-Garcia et al., 2004; Kutner et al., 2005; Marrero-Ponce et al., 2004a; Marrero-Ponce et al., 2005b; Marrero-Ponce et al., 2004b; Meneses-Marcel et al., 2005; Santana et al., 2006). In order to decide if a DNA sequence is classified as mycobacterial promoter (Prom) or not (nProm), we added an extra dummy variable named Prom/nProm (binary values of 1/-1 for LN and 1/0 for SG) and a cross-validation variable (CV). The best cross-validation methods used are practice is the independent dataset test, the subsampling test and the jackknife test (Chou and Zhang, 1995). The jackknife test has been increasingly used by investigators to examine the accuracy of various predictors (Chen and Li, 2007; Chou and Shen, 2007a; Chou and Shen, 2008; Diao et al., 2007; Ding et al., 2007; Lin, 2008; Xiao and Chou, 2007). In the actual work, the independent data test is used by splitting the data at random in a training series (train, 75%) used for model construction and a prediction one (val, 25%) for model validation (the CV column is filled by repeating 3 train and 1 val). All of the variables included in the models were standardized in order to bring them onto the same scale. Subsequently, standardized linear discriminant equations that allow comparison of their coefficients were obtained (Chiti et al., 2003; Pawar et al., 2005).

In the case of LN, the general QSAR formula is the following:

^{LN} Mps - score =
$$a_0 + \sum_{k=0}^{5} b_k \times^{LN} \pi_k + \sum_{k=0}^{5} c_k \times^{LN} \theta_k + \sum_{k=0}^{5} d_k \times^{LN} \xi_k$$
 (19)

where ^{LN}Mps-score is the continue score value for the DNA mycobacterial promoter classification corresponding to the lattice representation, ^{LN} π_k are Markov spectral moments (traces), ^{LN} θ_k are the Markov entropies, ^{LN} ξ_k the mean stochastic electrostatic potential, b_k , c_k , d_k are the coefficients of the previous indices and a_0 is the independent term. A similar formula is defining the SG QSAR model in Eq. 20.

$$^{\text{SG}}\text{Mps-score} = e_0 + \sum_{k=0}^5 f_k \times^{\text{SG}} \pi_k + \sum_{k=0}^5 f_k^{e_k} \times^{\text{SG}} \pi_k^{e_k} + \sum_{k=0}^5 g_k \times^{\text{SG}} \theta_k + \sum_{k=0}^5 g_k^{e_k} \times^{\text{SG}} \theta_k^{e_k} + \sum_{k=0}^{10} e_k \times II_k + \sum_{k=0}^{10} e_k \times II_k^{e_k}$$
(20)

where ^{SG}Mps-score is the continue score value for the DNA mycobacterial promoter classification corresponding to the SG representation, ${}^{SG}\pi_{k}{}^{e} / {}^{SG}\pi_{k}$ and ${}^{SG}\theta_{k}{}^{e} / {}^{SG}\theta_{k}$ are embedded/non-embedded traces (Markov spectral moments) and the Shannon entropies, TI_{*k*}^{*e*} / TI_{*k*} are the other 22 standard SG embedded and non-embedded TIs (H, W, S₆, S, J, ${}^{0}X$, ${}^{2-5}X$, ${}^{1}XR$, H^e, W^e, S₆^e, S^e, J^e, ${}^{0}X^{e}$, ${}^{2-5}X^{e}$, ${}^{1}XR^{e}$), f_{k}^{e} / f_{k} , g_{k}^{e} / g_{k} and e_{k}^{e} / e_{k} are the TIs coefficients and e_{0} is the independent term. Accuracy, specificity, sensitivity, F, Wilk's (λ) statistic ($\lambda = 0$ perfect discrimination, being $0 < \lambda < 1$) were examined in order to assess the discriminatory power of the model.

3. Results and discussion

Many different parameters can be used to encode RNA sequence information and further assign or predict the function or physical properties (González-Díaz and Uriarte, 2005). The present approach involves the calculation of different sequence parameters, which can be applied to different types of molecular graphs (Aguero-Chapin et al., 2006), including DNA, RNA and proteins (Di Francesco, 1999; González-Díaz et al.,

2005c). MM has been applied successfully to Genomics and Proteomics and represents an important tool for analyzing biological sequence data. In particular, MM has been used for protein folding recognition (Chou, 2001b) and for prediction of protein signal sequences (Chou and Shen, 2007b; Van Waterbeemd, 1995). This work compared two models based on different TIs including π_k and θ_k values of the stochastic matrices ¹II(LN) and ¹II(SG) (^{SG}M) associated with LN and SG, ^{LN} ξ_k parameters of ¹II(LN) as well as classic TIs for ¹II(SG). These parameters describe the distribution of the nucleotides of the DNA sequence in the above graphs/networks. This calculation was carried out for two groups of DNA sequences, one made up of Mps and the other formed by Cgs. In addition, previous results of the RNA secondary structure (2S) QSAR are compared.

3.1. Results for DNA LN indices

In the first study of the DNA LN representations, the best QSAR equation that classifies a novel sequence as Mps or not is the following (**Table 3**):

^{LN} Mps - score =
$$-1.2 - 4.1 \times^{LN} \xi_1 + 2.1 \times^{LN} \xi_5$$
 (21)

The statistical parameters of this equation were Wilk's statistic (λ =0.95) and error level (p-level<0.001). This discriminant function misclassified only 36 cases out of 511 Cgs used, reaching a high level of accuracy of 90.87%. More specifically, the model classified correctly 112/135 (82.9%) of Mps and 475/511 (92.9%) of the control group. Conversely, the remains four descriptors ${}^{LN}\xi_{0}$, ${}^{LN}\xi_{2}$, ${}^{LN}\xi_{3}$ and ${}^{LN}\xi_{4}$ do not have a significant relationship with the Mps characteristic. The use of only six molecular descriptors to model a data set of 585 sequences prevents us by large from chance correlation. In physical terms, the above results confirm other studies about the relationship between the electrostatic potential of the DNA molecule and its biological

activity. However, in this case not all the electrostatic interactions affect the activity in the same way. Finally, long-term electrostatic interaction potentials $({}^{LN}\xi_{0}, {}^{LN}\xi_{2}, {}^{LN}\xi_{3})$ and ${}^{LN}\xi_{4}$ do not correlate with the Mps activity. The detailed results of the forward stepwise analysis are given in **Table 3**.

Table 3 comes about here

Analyzing the above equations, it is important to highlight that, the combination of a negative contribution of ${}^{LN}\xi_1$ and a positive contribution of ${}^{LN}\xi_5$ in **Eq. 21** points to a pseudo-folding rule for the biological activity. A validation procedure was subsequently performed in order to assess the model predictability. This validation was carried out with an external series of Mps and randomized control sequences (Cgs). The present model showed accuracy of 90.87%, which is similar in comparison to results obtained by other researchers on using the LDA method in QSAR studies (González-Díaz et al., 2007b). These results are also consistent with many others we have recently reviewed in-depth and published in the form of review article where we used different network-like indices in small-sized, nucleic acid, and protein QSAR (González-Díaz et al., 2007b; González-Díaz et al., 2005d; González-Díaz et al., 2005d; Marrero-Ponce et al., 2005b; Van Waterbeemd, 1995).

3.2. Results for DNA SG indices

The second study used the SG-QSAR models in order to evaluate the same mycobacterial DNA promoter property (see **Table 3**). The grouping of the embedded and non-embedded TIs was done similar to the lattice models: the traces (${}^{SG}\pi_{k}{}^{e} / {}^{SG}\pi_{k}$), the Shannon entropies (${}^{SG}\theta_{k}{}^{e} / {}^{SG}\theta_{k}$), the rest of embedded and non-embedded TIs (H, W, S₆, S, J, ${}^{0}X$, ${}^{2-5}X$, ${}^{1}XR$, H^e, W^e, S₆^e, S^e, J^e, ${}^{0}X^{e}$, ${}^{2-5}X^{e}$, ${}^{1}XR^{e}$) and all SG TIs (pool). The *Forward Stepwise* selection variable method, conjugated with the nE & E TIs of the virtually translated DNA sequences, provides better results for the codon grouping of

the nucleotides, with accuracy, sensitivity and specificity greater than 70% for the ${}^{SG}\pi_k^e$ / ${}^{SG}\pi_k$ and for the pool (**Table 3**). Even if the accuracy of the simple nucleotide sequences are up to 81.58% (pool), the selectivity and the specificity have values less than 70%. The best QSAR model using the SG based on the codon sequences is defined with the ${}^{SG}\pi_k^e$ / ${}^{SG}\pi_k$ group of indices in **Eq. 22** and is characterized by 74.77% accuracy, 82.96% sensitivity and 72.60% specificity.

^{SG}Mps-score=
$$-1.9+1.3 \times^{SG} \pi_4 - 1.9 \times^{SG} \pi_4^{e} - 1.2 \times^{SG} \pi_5^{e}$$
 (22)

Despite the good values of accuracy, sensitivity and specificity (80.80%, 74.81%, 82.39%) for the pool group of TIs (${}^{SG}\theta_0$, ${}^{SG}\theta_4{}^e$, ${}^{SG}\pi_4{}^e$, ${}^{SG}\pi_5{}^e$, W), the QSAR model cannot be considered due to the low sensitivity for the CV set (66.67%). Thus, the results based on the traces (spectral moments) are similar in the case of LN and SG representations, maintaining the ${}^{SG}\pi_5{}^e$ / ${}^{LN}\pi_5{}$ in the equations.

3.3. Comparison with RNA 2S and other indices

In previous works, we published QSAR models to predict Mps using RNA electrostaticdriven 2S folding representations. These models were based on the ^{2S} θ_k (González-Díaz et al., 2007c), ^{2S} π_k (González-Díaz et al., 2005a) and ^{2S} ξ_k (González-Díaz et al., 2006a) values for the ¹**II**(2S) matrix associated to RNA 2S folding representations. In **Table 3** we illustrate that the best values of accuracy, sensitivity and specificity of 97.60%, 93.30% and 100% were found for ^{2S} θ_0 . This TI is present in the QSAR equations for DNA LN/SG and RNA 2S folding representations. All these observations pointed out the importance of the spectral moments, entropies and in the stochastic electrostatic potentials in the DNA/RNA QSAR models. In general, the results for RNA 2S folding representation are better, but require additional calculations for optimization of the RNA 2S. Therefore, more RNA 2S are possible for the same DNA sequence (theoretically

because the promoters are have not correspondent RNA) introducing an indeterminacy in the final model prediction. In the **Figure 3** we depicts a possible 2S for the RNA sequence corresponding to the DNA sequences used in **Figures 1** and **2** (dG is the free energy). This RNA 2S was obtained with the online DINAMelt server (Markham and Zuker, 2005). The SG TIs that show to not be important for the DNA/RNA models (H, W, S, J) can successfully describe protein QSAR models (Munteanu et al., 2008b). This work pointed out the conclusion that the models based on SG, LN and also 2S, which are linear and have few variables, compares very favourably in terms of complexity with other models previously reported by Kalate et al. - these authors used a non-linear artificial neural network and a large parameter space (Kalate et al., 2003).

Figure 3 comes about here

5. Conclusions

The work presents a comparative study of the parameters associated with LN and SG representations in order to predict the mycobacterial DNA promoters. LN QSAR classifier successfully discriminates between Mps and a control group, with values significatively better than the SG-QSAR results based on the DNA codon sequences. In addition, the DNA nucleotide sequences (used for LN) were not able to create a good model based on SG representations. The work promotes the use of the experience accumulated in small-molecules QSAR with spectral moments and other kind of indices (entropies and spectral moments) in new types of DNA QSAR studies, now in the focus of interest many researchers worldwide.

Acknowledgments

Cristian R. Munteanu thanks the FCT (Portugal) for support from grant SFRH/BPD/24997/2005. González-Díaz H. acknowledges program Isidro Parga Pondal

16

of Xunta de Galicia by financial support of a tenure-eligible research position at the Faculty of pharmacy, University of Santiago de Compostela (Spain). The authors thank financial support from grant INCITE07PXI203141ES, Conselleria de Industria, Xunta de Galicia, Spain.

Accepted manuscript

References

- Aguero-Chapin, G., Gonzalez-Diaz, H., Molina, R., Varona-Santos, J., Uriarte, E., and Gonzalez-Diaz, Y., 2006. Novel 2D maps and coupling numbers for protein sequences. The first QSAR study of polygalacturonases; isolation and prediction of a novel sequence from Psidium guajava L. FEBS Lett 580, 723-30.
- Althaus, I.W., Chou, J.J., Gonzales, A.J., Diebel, M.R., Chou, K.C., Kezdy, F.J., Romero, D.L., Aristoff, P.A., Tarpley, W.G., and Reusser, F., 1993a. Steadystate kinetic studies with the non-nucleoside HIV-1 reverse transcriptase inhibitor U-87201E. J Biol Chem 268, 6119-6124.
- Althaus, I.W., Gonzales, A.J., Chou, J.J., Diebel, M.R., Chou, K.C., Kezdy, F.J., Romero, D.L., Aristoff, P.A., Tarpley, W.G., and Reusser, F., 1993b. The quinoline U-78036 is a potent inhibitor of HIV-1 reverse transcriptase. J Biol Chem 268, 14875-14880.
- Althaus, I.W., Chou, J.J., Gonzales, A.J., Diebel, M.R., Chou, K.C., Kezdy, F.J., Romero, D.L., Aristoff, P.A., Tarpley, W.G., and Reusser, F., 1993c. Kinetic studies with the nonnucleoside HIV-1 reverse transcriptase inhibitor U-88204E. Biochemistry 32, 6548-6554.
- Andraos, J., 2008. Kinetic plasticity and the determination of product ratios for kinetic schemes leading to multiple products without rate laws: new methods based on directed graphs. Canadian Journal of Chemistry 86, 342-57.
- Arnvig, K.B., Gopal, B., Papavinasasundaram, K.G., Cox, R.A., and Colston, M.J., 2005. The mechanism of upstream activation in the rrnB operon of Mycobacterium smegmatis is different from the Escherichia coli paradigm. Microbiology 151, 467-73.
- Caballero, J., Fernandez, L., Abreu, J.I., and Fernandez, M., 2006. Amino Acid Sequence Autocorrelation vectors and ensembles of Bayesian-Regularized Genetic Neural Networks for prediction of conformational stability of human lysozyme mutants. J Chem Inf Model 46, 1255-68.
- Chen, Y.L., and Li, Q.Z., 2007. Prediction of apoptosis protein subcellular location using improved hybrid approach and pseudo amino acid composition. J Theor Biol 248, 377-81.
- Chiti, F., Stefani, M., Taddei, N., Ramponi, G., and Dobson, C.M., 2003. Rationalization of the effects of mutations on peptide and protein aggregation rates. Nature 424, 805-8.
- Chou, K.C., 1981. Two new schematic rules for rate laws of enzyme-catalyzed reactions. Journal of Theoretical Biology 89, 581-92.
- Chou, K.C., 1989. Graphical rules in steady and non-steady enzyme kinetics. J Biol Chem 264, 12074-79.
- Chou, K.C., 1990. Review: Applications of graph theory to enzyme kinetics and protein folding kinetics. Steady and non-steady state systems. Biophysical Chemistry 35, 1-24.
- Chou, K.C., 2001a. Prediction of protein cellular attributes using pseudo-amino acid composition. Proteins 43, 246-55.

- Chou, K.C., 2001b. Prediction of signal peptides using scaled window. Peptides 22, 1973-9.
- Chou, K.C., 2004. Review: Structural bioinformatics and its impact to biomedical science. Curr Med Chem 11, 2105-34.
- Chou, K.C., 2005. Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. Bioinformatics 21, 10-9.
- Chou, K.C., Jiang, S.P., Liu, W.M., and Fee, C.H., 1979. Graph theory of enzyme kinetics: 1. Steady-state reaction system. Scientia Sinica 22, 341-58.
- Chou, K.C., and Forsen, S., 1980. Graphical rules for enzyme-catalyzed rate laws. Biochemical Journal 187, 829-35.
- Chou, K.C., and Liu, W.M., 1981. Graphical rules for non-steady state enzyme kinetics. Journal of Theoretical Biology 91, 637-54.
- Chou, K.C., and Zhang, C.T., 1992. Diagrammatization of codon usage in 339 HIV proteins and its biological implication. AIDS Research and Human Retroviruses 8, 1967-76.
- Chou, K.C., Kezdy, F.J., and Reusser, F., 1994. Review: Steady-state inhibition kinetics of processive nucleic acid polymerases and nucleases. Analytical Biochemistry 221, 217-30.
- Chou, K.C., and Zhang, C.T., 1995. Prediction of protein structural classes. Crit Rev Biochem Mol Biol 30, 275-349.
- Chou, K.C., and Shen, H.B., 2006. Hum-PLoc: a novel ensemble classifier for predicting human protein subcellular localization. Biochem Biophys Res Commun 347, 150-7.
- Chou, K.C., and Shen, H.B., 2007a. Recent progress in protein subcellular location prediction. Anal Biochem 370, 1-16.
- Chou, K.C., and Shen, H.B., 2007b. Signal-CF: a subsite-coupled and window-fusing approach for predicting signal peptides. Biochem Biophys Res Commun 357, 633-40.
- Chou, K.C., and Shen, H.B., 2008. Cell-PLoc: A package of web-servers for predicting subcellular localization of proteins in various organisms. Nature Protocols 3, 153-62.
- Cornish-Bowden, A., 1979. Fundamentals of Enzyme Kinetics, Chapter 4. Butterworths, London.
- Di Francesco, V.M., P. J.; Garnier, J, 1999. FORESST: fold recognition from secondary structure predictions of proteins. Bioinformatics 15, 131-40.
- Diao, Y., Li, M., Feng, Z., Yin, J., and Pan, Y., 2007. The community structure of human cellular signaling network. J Theor Biol 247, 608-15.
- Ding, Y.S., Zhang, T.L., and Chou, K.C., 2007. Prediction of protein structure classes with pseudo amino acid composition and fuzzy support vector machine network. Protein Pept Lett 14, 811-5.
- Dobson, P.D., and Doig, A.J., 2005. Predicting enzyme class from protein structure without alignments. J Mol Biol 345, 187-99.

- Dobson, P.M., Boyle, M., and Loewenthal, M., 2004. Home intravenous antibiotic therapy and allergic drug reactions: is there a case for routine supply of anaphylaxis kits? J Infus Nurs 27, 425-30.
- Dobson, P.S., Weaver, J.M., Holder, M.N., Unwin, P.R., and Macpherson, J.V., 2005. Characterization of batch-microfabricated scanning electrochemical-atomic force microscopy probes. Anal Chem 77, 424-34.
- Du, Q.S., Wei, Y.T., Pang, Z.W., Chou, K.C., and Huang, R.B., 2007a. Predicting the affinity of epitope-peptides with class I MHC molecule HLA-A*0201: an application of amino acid-based peptide prediction. Protein Eng Des Sel 20, 417-23.
- Du, Q.S., Huang, R.B., Wei, Y.T., Wang, C.H., and Chou, K.C., 2007b. Peptide reagent design based on physical and chemical properties of amino acid residues. J Comput Chem 28, 2043-50.
- Du, Q.S., Huang, R.B., Wei, Y.T., Du, L.Q., and Chou, K.C., 2008. Multiple field three dimensional quantitative structure-activity relationship (MF-3D-QSAR). J Comput Chem 29, 211-9.
- Estrada, E., 2000. On the topological sub-structural molecular design (TOSS-MODE) in QSPR/QSAR and drug design research. SAR QSAR Environ Res 11, 55-73.
- Estrada, E., 2002. Characterization of the folding degree of proteins. Bioinformatics 18, 697-704.
- Estrada, E., and González-Díaz, H., 2003. What are the limits of applicability for graph theoretic descriptors in QSPR/QSAR? Modeling dipole moments of aromatic compounds with TOPS-MODE descriptors. J Chem Inf Comput Sci 43, 75-84.
- Gao, L., Ding, Y.S., Dai, H., Shao, S.H., Huang, Z.D., and Chou, K.C., 2006. A novel fingerprint map for detecting SARS-CoV. Journal of Pharmaceutical and Biomedical Analysis 41, 246-250.
- Garcia-Garcia, A., Galvez, J., de Julian-Ortiz, J.V., Garcia-Domenech, R., Munoz, C., Guna, R., and Borras, R., 2004. New agents active against Mycobacterium avium complex selected by molecular topology: a virtual screening method. J Antimicrob Chemother 53, 65-73.
- González-Díaz, H., and Uriarte, E., 2005. Biopolymer stochastic moments. I. Modeling human rhinovirus cellular recognition with protein surface electrostatic moments. Biopolymers 77, 296-303.
- González-Díaz, H., de Armas, R.R., and Molina, R., 2003. Markovian negentropies in bioinformatics. 1. A picture of footprints after the interaction of the HIV-1 Psi-RNA packaging region with drugs. Bioinformatics 19, 2079-87.
- González-Díaz, H., Pérez-Bello, A., and Uriarte, E., 2005a. Stochastic molecular descriptors for polymers. 3. Markov electrostatic moments as polymer 2Dfolding descriptors: RNA–QSAR for mycobacterial promoters. Polymer 46 6461–6473.
- González-Díaz, H., Cruz-Monteagudo, M., Molina, R., Tenorio, E., and Uriarte, E., 2005b. Predicting multiple drugs side effects with a general drug-target interaction thermodynamic Markov model. Bioorg Med Chem 13, 1119-29.

- González-Díaz, H., Aguero-Chapin, G., Varona-Santos, J., Molina, R., de la Riva, G., and Uriarte, E., 2005c. 2D RNA-QSAR: assigning ACC oxidase family membership with stochastic molecular descriptors; isolation and prediction of a sequence from Psidium guajava L. Bioorg Med Chem Lett 15, 2932-7.
- González-Díaz, H., Cruz-Monteagudo, M., Vina, D., Santana, L., Uriarte, E., and De Clercq, E., 2005d. QSAR for anti-RNA-virus activity, synthesis, and assay of anti-RSV carbonucleosides given a unified representation of spectral moments, quadratic, and topologic indices. Bioorg Med Chem Lett 15, 1651-7.
- González-Díaz, H., Aguero, G., Cabrera, M.A., Molina, R., Santana, L., Uriarte, E., Delogu, G., and Castanedo, N., 2005e. Unified Markov thermodynamics based on stochastic forms to classify drugs considering molecular structure, partition system, and biological species: distribution of the antimicrobial G1 on rat tissues. Bioorg Med Chem Lett 15, 551-7.
- González-Díaz, H., Perez-Bello, A., Uriarte, E., and Gonzalez-Diaz, Y., 2006a. QSAR study for mycobacterial promoters with low sequence homology. Bioorg Med Chem Lett 16, 547-53.
- González-Díaz, H., Vina, D., Santana, L., de Clercq, E., and Uriarte, E., 2006b. Stochastic entropy QSAR for the in silico discovery of anticancer compounds: prediction, synthesis, and in vitro assay of new purine carbanucleosides. Bioorg Med Chem 14, 1095-107.
- Gonzalez-Diaz, H., Sanchez-Gonzalez, A., and Gonzalez-Diaz, Y., 2006c. 3D-QSAR study for DNA cleavage proteins with a potential anti-tumor ATCUN-like motif. J Inorg Biochem 100, 1290-7
- González-Díaz, H., Molina-Ruiz, R., and Hernandez, I., <u>MARCH-INSIDE</u> version 3.0 (<u>MARkov CHains INvariants for SImulation & DEsign</u>); Windows supported version under request to the main author contact email: <u>gonzalezdiazh@yahoo.es</u>, 2007a.
- González-Díaz, H., Vilar, S., Santana, L., and Uriarte, E., 2007b. Medicinal Chemistry and Bioinformatics – Current Trends in Drugs Discovery with Networks Topological Indices. Curr Top Med Chem 7, 1025-39.
- González-Díaz, H., Pérez-Bello, A., Cruz-Monteagudo, M., González-Díaz, Y., Santana, L., and Uriarte, E., 2007c. Chemometrics for QSAR with low sequence homology: Mycobacterial promoter sequences recognition with 2D-RNA entropies. Chemom Intell Lab Systs 85, 20-26.
- González-Díaz, H., Agüero-Chapin, G., Varona, J., Molina, R., Delogu, G., Santana, L., Uriarte, E., and Gianni, P., 2007d. 2D-RNA-Coupling Numbers: A New Computational Chemistry Approach to Link Secondary StructureTopology with Biological Function. J Comput Chem 28, 1049–56.
- González-Díaz, H., Molina, R.R., Uriarte, E, 2003a. Stochastic molecular descriptors for polymers. 1. Modelling the properties of icosahedral viruses with 3D-Markovian negentropies. Polymer, 3845-53.
- González-Díaz, H., de Armas, R.R., and Molina, R, 2003b. Vibrational Markovian modelling of footprints after the interaction of antibiotics with the packaging region of HIV type 1. Bull. Math. Biol 65, 991-1002.

- González-Díaz, H., Vilar, S., Santana, L., and Uriarte, E., 2007a. Medicinal chemistry and bioinformatics - current trends in drugs discovery with networks topological indices. Curr. Top. Med. Chem. 10, 1015-29.
- González-Díaz, H., Bonet, I., Teran, C., De Clercq, E., Bello, R., Garcia, M.M., Santana, L., and Uriarte, E., 2007b. ANN-QSAR model for selection of anticancer leads from structurally heterogeneous series of compounds. Eur J Med Chem 42, 580-5.
- González-Díaz, H., González-Díaz, Y., Santana, L., Ubeira, F.M., and Uriarte, E., 2008. Proteomics, networks, and connectivity indices. Proteomics 8, 750-778.
- González, M.P., and Moldes del Carmen Teran, M., 2004. A TOPS-MODE approach to predict adenosine kinase inhibition. Bioorg Med Chem Lett 14, 3077-9.
- González, M.P., Helguera, A.M., and Cabrera, M.A., 2005. Quantitative structureactivity relationship to predict toxicological properties of benzene derivative compounds. Bioorg Med Chem 13, 1775-81.
- González, M.P., Teran, C., and Teijeira, M., 2006. A topological function based on spectral moments for predicting affinity toward A3 adenosine receptors. Bioorg Med Chem Lett 16, 1291-6.
- Griffiths, A.J.F., Miller, J.H., Suzuki, D.T., Lewontin, R.C., and Gelbart, W.M., 1999. Introduction to Genetic Analysis. W. H. Freeman & Co., New York.
- Han, L., Cui, J., Lin, H., Ji, Z., Cao, Z., Li, Y., and Chen, Y., 2006. Recent progresses in the application of machine learning approach for predicting protein functional class independent of sequence similarity. Proteomics 6, 4023-37.
- Harary, F., 1969. Graph Theory, MA.
- Harshey, R.M., and Ramakrishnan, T., 1977. Rate of ribonucleic acid chain growth in Mycobacterium tuberculosis H37Rv. J Bacteriol 129, 616-22.
- Jacchieri, S.G., 2000. Mining combinatorial data in protein sequences and structures. Mol Divers 5, 145-52.
- Kalate, R.N., Tambe, S.S., and Kulkarni, B.D., 2003. Artificial neural networks for prediction of mycobacterial promoter sequences. Comput Biol Chem 27, 555-64.
- King, E.L., and Altman, C., 1956. A schematic method of deriving the rate laws for enzyme-catalyzed reactions. Journal of Physical Chemistry 60, 1375-1378.
- Kutner, M.H., Nachtsheim, C.J., Neter, J., and Li, W., Standardized Multiple Regression Model, Applied Linear Statistical Models, McGraw Hill, New York 2005, pp. 271-277.
- Kuzmic, P., Ng, K.Y., and Heath, T.D., 1992. Mixtures of tight-binding enzyme inhibitors. Kinetic analysis by a recursive rate equation. Anal Biochem 200, 68-73.
- Levine, M., and Tjian, R., 2003. Transcription regulation and animal diversity. Nature 424, 147-51.
- Liao, B., and Wang, T.M., 2004a. Analysis of similarity/dissimilarity of DNA sequences based on nonoverlapping triplets of nucleotide bases. J Chem Inf Comput Sci 44, 1666-70.

- Liao, B., and Wang, T.M., 2004b. New 2D graphical representation of DNA sequences. J Comput Chem 25, 1364-8.
- Liao, B., and Wang, T.M., 2004c. A 3D graphical representation of RNA secondary structures. J Biomol Struct Dyn 21, 827-32.
- Liao, B., and Ding, K., 2005. Graphical approach to analyzing DNA sequences. J Comput Chem 26, 1519-23.
- Liao, B., Ding, K., and Wang, T.M., 2005. On a six-dimensional representation of RNA secondary structures. J Biomol Struct Dyn 22, 455-63.
- Liao, B., Xiang, X., and Zhu, W., 2006. Coronavirus phylogeny based on 2D graphical representation of DNA sequence. J Comput Chem 27, 1196-202.
- Lin, H., 2008. The modified Mahalanobis Discriminant for predicting outer membrane proteins by using Chou's pseudo amino acid composition. J Theor Biol 252, 350-6.
- Liu, Y., Guo, X., Xu, J., Pan, L., and Wang, S., 2002. Some notes on 2-D graphical representation of DNA sequence. J Chem Inf Comput Sci 42, 529-33.
- Markham, N.R., and Zuker, M., 2005 DINAMelt web server for nucleic acid melting prediction. Nucleic Acids Res 33 W577-W581.
- Marrero-Ponce, Y., Diaz, H.G., Zaldivar, V.R., Torrens, F., and Castro, E.A., 2004a. 3D-chiral quadratic indices of the 'molecular pseudograph's atom adjacency matrix' and their application to central chirality codification: classification of ACE inhibitors and prediction of sigma-receptor antagonist activities. Bioorg Med Chem 12, 5331-42.
- Marrero-Ponce, Y., Medina-Marrero, R., Castillo-Garit, J.A., Romero-Zaldivar, V., Torrens, F., and Castro, E.A., 2005a. Protein linear indices of the 'macromolecular pseudograph alpha-carbon atom adjacency matrix' in bioinformatics. Part 1: prediction of protein stability effects of a complete set of alanine substitutions in Arc repressor. Bioorg Med Chem 13, 3003-15.
- Marrero-Ponce, Y., Castillo-Garit, J.A., Olazabal, E., Serrano, H.S., Morales, A., Castanedo, N., Ibarra-Velarde, F., Huesca-Guillen, A., Sanchez, A.M., Torrens, F., and Castro, E.A., 2005b. Atom, atom-type and total molecular linear indices as a promising approach for bioorganic and medicinal chemistry: theoretical and experimental assessment of a novel method for virtual screening and rational design of new lead anthelmintic. Bioorg Med Chem 13, 1005-20.
- Marrero-Ponce, Y., Castillo-Garit, J.A., Olazabal, E., Serrano, H.S., Morales, A., Castanedo, N., Ibarra-Velarde, F., Huesca-Guillen, A., Jorge, E., del Valle, A., Torrens, F., and Castro, E.A., 2004b. TOMOCOMD-CARDD, a novel approach for computer-aided 'rational' drug design: I. Theoretical and experimental assessment of a promising method for computational screening and in silico design of new anthelmintic compounds. J Comput Aided Mol Des 18, 615-34.
- Meneses-Marcel, A., Marrero-Ponce, Y., Machado-Tugores, Y., Montero-Torres, A., Pereira, D.M., Escario, J.A., Nogal-Ruiz, J.J., Ochoa, C., Aran, V.J., Martinez-Fernandez, A.R., and Garcia Sanchez, R.N., 2005. A linear discrimination analysis based virtual screening of trichomonacidal lead-like compounds: outcomes of in silico studies supported by experimental results. Bioorg Med Chem Lett 15, 3838-43.

- Mulder, M.A., Zappe, H., and Steyn, L.M., 1997. Mycobacterial promoters. Tuber Lung Dis 78, 211-23.
- Mulligan, M.E., and McClure, W.R., 1986. Analysis of the occurrence of promoter-sites in DNA. Nucleic Acids Res 14, 109-26.
- Mulligan, M.E., Hawley, D.K., Entriken, R., and McClure, W.R., 1984. Escherichia coli promoter sequences predict in vitro RNA polymerase selectivity. Nucleic Acids Res 12, 789-800.
- Munteanu, C.R., Gonzalez-Diaz, H., and Magalhaes, A.L., 2008a. Enzymes/nonenzymes classification model complexity based on composition, sequence, 3D and topological indices. J Theor Biol 254, 476-82.
- Munteanu, C.R., González-Díaz, H., Borges, F., and Magalhães, A.L., 2008b. Natural/random protein classification models based on star network topological indices. J Theor Biol., http://dx.doi.org/10.1016/j.jtbi.2008.07.018.
- Myers, D., and Palmer, G., 1985. Microcomputer tools for steady-state enzyme kinetics. Bioinformatics (original: Computer Applied Bioscience) 1, 105-10.
- Nandy, A., 1994. Recent investigations into global characteristics of long DNA sequences. Indian J Biochem Biophys 31, 149-55.
- Nandy, A., 1996. Two-dimensional graphical representation of DNA sequences and intron-exon discrimination in intron-rich sequences. Comput Appl Biosci 12, 55-62.
- Nandy, A., and Basak, S.C., 2000. Simple numerical descriptor for quantifying effect of toxic substances on DNA sequences. J Chem Inf Comput Sci 40, 915-9.
- O'Neill, M.C., and Chiafari, F., 1989. Escherichia coli promoters. II. A spacing classdependent promoter search protocol. J Biol Chem 264, 5531-4.
- Pawar, A.P., Dubay, K.F., Zurdo, J., Chiti, F., Vendruscolo, M., and Dobson, C.M., 2005. Prediction of "aggregation-prone" and "aggregation-susceptible" regions in proteins associated with neurodegenerative diseases. J Mol Biol 350, 379-92.
- Prado-Prado, F.J., Gonzalez-Diaz, H., de la Vega, O.M., Ubeira, F.M., and Chou, K.C., 2008. Unified QSAR approach to antimicrobials. Part 3: First multi-tasking QSAR model for Input-Coded prediction, structural back-projection, and complex networks clustering of antiprotozoal compounds. Bioorganic & Medicinal Chemistry 16, 5871-880.
- Qi, X.Q., Wen, J., and Qi, Z.H., 2007. New 3D graphical representation of DNA sequence based on dual nucleotides. Journal of Theroretical Biology 249, 681-90.
- Ramos de Armas, R., González-Díaz, H., Molina, R., Perez Gonzalez, M., and Uriarte, E., 2004. Stochastic-based descriptors studying peptides biological properties: modeling the bitter tasting threshold of dipeptides. Bioorg Med Chem 12, 4815-22.
- Randic, M., 2004. 2-D graphical representation of proteins based on virtual genetic code. SAR QSAR Environ Res 15, 147-57.
- Randic, M., and Vracko, M., 2000. On the similarity of DNA primary sequences. J Chem Inf Comput Sci 40, 599-606.

- Randic, M., and Balaban, A.T., 2003. On a four-dimensional representation of DNA primary sequences. J Chem Inf Comput Sci 43, 532-9.
- Randic, M., and Zupan, J., 2004. Highly compact 2D graphical representation of DNA sequences. SAR QSAR Environ Res 15, 191-205.
- Randic, M., Zupan, J., and Vikic-Topic, D., 2007. On representation of proteins by starlike graphs. J Mol Graph Model, 290-305.
- Randic, M., Vracko, M., Nandy, A., and Basak, S.C., 2000. On 3-D graphical representation of DNA primary sequences and their numerical characterization. J Chem Inf Comput Sci 40, 1235-44.
- Santana, L., Uriarte, E., Gonzalez-Diaz, H., Zagotto, G., Soto-Otero, R., and Mendez-Alvarez, E., 2006. A QSAR model for in silico screening of MAO-A inhibitors. Prediction, synthesis, and biological assay of novel coumarins. J Med Chem 49, 1149-56.
- Song, J., and Tang, H., 2005. A new 2-D graphical representation of DNA sequences and their numerical characterization. J Biochem Biophys Methods 63, 228-39.
- StatSoft.Inc., STATISTICA (data analysis software system), version 6.0, www.statsoft.com.Statsoft, 2002.
- Todeschini, R., and Consonni, V., 2002. Handbook of Molecular Descriptors. Wiley-VCH.
- Van Waterbeemd, H., 1995. Chemometric methods in molecular design. Wiley-VCH, New York.
- Vilar, S., Santana, L., and Uriarte, E., 2006. Probabilistic neural network model for the in silico evaluation of anti-HIV activity and mechanism of action. J Med Chem 49, 1118-124.
- Vilar, S., Estrada, E., Uriarte, E., Santana, L., and Gutierrez, Y., 2005. In silico studies toward the discovery of new anti-HIV nucleoside compounds through the use of TOPS-MODE and 2D/3D connectivity indices. 2. Purine derivatives. J Chem Inf Model 45, 502-14.
- Wang, M., Yao, J.S., Huang, Z.D., Xu, Z.J., Liu, G.P., Zhao, H.Y., Wang, X.Y., Yang, J., Zhu, Y.S., and Chou, K.C., 2005. A new nucleotide-composition based fingerprint of SARS-CoV with visualization analysis. Medicinal Chemistry 1, 39-47.
- Wolfram, S., 1984. Cellular automation as models of complexity. Nature 311, 419-24.
- Wolfram, S., 2002. A New Kind of Science. Wolfram Media Inc., Champaign, IL.
- Woodcock, S., Mornon, J.P., and Henrissat, B., 1992. Detection of secondary structure elements in proteins by hydrophobic cluster analysis. Protein Eng 5, 629-35.
- Wyrick, J.J., and Young, R.A., 2002. Deciphering gene expression regulatory networks. Curr Opin Genet Dev 12, 130-6.
- Xiao, X., and Chou, K.C., 2007. Digital coding of amino acids based on hydrophobic index. Protein Pept Lett 14, 871-5.
- Xiao, X., Shao, S., Ding, Y., Huang, Z., Chen, X., and Chou, K.C., 2005a. Using cellular automata to generate Image representation for biological sequences. Amino Acids 28, 29-35.

- Xiao, X., Shao, S., Ding, Y., Huang, Z., Chen, X., and Chou, K.C., 2005b. An Application of Gene Comparative Image for Predicting the Effect on Replication Ratio by HBV Virus Gene Missense Mutation. Journal of Theoretical Biology 235, 555-65.
- Xiao, X., Shao, S.H., and Chou, K.C., 2006a. A probability cellular automaton model for hepatitis B viral infections. Biochem. Biophys. Res. Comm. 342, 605-10.
- Xiao, X., Shao, S.H., Ding, Y.S., Huang, Z.D., and Chou, K.C., 2006b. Using cellular automata images and pseudo amino acid composition to predict protein subcellular location. Amino Acids 30, 49-54.
- Yuan, Z., 1999. Prediction of protein subcellular locations using Markov chain models. FEBS Lett 451, 23-6.
- Zhang, C.T., and Chou, K.C., 1993. Graphic analysis of codon usage strategy in 1490 human proteins. Journal of Protein Chemistry 12, 329-35.
- Zhang, C.T., and Chou, K.C., 1994. Analysis of codon usage in 1562 E. Coli protein coding sequences. Journal of Molecular Biology 238, 1-8.
- Zhou, G.P., and Deng, M.H., 1984. An extension of Chou's graphical rules for deriving enzyme kinetic equations to system involving parallel reaction pathways. Biochemical Journal 222, 169-76.
- Zupan, J., and Randic, M., 2005. Algorithm for coding DNA sequences into "spectrumlike" and "zigzag" representations. J Chem Inf Model 45, 309-13.

Legend for Figures

- Figure 1. LN for the Mps of the gene Alpha in Mycobacterum bovis (BCG).
- Figure 2. SG for the Mps of the gene Alpha in Mycobacterum bovis (BCG).
- Figure 3. RNA 2S for the Mps of the gene Alpha in Mycobacterum bovis (BCG).

Accepted manuscript

Legend for Tables

- Table 1. LN construction rules for the Mps of the gene Alpha in Mycobacterum bovis

 (BCG).
- Table 2. SG codifications for the virtually translated Mps of the gene Alpha in

 Mycobacterum bovis (BCG).
- Table 3. Summary of the LDA results for DNA LN and SG models vs. RNA 2S folding representations.

Accepted manuscrip

Table 1:

DNA Lattice Network

 $c_{1}g_{2}a_{3}c_{4}t_{5}t_{6}t_{7}c_{8}g_{9}c_{10}c_{11}c_{12}g_{13}a_{14}a_{15}t_{16}c_{17}g_{18}a_{19}c_{20}\\ a_{21}t_{22}t_{23}t_{24}g_{25}g_{26}c_{27}c_{28}t_{29}c_{30}c_{31}a_{32}c_{33}a_{34}c_{35}a_{36}c_{37}g_{38}g_{39}t_{40}\\ a_{41}t_{42}g_{43}t_{44}t_{45}c_{46}t_{47}g_{48}g_{49}c_{50}c_{51}c_{52}g_{53}a_{54}g_{55}c_{56}a_{57}c_{58}a_{59}c_{60}\\ g_{61}a_{62}c_{63}g_{64}a_{65} \\ \end{array}$

n	Nucleotide	X	у
1	$c_1 a_3 t_5 g_{25}$	0	0
2	$g_2 c_{10} g_{26}$	0	-1
3	$c_4 t_{16}$	-1	0
4	t_6c_8	1	0
5	t_7	2	0
6	g 9	1	-1
7	$c_{11}c_{27}t_{29}$	-1	-1
8	$c_{12}a_{14}g_{18}c_{28}c_{30}g_{48}$	-2	-1
9	g 13 g 49	-2	-2
10	$a_{15}c_{17}a_{19}t_{45}t_{47}$	-2	-0
11	$c_{20}a_{32}t_{44}c_{46}$	-3	0
12	a ₂₁	-3	1
13	t ₂₂	-2	1
14	t ₂₃	-1	1
15	t ₂₄	0	1
16	c ₃₁	-3	-1
17	c ₃₃ g ₄₃	-4	0
18	a34t42	-4	1
19	c 35 a 41	-5	1
20	a ₃₆	-5	2
21	c ₃₇	-6	2
22	g ₃₈	-6	1
23	g ₃₉	-6	0
24	t_{40}	-5	0
25	c ₅₀	-3	-2
26	c ₅₁	-4	-2
27	$c_{52}a_{54}$	-5	-2
28	g 53 g 55	-5	-3
29	C ₅₆	-6	-3
30	a ₅₇	-6	-2
31	C ₅₈	-7	-2
32	a 59	-7	-1
33	$c_{60}a_{62}$	-8	-1
34	g ₆₁	-8	-2
35	c ₆₃ a ₆₅	-9	-1
36	g ₆₄	-9	-2

$\begin{array}{c} {}_{14}a_{15}t_{16}c_{17}g_{18}a_{19}c_{20}a_{21}\\ {}_{2}c_{33}a_{34}c_{35}a_{36}c_{37}g_{38}g_{39}\\ {}_{3}a_{54}g_{55}c_{56}a_{57}c_{58}a_{59}c_{60}g_{61}a_{62}c_{63}\\ {}_{a_{5}tc}g_{6}aca_{7}\\ {}_{aca_{12}c}gg_{13}\\ {}_{gca_{19}cac_{20}gac_{21}}\end{array}$
$_{2}c_{33}a_{34}c_{35}a_{36}c_{37}g_{38}g_{39}$ $_{3}a_{54}g_{55}c_{56}a_{57}c_{58}a_{59}c_{60}g_{61}a_{62}c_{63}$ $a_{5}tcg_{6}aca_{7}$ $aca_{12}cgg_{13}$ $gca_{19}cac_{20}gac_{21}$
$_{3}a_{54}g_{55}c_{56}a_{57}c_{58}a_{59}c_{60}g_{61}a_{62}c_{63}$ $a_{5}tcg_{6}aca_{7}$ $aca_{12}cgg_{13}$ $gca_{19}cac_{20}gac_{21}$
$a_5 tcg_6 aca_7$ $aca_{12} cgg_{13}$ $gca_{19} cac_{20} gac_{21}$
$aca_{12}cgg_{13}$ $gca_{19}cac_{20}gac_{21}$
gca ₁₉ cac ₂₀ gac ₂₁
$Y_{14}V_{15}L_{16}A_{17}R_{18}A_{19}H_{20}D_{21}$

Table	e 3:
-------	------

ΤI	Ac (%)	Se (%)	Sp (%)	Final TIs	Vars.	λ	F	р	Ref.			
Primary structure of DNA nucleotide & LN												
${}^{LN}\pmb{\theta}_k$	78.33	72.59	79.84	$^{LN}\theta_0$	1	0.74	230.5	0.0001	a			
${}^{LN} \pi_k$	81.73	78.52	82.58	$^{\mathrm{LN}}\pi_{0}, {}^{\mathrm{LN}}\pi_{1}, {}^{\mathrm{LN}}\pi_{5},$	3	0.89	76.3	0.0001	a			
$^{LN}\xi_k$	90.87	82.96	92.95	$^{LN}\xi_1, ^{LN}\xi_5$	2	0.82	142.1	0.0001	a			
Pool	92.88	75.56	97.46	$^{LN}\theta_0$, $^{LN}\pi_0$, $^{LN}\xi_1$, $^{LN}\xi_5$	4	0.83	130.8	0.0001	a			
Primary structure of DNA nucleotide sequences & SG												
${}^{SG}\boldsymbol{\theta}_k$	66.25	81.48	62.23	${}^{SG}\theta_1{}^e, {}^{SG}\theta_4{}^e$	2	0.78	69.62	0.001	a			
${}^{SG} \pi_k$	71.21	85.19	67.51	${}^{\mathrm{SG}}\pi_0^{\mathrm{e}}, {}^{\mathrm{SG}}\pi_2^{\mathrm{e}}, {}^{\mathrm{SG}}\pi_5^{\mathrm{e}}$	3	0.76	49.54	0.001	a			
$\mathbf{T}\mathbf{I}_k$	75.39	68.15	77.30	W, J ^e , ⁰ X ^e	3	0.73	58.19	0.001	a			
Pool	81.58	68.15	85.13	${}^{SG}\pi_5^{e}$, H, ${}^{1}XR^{e}$	3	0.67	79.94	0.001	a			
Primary structure of DNA <i>codon</i> sequences & SG												
${}^{SG}\boldsymbol{\theta}_k$	70.43	76.30	68.88	${}^{SG}\theta_0, {}^{SG}\theta_1, {}^{SG}\theta_4{}^e$	3	0.75	52.31	0.001	a			
${}^{SG} \pi_k$	74.77	82.96	72.60	${}^{\mathrm{SG}}\pi_4, {}^{\mathrm{SG}}\pi_4^{\mathrm{e}}, {}^{\mathrm{SG}}\pi_5^{\mathrm{e}}$	3	0.74	56.37	0.001	a			
$\mathbf{T}\mathbf{I}_k$	76.16	59.26	80.63	S, ⁰ X, ¹ XR ^e	3	0.72	60.98	0.001	a			
Pool	80.80	74.81	82.39	${}^{SG}\theta_0, {}^{SG}\theta_4^{e}, {}^{SG}\pi_4^{e}, {}^{SG}\pi_5^{e}, W$	5	0.67	47.04	0.001	a			
RNA electrostatic-driven 2S folding												
${}^{2S}\theta_k$	97.60	93.30	100.00	$^{28}\theta_0$	1	0.34	724.47	0.001	b			
${}^{2S}\boldsymbol{\pi}_k$	93.83	83.70	98.89	$^{28}\pi_0, {}^{28}\pi_2$	2	0.44	515.03	0.05	c			
${}^{2S}\xi_k$	96.58	85.19	100.00	$28_{\xi_0}, 28_{\xi_1}$	2	0.41	38.8	0.001	d			
JR				10 00								

Note: the terms Ac, Se, and Sp mean accuracy, sensitivity and specificity, and measure the ratio of the number of total, Mps, or Cgs sequences correctly classified by the model with respect to the real classification; Vars. = no of variables in the QSAR equations; SG = star-graph; LN = lattice network; 2S = secondary structure; super index "e" represents the embedded calculations; References (Ref.) are a: this work, b: (González-Díaz et al., 2007c), c: (González-Díaz et al., 2005a) and d: (González-Díaz et al., Accel

2006a).

Figure 1:







Figure 3:

Output of sir_graph (®) by D. Stewart and M. Zuker

 \overline{v}



dG = -10.5 A