

A probabilistic approach to the visual exploration of G Protein-Coupled Receptor sequences

Alfredo Vellido¹, Martha Ivón Cárdenas¹, Iván Olier², Xavier Rovira³
and Jesús Giraldo³ *

1- Departament de Llenguatges i Sistemes Informàtics
Universitat Politècnica de Catalunya, 08034, Barcelona - Spain

2- School of Psychological Sciences
The University of Manchester, M13 9PL, Manchester - United Kingdom

3- Institut de Neurociències, Unitat de Bioestadística
Universitat Autònoma de Barcelona, 08193, Bellaterra (Barcelona) - Spain

Abstract. The study of G protein-coupled receptors (GPCRs) is of great interest in pharmaceutical research, but only a few of their 3D structures are known at present. On the contrary, their amino acid sequences are known and accessible. Sequence analysis can provide new insight on GPCR function. Here, we use a kernel-based statistical machine learning model for the visual exploration of GPCR functional groups from their sequences. This is based on the rich information provided by the model regarding the probability of each sequence belonging to a certain receptor group.

1 Introduction

The study of G protein-coupled receptors (GPCRs) is of great interest in pharmaceutical research. These receptors regulate the function of most cells in living organisms and it is estimated that they are targets for about one third of clinically used drugs.

The function of the proteins depends directly on their 3D structure, which is embodied in their amino acid sequence. GPCRs are membrane proteins, and this environment makes their 3D structure difficult to unravel through nuclear magnetic resonance or X-ray crystallography. Modern molecular biology methods, though, make their sequences easy to acquire. The grouping of GPCRs into classes and subclasses based on sequence analysis may significantly contribute to helping drug design and to a better understanding of the molecular processes involved in receptor signaling both in normal and pathological conditions [1].

In order to group GPCR sequences, we need a measure of similarity between them. A GPCR-specific kernel was recently defined in [2] to this purpose, as part of a kernel-based statistical machine learning model of the manifold learning family, namely the Kernel Generative Topographic Mapping (KGTm). This model describes multivariate data in terms of low dimensional representations, so as to achieve the visualization of high dimensional data that would otherwise be

*This research was partially supported by Catalan La Marató de TV3 Foundation project 070530 and Spanish MICINN projects TIN2009-13895-C02-01 and SAF2007-65913.

difficult to visualize. The visualization of the high-dimensional GPCR sequences would considerably help understanding their global grouping structure.

Recent research [2] using KGTM provided preliminary results of the existence of this structure, including GPCR subclass-specific groups and some level of subclass mixing. Here, we use the probabilistic properties of KGTM to explore these subclasses in more detail. For that we resort to the explicit calculation of the probability of each of the available sequences belonging to each of the model groupings. This provides us with a map of probability that can qualify the differences between sequences of either clear or dubious subclass ascription.

2 Kernel GTM

The GTM is a nonlinear statistical machine learning model of the manifold learning family. It performs simultaneous clustering (as a constrained mixture of distributions model) and low-dimensional visualization of multivariate data. It is defined as a nonlinear mapping from a latent space in \mathfrak{R}^ℓ (with ℓ being usually 1 or 2 for visualization purposes) onto a manifold embedded in the data \mathfrak{R}^D space. This is expressed as a generalized regression function: $\mathbf{y} = \mathbf{W}\phi(\mathbf{u})$, where $\mathbf{y} \in \mathfrak{R}^D$, $\mathbf{u} \in \mathfrak{R}^\ell$, \mathbf{W} is an adaptive matrix of weights, and ϕ is a vector with the images of S basis functions ϕ_s . The prior distribution of \mathbf{u} in latent space is constrained to form a uniform discrete grid of M centres, in the form of a sum of delta functions $p(\mathbf{u}) = \frac{1}{M} \sum_{m=1}^M \delta\mathbf{u} - \mathbf{u}_m$. Each component m in the mixture defines the probability of an observable data point \mathbf{x} given a latent point \mathbf{u}_m and the model parameters:

$$p(\mathbf{x}|\mathbf{u}_m, \Theta) = \left(\frac{\beta}{2\pi}\right)^{D/2} \exp\left\{-\frac{\beta}{2}\|\mathbf{x} - \mathbf{y}_m\|^2\right\} \quad (1)$$

where $\mathbf{y}_m = \mathbf{W}\phi(\mathbf{u}_m)$ and the adaptive parameters Θ are \mathbf{W} and the common inverse variance β . With these probabilities, a density model in data space can be generated for each component m of the mixture, leading to the definition of a complete model likelihood. The adaptive parameters of the model can be optimized by Maximum Likelihood (ML) using the Expectation-Maximization (EM) algorithm. Details can be found in [3].

Kernelization is a method originally defined for Support Vector Machines (SVM). In recent years it has been extended to other models, including those functionally similar to GTM [4]. The idea is that a method formulated in terms of kernels can use the one that best suits the problem and data type at hand. GTM was originally defined for quantitative data in the real domain. The type of data analyzed in this study though, which could be considered as a *text-like* sequence of symbols, should benefit from a kernel formulation of the model.

Observed data \mathbf{X} can be implicitly mapped into a high-dimensional feature space H via a nonlinear function: $\mathbf{x} \rightarrow \psi(\mathbf{x})$. A similarity measure can then be defined from the dot product in space H as follows:

$$K(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle \quad (2)$$

K is a kernel function that should satisfy Mercer's condition [5]. Data are expressed in the high-dimensional dot product space H , usually known as feature space. This use of the feature space reduces the computational cost by employing the kernel function K instead of directly computing the dot product in H .

The kernelization of GTM entails the redefinition of Eq.1 in feature space as:

$$p(\boldsymbol{\psi}(\mathbf{x})|\mathbf{u}_m, \boldsymbol{\Theta}) = \left(\frac{\beta}{2\pi}\right)^{D/2} \exp\left\{-\frac{\beta}{2}\|\boldsymbol{\psi}(\mathbf{x}) - \mathbf{y}_m\|^2\right\} \quad (3)$$

Note that the prototypes \mathbf{y}_m are now defined in the feature space and not in data space, as originally. The expression $\|\boldsymbol{\psi}(\mathbf{x}) - \mathbf{y}_m\|^2$ can be reformulated in terms of kernel functions by expanding the prototypes on the data in feature space. That is $\|\boldsymbol{\psi}(\mathbf{x}) - \mathbf{y}_m\|^2 = K_{nn} + (\boldsymbol{\Lambda}\boldsymbol{\phi}_m)^T \mathbf{K}\boldsymbol{\Lambda}\boldsymbol{\phi}_m - 2\mathbf{k}_n\boldsymbol{\Lambda}\boldsymbol{\phi}_m$, where \mathbf{K} is a kernel matrix with elements $K_{nn'} = \langle \boldsymbol{\psi}(\mathbf{x}_n), \boldsymbol{\psi}(\mathbf{x}_{n'}) \rangle$, and row vectors \mathbf{k}_n . The adaptive parameters of the model are now $\boldsymbol{\Lambda}$ (an adaptive weight matrix) and β , which can again be optimized by ML using EM (see details in [2]). Here we are specially interested in one of the results of the expectation step of EM, namely the estimation of the posterior $R_{mn} = p(\mathbf{u}_m|\boldsymbol{\psi}(\mathbf{x}_n), \boldsymbol{\Lambda}, \beta)$ as:

$$R_{mn} = \frac{p(\boldsymbol{\psi}(\mathbf{x}_n)|\mathbf{u}_m, \boldsymbol{\Lambda}, \beta)}{\sum_{m'=1}^M p(\boldsymbol{\psi}(\mathbf{x}_n)|\mathbf{u}_{m'}, \boldsymbol{\Lambda}, \beta)}. \quad (4)$$

R_{mn} measures the degree of responsibility (probability) of a point \mathbf{u}_m in the latent space for the generation of a $\boldsymbol{\psi}(\mathbf{x}_n)$ GPCR data subsequence. Each R_{mn} is an element of a $M \times N$ responsibility matrix \mathbf{R} .

3 Experiments

GPCRs are traditionally divided into five main classes (*rhodopsin-like* (class A), *secretin-like* (B), *glutamate-like* (C), *adhesion*, and *Frizzled/Taste2*) and, in turn, into a complex branched sub-structure. Seven subclasses of class C are modeled and visualized in this paper using KGTM, namely 1: *Metabotropic glutamate*, 2: *Calcium sensing*, 3: *GABA-B*, 4: *Vomer nasal*, 5: *Pheromone*, 6: *Odorant*, 7: *Taste*. The dataset consists of 232 protein sequences obtained from GPCRDB¹. Each position in a sequence is called a *residue*, which in turn may be one of 20 possible amino acids. Each amino acid has a standard one-letter code, and a sequence is therefore represented by a combination of these letters. The number of residues by sequence in the dataset is 253 (data dimensionality).

The kernel function designed to analyze such data with KGTM is a variation on that in [2], based on the mutations and gaps between sequences:

$$K(x, x') = \exp\left\{\nu \frac{\pi(x, x')}{\sqrt{\pi(x, x)\pi(x', x')}}\right\} \quad (5)$$

where x and x' are two sequences and ν is a prefixed parameter; $\pi(\cdot)$ is a score function commonly used in bioinformatics and expressed as: $\pi(x, x') =$

¹GPCRDB web site: <http://www.gpcr.org/7tm/>

$\sum_r s(x_r, x'_r) - \gamma$, where x_r and x'_r are the r^{th} residue in the sequences. The value of $s(x_r, x'_r)$ can be found in a mutation matrix [6] and γ is a gap penalty (usually the number of gaps in sequences). A normalization factor, defined as the geometric mean of the maximum scores for each of the sequences, is used in 5 instead of their sum, as used in [2]. The modified kernel function now has a proper delimitation of its range.

The first visualization results with KGTM are shown in Fig.1. There is quite clear separation between many of the GPCR subclasses, which are visualized in the latent space using the *mode-projection*, defined as: $m_{mode} = \underset{m}{\operatorname{argmax}} R_{mn}$. Many subclasses occupy a rather differentiated area on the map, showing little overlapping. A few of them, though, have overlapping representations. Both cases could be the source of insight on the peculiarities of subclass structure. *Metabotropic glutamate* (subclass 1), *GABA-B* (3), and *Taste* (7) are clearly differentiated from the rest of subclasses, which show significant overlapping between them.

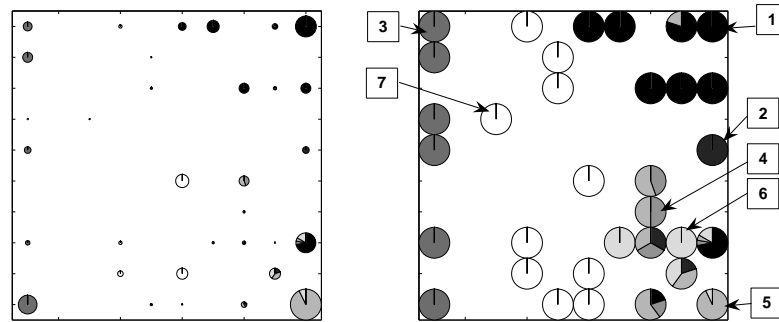


Fig. 1: Data visualization on a 10×10 representation map using the *mode-projection* as described in the text. Left) Pie charts represent latent points, and their size is proportional to the ratio of sequences assigned to them. Each portion of a chart corresponds to the percentage of sequences belonging to each subclass, coded in shades of gray. Right) The same map without sequence ratio size scaling, for better visualization. Labels as described in the text.

The *mode-projection* is an intuitive form of visualization that sacrifices detail in favour of clarity. By using only the maximum of the responsibilities in \mathbf{R} , though, it disposes of much of the rich information that might be contained in this matrix of probabilities.

There are different ways of visually representing this information. One of them is the display of *maps of probability* \mathbf{R}_i , for a given sequence i . Sequences clearly ascribed to a subclass are likely to have their responsibilities concentrated in only a few modes (latent points), whereas the probabilities of sequences with-

out clear subclass ascription may be more evenly spread across the map. Due to space limitations, representations of this level of detail are omitted here.

We may be also interested in the responsibilities of all sequences of a given subclass at once. In this case, we would aim to assess if each subclass has its responsibilities located in a well-defined area of the map or not. The *cumulative responsibility* of the sequences that belong to a given subclass c is defined as a vector $\mathbf{CR}_c = \sum_{\{n \in c\}} (R_{mn})$, for $m = \{1, \dots, M\}$. Figure 2 provides the visualization of the \mathbf{CR}_c for four subclasses, namely those with $c = 1, 3, 5, 6$.

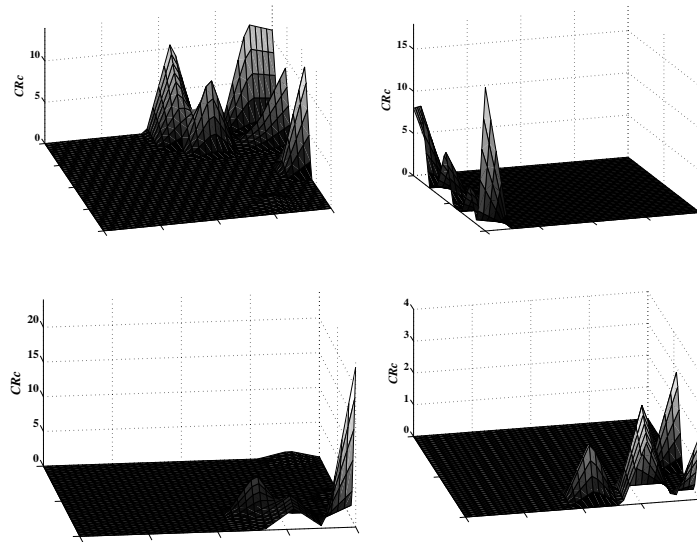


Fig. 2: \mathbf{CR}_c representation maps for 4 GPCR subclasses. Top row) left: subclass 1 (*Metabotropic glutamate*), the most populated, is well-defined on the top-right corner of the map; right: subclass 3 (*GABA-B*), also isolated and unmixed in the left hand-side of the map. Bottom row) left: subclass 5 (*Pheromone*), strongly focused on the bottom right corner of the map, but partially overlapping with right: subclass 6 (*Odorant*). The layout corresponds to that of Fig.1, although with its viewpoint slightly displaced to the left, to provide some perspective.

This takes us to the possibility of displaying the *cumulative responsibility* of all sequences in the database, defined as vector $\mathbf{CR} = \sum_{n=1}^N (R_{mn})$, for $m = \{1, \dots, M\}$. With this map of probability, the existence of \mathbf{CR} *peaks* and *valleys* can be explored. The latter are likely to define the boundaries between subclasses. The global \mathbf{CR} is displayed in Fig.3. Consistent with the subclass-specific representations in Fig.2, several local maxima are shown to correspond to each subclass, which could be an indication of heterogeneity within the subclasses. Some deep *valleys of probability* can be seen in the central parts of the map, drawing clear boundaries between subclasses represented in the periphery of the map and those around its center. Some amongst the latter are the ones

with a higher level of mixing and would merit further investigation.

Our results are consistent with early classification studies using other techniques such as Hidden Markov Models, thereby validating the present methodology. Importantly, the method herein presented reveals mixing between some receptor subclasses, suggesting its possible applicability to the study of heterodimerization between receptors. This finding paves the way for new strategies in drug discovery research. KGTM may help in the exploration of receptors susceptible of heterodimerization and thus be useful in the quest of more potent and safer drugs.

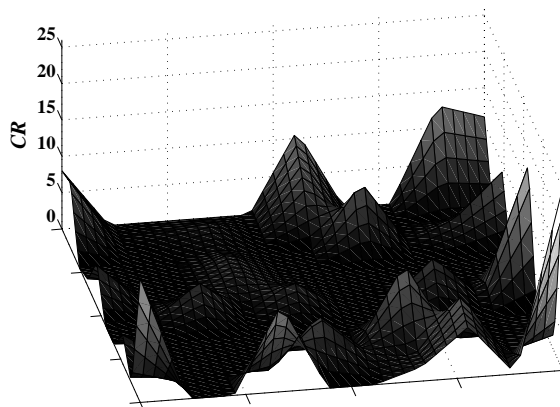


Fig. 3: Visualization of the global **CR** (on the vertical axis) of the data set in the representation map. Layout as in Fig.2.

References

- [1] M. C. Cobanoglu, Y. Saygin and U. Sezerman, Classification of GPCRs using family specific motifs. *IEEE-ACM T. Comput. Bi.*, In press , 2010.
- [2] I. Olier, A. Vellido and J. Giraldo, Kernel Generative Topographic Mapping. In M. Verleysen, editor, *proceedings of the 18th European symposium on artificial neural networks (ESANN 2010)*, 481-486.
- [3] C. M. Bishop, M. Svensén, and C. K. I. Williams, GTM: The Generative Topographic Mapping, *Neural Comput.*, 10(1):215–234, Elsevier, 1998.
- [4] N. Villa and F. Rossi, A comparison between dissimilarity SOM and kernel SOM for clustering the vertices of a graph. In *proceedings of the 6th workshop on self-organizing maps (WSOM 07)*, Bielefeld, Germany, 2007.
- [5] B. Schölkopf and A. Smola. *Learning with Kernels*. The MIT Press, Cambridge, Massachusetts, 2002.
- [6] R. Durbin, S. R. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge Univ. Press, 2004.
- [7] Z. R. Yang and R. Thomson, A novel neural network method in mining molecular sequence data, *IEEE T. Neural Networ.*, 16:263–274, 2005.