Generative Topographic Mapping: universal tool for chemical space analysis

Alexandre Varnek University of Strasbourg

BigChem school, Modena, 25 October 2017



Chemical universe: Big Data problem

- 10⁸ compounds are currently available
- 10³³ drug-like molecules
 could be synthesized *



How to represent this huge chemical space and to navigate in this space ?

* P. Polischuk, T. Madzidov, A. Varnek, J. Comp. Aided Mol. Des, 2013, 27, 675-679

Encoding chemical structures by molecular descriptors

Molecular graph



Descriptors

Constitutional descriptors Ring descriptors Topological indices Walk and path counts Connectivity indices Information indices 2D matrix-based descriptors 2D autocorrelations Burden eigenvalues P_VSA-like descriptors ETA indices Edge adjacency indices Geometrical descriptors 3D matrix-based descriptors 3D autocorrelations

Descriptor vector





> 5000 types of descriptors are used

Chemography: cartography of chemical space

Data visualization => dimensionality reduction problem



Data space (N-dimensional) Latent space (2-dimensional)

Dimensionality reduction requirements







- minimal information loss,
- topology preservation,
- distance preservation

Dimensionality Reduction: information loss



Taxonomy of Dimensionality Reduction Techniques





Dimensionality reduction methods

Acetylcholinesterase dataset (DUD) : 100 actives and 100 inactives



Multi-Dimensional Scaling



Canonical Correlation Analysis



Independent **Component Analysis**



Exploratory Factor Analysis



Sammon map









Isomap



Locally Linear Embedding







t-SNE



Autoencoder dimensionality reduction



SOM

Laplacian Eigenmaps

Generative Topographic Mapping approach





- Nonlinear unsupervised approach
- Simple interpretation
- Topology preservation
- Can be used for classification purposes

Limitations of SOMs

- \succ the absence of a cost function to be optimized in training;
- ➤ the lack of a theoretical basis for choosing methods parameters;
- ➤ the absence of any general proofs of convergence;
- > models do not define a probability density

C.M.Bishop, M.Svensen, C.K.I.Williams, « The Generative Topographic Mapping», *Neural Computation*, 10, No. 1, 215-234 (1998)

Self-Organizing Maps (SOM)







Teuvo Kohonen

Generative Topographic Mapping (GTM)





Christopher Bishop

GTM overcomes most of limitations of SOMs without introducing disadvantages

Generative Topographic Mapping : algorithm



Generative Topographic Mapping (GTM)



GTM generates a data probability distribution in *both initial and latent data spaces*.

This opens an opportunity to use GTM not only to visualize the data but also for structure-property modeling tasks

• C. M. Bishop Pattern Recognition and Machine Learning, 2006 Springer

• N. Kireeva, I.I. Baskin, H. A. Gaspar, D. Horvath, G. Marcou and A. Varnek, Mol. Informatics, 2012, 31, 201-312



Projection of an object on GTM is described by the probability distribution (*responsibilities*) over the lattice nodes.



Probabilities (responsabilities) of acetylcholinesterase ligand projected into latent space

GTM descriptors for molecules and datasets





Map resolution: $N_{nodes} = K^*K$ Standard setting: K = 25, $N_{arid} = 625$

Molecule \longrightarrow responsibilities' vector $\{R_{tk}\}$ of N_{nodes} length

Dataset — Inormalized cumulated responsibilities' vector of **N**_{nodes} length

GTM : areas of application



Chemical data analysis and activity prediction

Properties mapping





political map

physical map







population density

GTM property (activity) landscape



Dataset \longrightarrow property lanscape vector $\{\bar{A}_k\}$ of N_{nodes} length

GTM activity landscape



H. A. Gaspar, I. I. Baskin, G. Marcou, D. Horvath, A. Varnek Mol. Informatics, 2015, 34 (6-7), 348-356



Activity landscape for Lu³⁺ complexation



H. Gaspar, I. Baskin, D. Horvath, G. Marcou, A. Varnek Mol. Informatics, 2015, 34 (6-7), 348-356

Activity landscape: prediction on a test



Responsabilities' distribution of a test set cmpd

GTM for the training set

Predictions on a test set

$$\hat{A}_{j}(\text{test}) = \sum_{k} \overline{A}_{k}(\text{training})R_{jk}(\text{test})$$

Performance of GTM-based QSAR models

Regression models for LogK (Lu³⁺) using ISIDA descriptors



GTM-based models preform similarly to those obtained with popular machine-learning methods

Class landscape for antiviral activity





Class landscape can be used to predict a class ("active"/"inactive") for any new compound

Class landscape for antiviral activity

inactives

actives



Responsabilities' distribution of a test set cmpd

Probabilities to be "active" or "inactive") for a new compound is estimated using the Bayes equation

$$P(c_i | \mathbf{x}_k) = \frac{P(\mathbf{x}_k | c_i) \times P(c_i)}{\sum_i P(\mathbf{x}_k | c_i) \times P(c_i)}$$

K. Klimenko, G. Marcou, D. Horvath, A. Varnek J. Chem. Inf. Model. 2016, 56, 1438–1454

Toward « universal » map of chemical space

What do we expect from an "universal" map of the Chemical Space?





Map of the chemical space should:

- be representative with respect to the variety of known chemotypes;
- be able to distuinguish different activity classes and different chemotypes;
- be able to accomodate novel structures and activities in agreement with the neighborhood behavior principle

Universal map of chemical space

 As any machine-learning method, GTM is limited by the data size. For very large data sets, only part of molecules can be used for the manifold construction. Thus, a representative subset (a "Frame" Set) must be selected for this purpose.

Examples of Frame Sets



Non-representative





contour lines world map

Hecataeus of Miletus (c. 550 – 476 BCE)

Choice of descriptors – a vital issue for chemical space construction



- Positioning of objects in the initial and latent spaces depends on the choice of molecular descriptors
- Is there a way to select some « optimal » descriptors for maps construction ?

DUD dataset: GTM as a function of type of descriptors



 Γ -score – measures the ability of a model to produce clustering of similar structures in latent space. For ideal classes separation, $\Gamma = 1$.

Chemical space construction driven by (Q)SAR models



Optimal descriptors are supposed to provide with the best GTM-based regression or classification models built on some « scoring » dataset(s).

ISIDA fragment descriptors



ISIDA fragment descriptors



Several hundreds types of fragment descriptors can be generated for one same data set

Toward universal map(s) of chemical space



More activities are used, more universal is a map

Toward universal map(s) of chemical space

Descriptors type	МАР	SCORE		
Descr ₁	Map ₁	SCORE ₁		
Descr ₂	Map ₂	SCORE ₂		
		•••••		
Descr _N	Map _N	SCORE _N		

Decriptors leading to largest score are selected

Chemical space of druglike compounds



Basic assumption:

If the manifold is trained on activity predictions of ligands for >100 different biological targets, it may also accommodate other biological activities and compounds.

Chemical space of druglike compounds



Universal manifold:

- Optimized for 144 sets (GPCR, kinases, ...)
- Validated on > 450 sets

ChEMBL: chemical space of antiviral compounds

3 "universal" maps based on different types of ISIDA descriptors



actives



K. Klimenko, G. Marcou, D. Horvath, A. Varnek J. Chem. Inf. Model. 2016, 56, 1438–1454

Privilege patterns extraction

From responsibility vector to binned pattern



Responsibility Vector



Node Number n		183	184	185	186	187	188	
Probability of Residence R _n	0.0	0.005	0.020	0.022	0.031	0.021	0.002	0.0
Node Number n		183	184	185	186	187	188	
Binned "Pattern" vector	0	1	2	2	3	2	0	0

ChEMBL: chemical space of antiviral compounds

3 different maps are needed in order to extract the privilege patterns for main classes of antivirals



3 "universal" maps based on different types of ISIDA descriptors

GTM: extraction of privilege structural patterns

ChEMBL: class landscape for GPCR ligands





Evolution of privileged structural motifs of GPCR ligands in ChEMBL

Big Data challenge

Comparison of databases: GDB vs PubChem

GDB-17 – computer-generated virtual molecules containing up to 17 heavy atoms *

- whole set: 1.66 *10¹¹ virtual molecules
- « lead-like » subset used in this study: **10 M** molecules

PubChem-17 – subset of **10.8** M *real* molecules containing up to 17 heavy atoms exracted from the PubChem database

* L. Ruddigkeit et al. J Chem Inf Model 2012, 52, 2864–2875

Comparison of databases: GDB vs PubChem



Chemography: hierarchical GTM

PubChem-17 vs GDB-17



No analogues in PubChem

Databases comparison and properties profiling

Comparison of suppliers databases



* Data from: T. Petrova et al. *Med. Chem. Commun.*, 2012, 3, 571-579

GTM for the Suppliers DB (> 2 M cmpds)

Data density distributions built on responsibility vectors



H.Gaspar, Igor I. Baskin, G.Marcou, D.Horvath and A.Varnek J. Chem. Inf. Model., 2015, 55 (1), 84–94



Suppliers DB: GTM property landscapes





Suppliers DB: GTM property landscapes

High molecular weight & High chirality & Low solubility





Min

Suppliers DB map: regions of interest



Meta-GTM: individual libraries as datapoints



each library is described by GTM descriptors



Stargate GTM (S-GTM) setup



Descriptors space

	Descr1	Descr2	Descr3	
Mol 1	0.22	5.43	1.12	
Mol 2	7.18	0.96	-5.42	
Mol 3	-1.01	7.41	10.63	

Stargate GTM (S-GTM) setup



Descriptors space

	Descr1	Descr2	Descr3	
Mol 1	0.22	5.43	1.12	
Mol 2	7.18	0.96	-5.42	
Mol 3	-1.01	7.41	10.63	

S-GTM: prediction of pharmacological profile



S-GTM: discovery of structures corresponding to a given pharmacological profile



	Descr1	Descr2	Descr3
Mol 1	??	??	??
Mol 2	??	??	??
Mol 3	??	??	??

ID	Target name	Mol 1	Mol 2		Mol 1325
A2a	Alpha-2a adrenergic receptor	AFFINITY		VALUES	
D2	Dopamine D2 receptor	AFFINITY		VALUES	
D3	Dopamine D3 receptor	AFFINITY		VALUES	
D4	Dopamine D4 receptor	AFFINITY		VALUES	
S1a	Serotonin 1a (5-H1a) receptor	AFFINITY		VALUES	
S2a	Serotonin 2a (5-H2a) receptor	AFFINITY		VALUES	
S7	Serotonin 7 (5-HT7) receptor	AFFINITY		VALUES	
ST	Serotonin transporter	AFFINITY		VALUES	

S-GTM: retrieval of structures corresponding to a given pharmacological profile



H. A. Gaspar , I. I. Baskin, G. Marcou, D. Horvath, A. Varnek J. Chem. Inf. Model., 2015, 55 (11), 2403–2410

S-GTM: retrieval of structures corresponding to a given pharmacological profile





ISIDA-GTM Software

• GTM Algorithms:

« classical », incremental and kernel,

• Modeling:

regression and classification models

• Visualization:

data points, data probability distribution, activity landscapes, chemical structures































