

Interpretation of QSAR models

Pavel Polishchuk

Institute of Molecular and Translational Medicine Faculty of Medicine and Dentistry Palacky University

> pavlo.polishchuk@upol.cz pavel_polishchuk@ukr.net qsar4u.com

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Interpretability vs. predictivity of QSAR models





Machine learning



Representative learning

Why interpretation is important?

Found active/inactive patterns which can be used for optimization of compound properties

Retrieve trends of stricture-activity relationships which can be used for knowledge-base model validation and for better understanding of the studied end-point (formulate hypothesis)

Regulatory purposes

OECD principles for the validation, for regulatory purposes, of (Q)SAR models

- 1) a defined endpoint
- 2) an unambiguous algorithm
- 3) a defined domain of applicability
- 4) appropriate measures of goodness-of-fit, robustness and predictivity
- 5) a mechanistic interpretation, if possible

Model should be predictive

Interpretation is valid within the applicability domain of the model

Interpretation results are data set dependent

Local vs. global interpretation

Local interpretation

considers mutual influence of atoms (fragments) in a single compound



Global interpretation

summarizes interpretation results across studied compounds to reveal general structure-activity relationship trends

$F \rightarrow CI \rightarrow Br \rightarrow I$

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Learning and interpretation

Learning



"model → descriptor → (structure)" paradigm

"model \rightarrow descriptor \rightarrow (structure)" paradigm

Machine learning-specific interpretation approaches

- Linear models (LR, PLS, OPLS, etc)
- **Decision tree**
- Random Forest
- Neural nets
- Support vector machine
- Rule-extraction

Machine learning-independent interpretation approaches Variable importance Sensitivity analysis

Partial derivatives

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Hansch approach

plant growth inhibition activity of phenoxyacetic acids $1/C = 4.08\pi - 2.14\pi^2 + 2.78\sigma + 3.38$ rate of penetration of membranes in the plant cell

 $\pi = \log P_X - \log P_H$ σ - Hammet constant

Hansch, C.; Maloney, P. P.; Fujita, T.; Muir, R. M. *Nature* **1962**, *194*, 178-180

Free and Wilson approach

Inhibition activity of compounds against *Staphylococcus aureus*



R is H or CH_3 ; X is Br, Cl, NO_2 and Y is NO_2 , NH_2 , $NHC(=O)CH_3$

$$Act = 75R_{H} - 112R_{CH3} + 84X_{CI} - 16X_{Br} - 26X_{NO2} + 123Y_{NH2} + 18Y_{NHC(=0)CH3} - 218Y_{NO2}$$

Free, S. M.; Wilson, J. W. Journal of Medicinal Chemistry 1964, 7, 395-399

Interpretation of decision tree models

Solid-phase fluorescence enhancement of 2-(diphenylacetyl)-1,3-indandione 1-(p-(dimethylamino)benzaldazine)

ANALYTICAL CHEMISTRY, VOL. 57, NO. 9, AUGUST 1985 . 1953



Figure 5. Decision tree: compound/II fluorescence enhancement.

Ashman, W.P., Lewis, J.H., Poziomek, E.J., Analytical Chemistry, **1985**, 57(9): p. 1951-1955.

14

Rule-extraction approaches



Martens, D., Baesens, B., Gestel, T.V., IEEE Transactions on Knowledge and Data Engineering, **2009**, 21(2): p. 178-191.

Rule-extraction approaches



Raccuglia, P., Elbert, K.C., Adler, P.D.F., Falk, C., Wenny, M.B., Mollo, A., Zeller, M., Friedler, S.A., Schrier, J., Norquist, A.J., Nature, **2016**, 533(7601): p. 73-76.

Interpretation of random forest models



- $S_{k,i}$ contribution of i-th descriptor in k-th compound
- T number of trees

LS_{i,i} - local contribution of i-th descriptor where compound k fits the node

17

Kuz'min, V.E., Polishchuk, P.G., Artemenko, A.G., Andronati, S.A., Molecular Informatics, **2011**, 30: p. 593-603.

Visualization of descriptor contributions



Interpretation of random forest models



347 agonists of 5-HT_{1A} receptor

- Ar substituted (hetero)aryls
- L polymethylene chain
- R various (poly)cyclic residues



Kuz'min, V.E., Polishchuk, P.G., Artemenko, A.G., Andronati, S.A., Molecular Informatics, **2011**, 30: p. 593-603.

Message

There are a lot of approaches which can be applied for interpretation of specific models

Interpretation results converge regardless machine learning method and interpretation approach used

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Variable importance

Gives information about relative importance of descriptors used for model building but not about the direction of their influence (positive or negative)

- add noise to input variables (descriptors)¹
- permute variable values²

The more prediction accuracy is dropped the more important variable is.

¹ Györgyi, G., Physical Review Letters, **1990**, 64(24): p. 2957-2960. ² Breiman, L., Machine Learning, **2001**, 45(1): p. 5-32.

Variable importance

toxicity on Tetrahymena Pyriformis



Polishchuk, P.G., Muratov, E.N., Artemenko, A.G., Kolumbin, O.G., Muratov, N.N., Kuz'min, V.E., Journal of Chemical Information and Modeling, **2009**, 49(11): p. 2481-2488.

Sensitivity analysis

Explores the dependence of output values to systematic changes in descriptor values while values of all other descriptors remain constant

Sensitivity analysis

inhibitors of dihydrofolate reductase (substituents X are at 3, 4, and 5 positions)

 NH_2

linear model

	previous	updated
MR ₅ ³		11.79
MR ₅ ²		15.74
MR ₅	0.95	6.55
MR ₃	0.89	0.89
MR ₄	0.80	0.80
MR ₄ ²	-0.21	-0.21
π ₃	1.58	1.58
Log(β*10 ^{π3} + 1)	-1.77	-1.77
intercept	6.65	6.24
RMSE	0.093	0.074

So, S.S., Richards, W.G., Journal of Medicinal Chemistry, **1992**, 35(17): p. 3201-3207.

Partial derivatives (or local sensitivity analysis) estimates "regression coefficients" from nonlinear models for a particular data point (compounds) in a chemical space

$$f = ax + by + c$$
$$\frac{\partial f}{\partial x} = a$$
$$\frac{\partial f}{\partial y} = b$$

Aoyama, T., Ichikawa, H., Journal of Chemical Information and Computer Sciences, **1992**, 32(5): p. 492-500.

Numerical solution

Hasegawa, K., Keiya, M., Funatsu, K., Molecular Informatics, 2010, 29(11): p. 793-800.

- + regression/classification
- + linear, non-linear, consensus
- k-NN models
- proper estimation of differentiation error

"model \rightarrow descriptor \rightarrow (structure)" paradigm

All models regardless machine learning method used are interpretable

Limitation: use of only interpretable descriptors

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Learning and interpretation

Learning

Universal structural interpretation of QSAR models

Activity _{pred} (A)	Activity _{pred} (B)	Contribution(C)
<i>f</i> (A) = x	<i>f</i> (B) = y	W(C) = x - y

Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. *Molecular Informatics* **2013**, *32*, 843-853

"Model \rightarrow structure" interpretation approaches

1. Similarity maps

Riniker, S.; Landrum, G. Similarity maps - a visualization strategy for molecular
fingerprints and machine-learning methods. *Journal of Cheminformatics* 2013, 5,
43

2. Universal structural interpretation

Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. Universal Approach for Structural Interpretation of QSAR/QSPR Models. *Molecular Informatics* **2013**, *32*, 843-853

3. Computational matched molecular pairs

Sushko, Y.; Novotarskyi, S.; Korner, R.; Vogt, J.; Abdelaziz, A.; Tetko, I. Predictiondriven matched molecular pairs to interpret QSARs and aid the molecular optimization process. *Journal of Cheminformatics* **2014**, *6*, 48

Similarity maps (per atom interpretation)

Dopamine D3 inhibitors

NB

+ implemented in RDKit (<u>http://rdkit.org</u>)

- implementation supports only classification models built with atom pairs, topological torsions or Morgan fingerprints

Riniker, S., Landrum, G., Journal of Cheminformatics, **2013**, 5(1): p. 43.

RF

Universal structural interpretation vs. Free-Wilson

Comparison with Free-Wilson

Polishchuk P, Tinkov O, Khristova T, Ognichenko L, Kosinskaya A, Varnek A, Kuz'min V. *Journal of Chemical Information and Modeling* **2016**, 56, 1455-1469.

Universal structural interpretation (Ames)

Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. *Molecular Informatics* **2013**, *32*, 843-853

Universal structural interpretation

Polishchuk P, Tinkov O, Khristova T, Ognichenko L, Kosinskaya A, Varnek A, Kuz'min V. *Journal of Chemical Information and Modeling* **2016**, 56, 1455-1469.

Message

Interpretation results converge regardless machine learning method and descriptors used

Implemented in open-source SPCI software <u>https://github.com/DrrDom/spci</u> and R package <u>https://github.com/DrrDom/rspci</u>

76 SPCI - structural and physico-chemical interpretation of QSAR models			X	
Build models Calc contributions Plot contributions Predict				
Structural/physico-chemical interpretation (result in different descriptors)				
 Structural & functional (Chemaxon required) 				
Structural only (no Chemaxon usage)				
SDF with compounds				
Path to SDF-file	property	field name		
	Browse		•	
Optional. Compound names. External text file with compound property values				
 Automatically generate compound names 				
O Use compound titles from SDF file				
Use field values as compound names from SDF file	-			
Path to text file with property values				
· · · · · · · · · · · · · · · · · · ·	Browse			
Models				
 Regression (RF, GBM, SVM, PLS) 	 Binary classification (0-1) (RF, GBM, SVM) 			
Random Forest (RF)	Random Forest (RF)			
Support vector regression (SVR)	✓ Support vector classification (SVC)			
Gradient boosting regression (GBR)	Gradient boosting classification (GBC)			
Partial least squares (PLS)	k-Nearest neighbors (kNN)			
k-Nearest neighbors (kNN)				
Number of cores to use				
3 🔹				
Build models Show statistics				
(c) Pavel Polishchuk 2014-2016				

Prediction-driven MMP

Sushko, Y., Novotarskyi, S., Korner, R., Vogt, J., Abdelaziz, A., Tetko, I. Journal of Cheminformatics, **2014**, 6(1): p. 48.

Message

Prediction-driven MMP can improve coverage of chemical space (fill gaps)

Ignores molecular context which can be important

Implemented on http://ochem.eu

Design of molecular series

1,3 dinitrobenzene has higher toxicity than 1,2- and 1,4-dinitrobenzes that may be a result of a different mechanism of action, e.g. it may act as uncoupler of oxidative phosphorylation like 2,4 dinitrophenol

Tinkov, O.V., Ognichenko, L.N., Kuz'min, V.E., Gorb, L.G., Kosinskaya, A.P., Muratov, N.N., Muratov, E.N., Hill, F.C., Leszczynski, J., Structural Chemistry, **2016**, 27(1): p. 191-198.

42

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Analysis of context dependence of fragment contributions

Analysis of context dependence of fragment contributions

Toxicity on Tetrahymena pyriformis

unpublished results

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Structural interpretability of data sets (global interpretation)

Structural interpretability of data sets. Case 1

Blood-brain barrier permeability

Polishchuk, P., Tinkov, O., Khristova, T., Ognichenko, L., Kosinskaya, A., Varnek, A., Kuz'min, V., Journal of Chemical Information and Modeling, **2016**, 56(8): p. 1455-1469.

48

Structural interpretability of data sets. Case 2

 $\alpha_{IIb}\beta_3$ antagonists – RGD mimetics

Polishchuk, P., Tinkov, O., Khristova, T., Ognichenko, L., Kosinskaya, A., Varnek, A., Kuz'min, V., Journal of Chemical Information and Modeling, **2016**, 56(8): p. 1455-1469.

49

Structural interpretability of data sets. Case 3

347 agonists of 5-HT_{1A} receptor

- Ar substituted (hetero)aryls
- L polymethylene chain
- R various (poly)cyclic residues

Kuz'min, V.E., Polishchuk, P.G., Artemenko, A.G., Andronati, S.A., Molecular Informatics, **2011**, 30: p. 593-603.

Outline

- 1. Introduction. Importance, principles and issues of interpretation of QSAR models. Global and local interpretation.
- "Model → descriptor → (structure)" interpretation paradigm
 - a) Machine learning-specific interpretation approaches
 - b) Machine learning-independent interpretation approaches
- 3. "Model \rightarrow structure" interpretation paradigm
- 4. Context dependence of interpretation results
- 5. Interpretability of data sets
- 6. Conclusion

Applicability of interpretation approaches

Madala	Descriptors		
Models	interpretable	non-interpretable	
linear regression	regression coefficients (Hansch, Free-Wilson)		
PLS (OPLS, O2PLS, etc)	regression coefficients, X- and Y-scores, variable importance		
decision trees	logical rules	universal structural interpretation.	
NN	variable importance based on weights and biases, variable contributions		
RF	variable importance based on permutation, variable contributions	similarity maps, computational matched	
NN, SVM, RF	rule extraction	molecular pairs and series	
any model including consensus ones	partial derivatives, variable importance based on permutation, sensitivity analysis		
Interpretation paradigm	model → descriptors → (structure) or model → structure	model → structure	

Conclusion

Interpretation results of valid predictive models should converge independent of:

interpretation approach descriptors machine learning method

All models are interpretable but not all end-points

OURNAL OF CHEMICAL INFORMATION AND MODELING

pubs.acs.org/jcim

Interpretation of Quantitative Structure–Activity Relationship Models: Past, Present, and Future

Pavel Polishchuk*®

Institute of Molecular and Translational Medicine, Faculty of Medicine and Dentistry, Palacký University and University Hospital in Olomouc, Hněvotínská 1333/5, 779 00 Olomouc, Czech Republic

