

## [P51] GTM Applicability Domain for Classification and Regression models

Hélène Gaspar<sup>1</sup>, Timur Gimadiev<sup>1,2</sup>, Gilles Marcou<sup>1</sup>, Dragos Horvath<sup>1</sup>, Alexandre Varnek<sup>1</sup>

<sup>1</sup>Laboratoire de Chémoinformatique, UMR 7140 CNRS, Université de Strasbourg, 1 rue B. Pascal, 67000 Strasbourg, France.

<sup>2</sup>Butlerov Institute of Chemistry, Kazan Federal University, Kazan, Russia

Generative Topographic Mapping (GTM) is a dimensionality reduction method, and the probabilistic counterpart of Kohonen maps. Each  $i$ -th molecule in  $N$ -dimensional initial space is projected into the  $k$ -th node of 2D latent space with a probability  $R_{ik}$ , so that each compound is represented both by a mean position (a point) on a 2D map, and a probability distribution  $R_i$ , which may be used for predictions of activity (property) of new compounds.

Here, we suggest several different GTM-based definitions of applicability domain (AD) of both regression and classification models. This concerns the approaches involving: (i) a likelihood threshold, (ii) relative population of nodes, (iii) class entropy, and, (iv) ratio of classes' probabilities. These approaches are demonstrated for regression models for stability constants of ligand-metal complexes, and GTM-based classification models [1] for Biopharmaceutics Drug Disposition Classification System (BDDCS) [2] and inhibitors of P-glycoprotein 1 (Pgp), an ATP-dependent efflux pump.

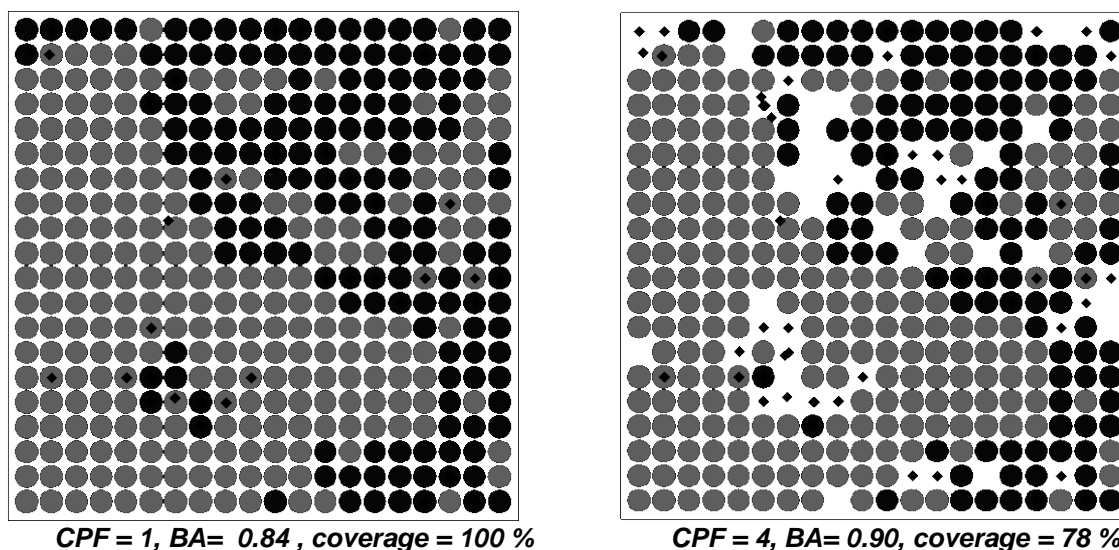


Figure 1. Graphical interpretation of the applicability domain for GTM classification models. On the map prepared for the entire set of 1568 molecules of inhibitors (dark grey) and non inhibitors (light grey) of P-glycoprotein 1, the color stands for the class having the highest probability compared to the other in a given node. Black points correspond to incorrectly classified molecules. The increase the class prevalence factor ( $CPF = \text{Probability of major class} / \text{Probability of another class}$ ) from  $CPF=1$  (left) to 4 (right) results in shrinking the AD area. This leads to the decrease of the number of molecules inside AD (coverage), on one hand, and to the increase of the model's performance (Balance Accuracy, BA), on the other hand.

[1] Kireeva, N., Baskin, I. I., Gaspar, H. A., Horvath, D., Marcou, G. and Varnek, A., Mol. Inf., 31 (2012): 301–312. doi: 10.1002/minf.201100163

[2] Gaspar, H. A., Marcou, G., Horvath, D., Arault, A., Lozano, S., Vayer, P., and Varnek, A. J. Chem. Inf. Model. 53 (2013), 3318–3325.