

QSPR Models for Predictions and Data Quality Assurances: Melting Point and Boiling Point of Perfluorinated Chemicals

Barun Bhatarai^a, Wolfram Teetz^d, Tomas Öberg[†], Tao Liu[†], Nina Jeliazkova[‡], Nikolay Kochev[§], Ognyan Pukalov[§], Igor V. Tetko^d, and Paola Gramatica^{*,†,‡}

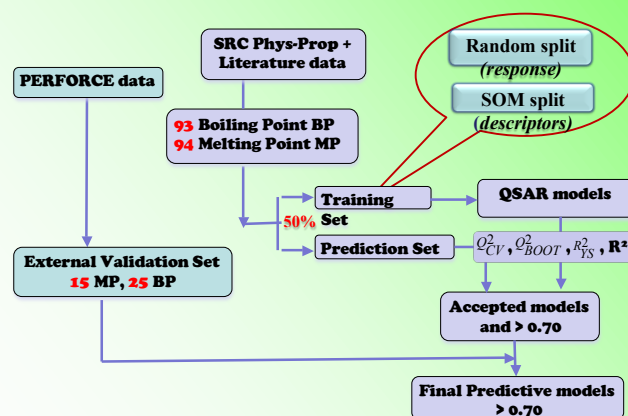
^aQSPAR Research Unit in Environmental Chemistry and Ecotoxicology, Department of Structural and Functional Biology (DBSF), University of Insubria, via JH Dunant 3, Varese, 21100, Italy. Email: paola.gramatica@uninsubria.it; [†]School of Pure and Applied Natural Sciences, Linnaeus University, SE-391 82, Kalmar, Sweden. Email: Tomas.Oberg@lnu.se; [‡]Ideaconsult Ltd, 4 A. Kanchev str., Sofia 1000, Bulgaria; [§]Department of Analytical and Computer Chemistry, University of Plovdiv, 24 Tsar Assen Str., Plovdiv 4000, Bulgaria. Email: nina@idea.bg; [†]Institute of Bioinformatics and Systems Biology, Helmholtz Zentrum Muenchen - German Research Center for Environmental Health, Ingolstaedter Landstrasse 1, D-85764 Neuherberg, Germany. Email: itetko@vclab.org

ABSTRACT

Quantitative structure-property relationship (QSPR) studies on Melting Point and Boiling Point of Perfluorinated Chemicals (PFCs) are presented. PFCs are studied under the EU-FP7 funded **CADASTER** project to understand its behavior in biota and environment. They are considered as 'emerging pollutants' and found widely distributed in the environment, released due to their widespread use in different household and industrial products as cleansers, fire-fighting foams, micelles, repellants for leather, paper, and textiles etc. Continues exposure of these chemicals is found to be the source of bio-accumulation in body parts of human, wildlife and is ultimately becoming the cause of toxic reactions and poisoning. Models are developed using SRC PhysProp data as described below. In addition, the predictive performances of the developed models were verified on a blind external validation set (EV-set) prepared from experimental values available from PERFORCE database. This database contains only long chain perfluoro-alkylated chemicals, particularly monitored by regulatory agencies like US-EPA and EU-REACH. QSPR modeling using different approaches, internal and external validation on two different prediction sets and studies of the applicability domain highlight the robustness and high accuracy of the proposed models. Finally, Melting Point for additional 397 PFCs and Boiling Point for 364 PFCs for which experimental measurements are unknown were predicted, verifying their applicability domain. The set of descriptors which best describes the structure-property relationship, the similarities, and the differences observed will be discussed as well as the consensus model predictions.

MATERIALS AND METHODS

	HMGU, Germany	IDEA Consult, Bulgaria	UI, Italy	LNU, Sweden
Descriptors	E-State indices	Fragment based	DRAGON (oD - 2D)	
Descriptor Selections	Pearson Pairwise Correlation	Exhaustive isomorphism search of fragment against structure	Pearson Pairwise Correlation & Genetic Algorithm	variable influence on projection (VIP)
Descriptors used for Modeling	MP = 87 indices BP = 66 indices	MP = 3 descriptors BP = 8 descriptors	MP = 3 descriptors BP = 3 descriptors	MP = 37 descriptors at 3 components BP = 149 descriptors at 4 components
Methods	Associative Neural Network (ASNN) Architecture: 10x3x1	Multiple Linear Regression (MLR) using ordinary-least-squares (OLS)		Partial least squares regression (PLSR)
External validation	Double: Prediction sets by splitting and blind External Validation set			Single: External Validation Set
Structural Applicability Domain	Distance to model (DM) on standard deviation of ensemble prediction, 5x-cross-validation	Williams plot for response outliers Leverage approach (H matrix) for structural chemical domain		residual standard deviation (Euclidean distance) and leverage (Mahalanobis distance)



Data Quality Assurance

CAS	Endpoint reported	Data from PhysProp (°C) used by UI, LNU, IDEA	UI Predictions	LNU Predictions	Data (°C) used by HMGU and the references	HMGU Predictions
76-16-4	MP	-101.00	-155.71	-138.33	-155.60	-111.656
307-34-6	MP	-42.0	-50.36	-54.73	-56.80	-57.435
307-55-1	MP	108.0	75.63	107.29	111.0 [33]	66.166
354-32-5	MP	146	2.45	-91.56	-146.0 [34]	-86
375-22-4	MP	-17.5	13.40	-1.99	-18.0 [33]	13.248
423-55-2	MP	25*	-22.74	-40.99	-6.0 [35]	-59.167
1493-13-6	MP	25*	50.26	-31.38	-40.0 [36]	-12.567
426-65-3	MP → BP	75.5	32.29	-21.43	n/a [37]	n/a
355-46-4	BP	238.5	228.71	241.87	225.0 [38]	221.176
375-73-5	BP	211.0	196.93	207.33	200.0 [39]	191.358

RMSE Comparison

	RMSE			
	EPI	UI	LNU	HMGU
Melting Point (94)	47.97	44.09	32.81	36
Boiling Point (93)	24.80	23.24	14.12	32

RESULTS AND DISCUSSION

Dragon descriptors allow a complex and differentiated view on the molecule, while E-State indices give a more uniform description and Fragment based descriptors provide an easily interpretable base for modeling. For simple properties like boiling point, fitting a small, variable-selected MLRA model to the data subset provides excellent results. This approach is also robust against erroneous data. At the complexity level of melting point (or e.g. vapor pressure), this approach competes in quality with E-State-ASNN models, that are easily obtainable from scratch. The well interpretable but quite tedious approach using selected fragment descriptors results in a slight drop in model quality. Also, since literature data is often published for classes of compounds that are directly connected to fragments, common systematic errors (such as pressure variations for BP) give low RMSE of the models but inadequate models, so extra care has to be taken here in the validation step. It is remarkable that data collected from the databases has a high number of errors like mixed up algebraic signs or approximated values, so that data validation and overlap is necessary. Here, the relation between BP and MP gives valuable information that can be employed. As expected, the accuracy of the prediction models is better than for 'generic' boiling and melting point models.

CONCLUSIONS

The results fit our experience that a consensus model, built from independently developed models using different descriptors and using different algorithms, delivers the best prediction results. In the special case of PFCs, simple statistical algorithms applied to complex descriptors perform about as good as complex algorithms applied to simple descriptors. Developing both types of models enables a more specialized and also more detailed look on outliers and opens lots of possibilities to analyze them. Chemical interpretation of and experimental design emerging from the models benefit from having a set of models representing different views of the underlying mechanics.

REFERENCES

- [1] ESTATE: Kier, L. B.; Hall, L. H. An Electrotopological State Index for Atoms in Molecules. *Pharm. Res.* 1990, 7, 801-807.
- [2] DRAGON: TALETTI, E.; Via, V.; Pisani, F. 2012, Milano - Italy
- [3] ASNN, Tetko, I. V. Associative neural network. *Neural Processing Letters*, 2002, 16, 187-199.
- [4] MAHALANOBIS, Mahalanobis, P. C. (1936). "On the generalised distance in statistics". *Proceedings of the National Institute of Sciences of India* 2 (1): 49-55. <https://ic.isical.ac.in/dspace/handle/123456789>.
- [5] AD, Jaworska, J., Nikolova-Jeliazkova, N., Aldenberg, T. (2005). QSAR applicability domain estimation by projection of the training set in descriptor space: A review. *ATLA-Alternatives to Laboratory Animals* 33(3): 445-459.
- [6] AD, Paps, E.; Kevacik, S.; Gramatica P. Development, Validation and Inspection of the Applicability Domain of QSPR Models for physico-chemical properties of Polybrominated DiphenylEthers QSAR and Combinatorial Science, 2009, 28, 790-796.

Financial support by European Union through the project CADASTER FP7-ENV-2007-1-212668