

POSTER PRESENTATION

Open Access

The influence of training actives/inactives ratio on machine learning performance

Rafał Kurczab^{1*}, Sabina Smusz^{1,2}, Andrzej J Bojarski¹

From 8th German Conference on Chemoinformatics: 26 CIC-Workshop
Goslar, Germany. 11-13 November 2012

In drug discovery, machine learning is widely used to classify molecules as active or inactive against a particular target. The vast majority of these methods (supervised learning) needs a training set of objects (molecules) to develop a decision rule that can be used to classify new entities (the test set) into one of the two mentioned classes [1].

A lot of studies, searching an optimal learning parameters and their impact on classification effectiveness were performed [2,3]. Unfortunately, there is no data showing the influence of actives/inactives ratio, used to model training, on the efficiency of new active compounds identification. Therefore, the main goal of this study was to examine the impact of changing the number of inactives in the training set with fixed amount of actives. For a given ratio, the inactives were randomly selected from ZINC database (10-times to prevent an overestimations error). This concept was verified on three different protein targets (i.e. 5-HT_{1A}, HIV-1 protease and matrix metallo-proteinase) and a set of algorithms (SMO, Naïve Bayes, Ibk, J48 and Random Forest) implemented in WEKA package [4]. To compounds representation, two types of molecular fingerprints were used (MACCS and hashed fingerprint), to determine their possible impact on machine learning performance.

Acknowledgements

The study was supported by a grant PRELUDIUM 2011/03/N/NZ2/02478 financed by the National Science Centre.

Author details

¹Department of Medicinal Chemistry, Institute of Pharmacology Polish Academy of Sciences, Kraków, 31-343, Poland. ²Faculty of Chemistry, Jagiellonian University, Kraków, 30-060, Poland.

* Correspondence: kurczab@if-pan.krakow.pl

¹Department of Medicinal Chemistry, Institute of Pharmacology Polish Academy of Sciences, Kraków, 31-343, Poland

Full list of author information is available at the end of the article

Published: 22 March 2013

References

1. Melville JL, Burke EK, Hirst JD: **Machine learning in virtual screening.** *Comb Chem & High Thr Scr* 2009, **12**:332-343.
2. Ma XH, Wang R, Yang SY, Li R, Xue Y, Wei YC, Low BC, Chen YZ: **Evaluation of virtual screening performance of support vector machines trained by sparsely distributed active compounds.** *J Chem Inf Mod* 2008, **48**:1227-1237.
3. Plewczynski D, Spieser SH, Koch U: **Assessing different classification methods for virtual screening.** *J Chem Inf Mod* 2006, **46**:1098-106.
4. Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH: **The WEKA data mining software: an update.** *SIGKDD Explorations* 2009, **11**(1):10-18.

doi:10.1186/1758-2946-5-S1-P30

Cite this article as: Kurczab et al.: The influence of training actives/inactives ratio on machine learning performance. *Journal of Cheminformatics* 2013 **5**(Suppl 1):P30.

Publish with **ChemistryCentral** and every scientist can read your work free of charge

"Open access provides opportunities to our colleagues in other parts of the globe, by allowing anyone to view the content free of charge."

W. Jeffery Hurst, The Hershey Company.

- available free of charge to the entire scientific community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
<http://www.chemistrycentral.com/manuscript/>


ChemistryCentral