On-line chemoinformatics tools for SAR/QSAR model development

Igor V. Tetko Institute of Structural Biology Helmholtz Zentrum Muenchen

13th October 2016, on-line course of the BIGCHEM project

HelmholtzZentrum münchen German Research Center for Environmental Health

Outline

Overview of existing on-line services on the web

- Individual groups web servers
- Projects web servers
- Model repositories
- Model development platforms

OCHEM

- Typical uses of OCHEM
- Model development
- Model interpretation
- Model evaluation

Individual groups web servers

http://cdb.ics.uci.edu

ChemDB Chemoinformatics Portal

Home

Molecules

ChemicalSearch

Find a molecule by its name, structure, or similarity to another molecule and filter the results.

Virtual Chemical Space

Interactively deconstruct a target molecule into possible chemical precursors and reassemble them into a combinatorial library of real or virtual molecules around the target.

MOLpro

Calculate or predict molecular properties other than 3D structure.

AquaSol

Predict aqueous solubility of small molecules using deep learning and ensembles.

COSMOS (downloadable)

Predict 3D molecular structures using open crystallography libraries.

COSMOS Predict 3D molecular structures.

Reactions

Reaction Explorer

Learn and practice reactions, syntheses, and mechanisms interactively with support for: automated generation of problems, curved-arrow mechanism diagrams, and inquiry-based learning.

Reaction Predictor

Predict reaction outcomes and mechanisms using machine learning.

Prof. P. Baldi, UCI

Tools

Smi2Depict Generate 2D images from SMILES.

Babel

Convert between molecule file formats.

Reaction Processor Generate product libraries.

Pattern Match Counter Count functional groups (sub-structures).

Pattern Count Screen Screen molecules by functional group count.

MSFragment Fragment molecules for mass spec analysis.

Mass2Structure Search ChemDB by monoisotopic mass and substructure filtering.

Protein Target Predictor

Predict activities of small molecules against a large set of protein targets

Datasets

Chemical datasets Datasets for training and testing machine learning and other algorithms.

http://infochim.u-strasbg.fr/webserv/VSEngine.html



Prof. A. Varnek, Strasbourg

http://tox.charite.de/tox



Dr. R. Preissner, Charite University of Medicine

Projects web servers

http://www.vcclab.org

http://www.vcclab.org



Home About Partners Software Articles Servers Jobs Web Services How to cite? Contact

Copyright 2001 -- 2016 http://www.vcclab.org. All rights reserved.

A first QSAR web tool: Polynomial Neural Networks

Virtual Computational Chi	enistry Laboratory
	http://www.vcclab.org

Welcome to the PNN program!

© MIPS/VCCLAB PNN parameters	Login	submit your task		
Type of data ANALYSIS	LOO of the training set	no ᅌ		
Maximum DEGREE of the model + 3-order	Number of <u>ITERATIONS</u>	15 🗘		
Maximal NUMBER of terms in model	Number of <u>STORED</u> models	3 🗘		
CRITERION to select the best models FPE	VALIDATION set in RR criterion			
Details of calculated results (PRINT)				
🧹 display calculated vs experimental values	save input data in stdout			
save calculated values in stdout	save detailed statistics of a	nalysis in stdout		
save statistics of input data in stdout	select all options	select all options		
Specify polynomial neural network parameters. Cl	lick the uderlined links for more informa	ition.		

I. V. Tetko, T. I. Aksenova, V. V. Volkovich, T. N. Kasheva, D. V. Filipov, W. J. Welsh, D. J. Livingstone, A. E. P. Villa, SAR QSAR Environ. Res. 2000, 11, 263-280.

http://knimewebportal.cosmostox.eu

C SMOS KNIME WEBPORTAL

Integrated In Silico Models for the Prediction of Human Repeated Dose Toxicity of COSMetics to Optimise Safety

Welcome to the COSMOS KNIME WebPortal

COSMOS is a unique collaboration addressing the safety assessment needs of the cosmetics industry, without the use of animals. The main aim of COSMOS is to develop freely available tools and workflows to predict the safety to humans following the use of cosmetic ingredients.

The models developed within COSMOS have been implemented into KNIME workflows. The KNIME Analytics Platform integrates access to chemical databases, data processing and analysis, modelling approaches, profiling of structures and calculation of properties in a flexible way.

The COSMOS KNIME WebPortal versions of these workflows allow users, not experienced with the software, to execute the workflows in a web browser, without local software installation required. The results are downloadable as pdf reports and other formats e.g. Excel sheets.

The models are documented and user guidance is available through COSMOS Space. A list of all workflows available can be found here.

Login		
Username		
Password		
Login		
How to get your login		
Registration is free. To get your login details from the	COSMOS Space follow the folling	ng steps:
1. Visit http://cosmosspace.cosmostox.eu		

2. Click on Login

3. Click on register

4. Fill in your email address and choose a password.

http://cosmosspace.cosmostox.eu/app/space/cosmosshare/ckwdprojects?pn=2

-

	C	SMC	DSOS	pace	\bigcirc	
ABOUT	CONTACT					LOGOUT
Home	COSMOS Share	COSMOS DB v1 dump	COSMOS KNIME Workflows	My COSMOS Space	My Gallery	My Profile
(L	COSMOS Sh .ist of all availa	are - COSMOS H ble workflows with	KNIME Workflow Do	cumentations s and user guidar	ice.	
т У	Title	reely available for direct of ceredentials.	execution in a web browser at	COSMOS KNIME Wel	oPortal - login	with
	Hepatotoxicity Aler	ts (WebPortal version)		Hepatotoxicity, Me Events (MIEs), gro structural alert	plecular Initiatin uping chemical	g s,
	LXR binding predic	tion (PLS-DA/RDKit) (Deskt	op version)	LXR; QSAR; nucle classification	ar receptor;	
	Mitochondrial Toxic	city Alerts (WebPortal versio	in)	Mitochondrial toxi Initiating Events (N chemicals, structu	city, Molecular /IEs), grouping ıral alert	
	Potential NR ligand	ls and alerts towards hepat	osteatosis (Desktop Version)	NR, endocrine dis hepatosteatosis, r	ruptor, fatty live nuclear receptor	r, r
	Chemical Space Ar	nalysis (Report)		Chemical Space; I Analysis; PCA	Explorative	
	Covalent Protein B	inding Alerts (WebPortal ver	rsion)	Molecular Initiatin Protein binding, g chemicals, structu	g Events (MIEs) rouping ıral alert	,
	Covalent DNA Bind	ling Alerts (WebPortal version	on)	DNA binding, Mol Events (MIEs), gro structural alert	ecular Initiating uping chemical	s,
	Chemical Space Ar	nalysis		Chemical Space; I Analysis; PCA	Explorative	
	PAMPA logPm prec	dictor (Desktop version)		GIA, PAMPA perm logPm	eability, QSAR,	
	Ranking (WebPorta	l version)		ranking; priority se prediction	etting; screening	9;

http://www.vega-qsar.eu/



http://toxpredict.org



Unable to connect

Firefox can't establish a connection to the server at apps.ideaconsult.net:8080.

- The site could be temporarily unavailable or too busy. Try again in a few moments.
- If you are unable to load any pages, check your computer's network connection.
- If your computer or network is protected by a firewall or proxy, make sure that Firefox is permitted to access the Web.

Try Again

Model repositories

http://qsardb.jrc.it/qmrf/



#	<u>QMRF#</u> Θ	Title 9	Last updated 😐	View	Download 🥹
1	Q50-54-55-501	BIOVIA toxicity prediction model – Ames Mutagenicity	2016-6-17 14:58	, O	🎫 <mark>🎦 🖾</mark> 📾
2	Q51-54-55-502	BIOVIA toxicity prediction model – rat oral LD50	2016-6-17 14:58	, O	🎫 📴 🔟 📠
3	Q50-54-55-503	BIOVIA toxicity prediction model – NTP carcinogenicity call (male rat)	2016-6-17 14:58	, O	🎫 <mark>🎦 🖾</mark> 🛲

For information about this site please contact JRC-IHCP-COMPUTOX@ec.europa.eu

This page has been accessed 508795 times since 2008-07-03 15:25:48.0

Developed by Ideaconsult Ltd. (2007-2008) on behalf of JRC

http://qsardb.org



Recent submissions to the repository

07 Oct 2016	Ran, Y.; He, Y.; Yang, G.; Johnson, J. L. H.; Yalkowsky, S. H. Estimation of aqueous solubility of organic compounds by using the general solubility
	equation. Chemosphere 2002, 48, 487–509.
27 Sep 2016	Li, F.; Wu, H.; Li, L.; Li, X.; Zhao, J.; Peijnenburg, W. J. G. M. Docking and QSAR study on the binding interactions between polycyclic aromatic
	hydrocarbons and estrogen receptor. Ecotoxicology and Environmental Safety 2012, 80, 273 - 279.
20 Sep 2016	Papa, E.; Pilutti, P.; Gramatica, P. Prediction of PAH mutagenicity in human cells by QSAR classification. SAR and QSAR in Environmental Research
	2008, 19, 115–127.
19 Sep 2016	Kar, S.; Deeb, O.; Roy, K. Development of classification and regression based QSAR models to predict rodent carcinogenic potency using oral slope
	factor. Ecotoxicology and Environmental Safety 2012, 82, 85–95.

Model development platforms

http://chembench.mml.unc.edu



We thank the following commercial sponsors for their support:











Data storage and model development: http://ochem.eu



Typical use of OCHEM



Working with ToxAlerts

Online chemical database with modeling environment										
Home -	Database -	Models -	Ioderation -							
Ex Sea mea our i	Compound Molecules Properties Conditions Units Articles/Boo Journals	I properties Dks	•	our possible actions t data: experimentally ed to public access by your data.						
Cre	ToxAlerts		•	ToxAlerts home						
Build	Pathways		•	View alerts						
prop	Baskets			Screen compounds against alerts						
expe	Tags			Upload new alerts						
D	Set area of	interest	_							

particu technic can be	Structural alerts (also known as lar type of toxicity. The studies p ue to detect potentially toxic cho a good practice to complement	"toxicophores") are molecular performed last decade has sho emicals. Screening chemical c the QSAR models and to help	patterns known to be associated with wn that structural alerts is an efficient ompounds against known structural alerts interpreting their predictions.
search	ToxAlerts is a platform for scree	ning chemical compounds aga	ainst structural alerts. The platform allows
Scarch	Suddural alerts, introduce your		an ibranes for alerentary compounds.
		Lieland new slads	Screen your molecules
	View available alerts	Upload new alerts	Screen your molecules

ToxAlerts

- Screening of compounds against published toxicity alerts, groups, frequent hitters
- Filter alerts by endpoints or publications
- Create or upload custom SMARTS rules

Functional groups	÷
All endpoints	
Acute Aquatic Toxicity	
Dummy	
Skin sensitization	
Non-genotoxic carcinogenicity	
Genotoxic carcinogenicity, mutagenicity	
Reactive, unstable, toxic	
Potential electrophilic agents	
Idiosyncratic toxicity (RM formation)	
Custom filters	
Functional groups	
Promiscuity	
Developmental and mitochondrial toxicity	
PAINS compounds	
Biodegradable compounds	
Nonbiodegradable compounds	
His-tag frequent hitters	
AlphaScreenTM frequent hitters	
Chelating agents	

Article: All articles All articles 1988 Ashby 1990 Hermens 1992 Verhaar, H.J.M. 1994 Payne 1994 Barratt 2004 Gerner 2005 Kazius 2005 CheckMol 2005 Kalgutkar 2005 Bailey 2008 Enoch 2008 Benigni 2011 Maybridge 2011 Enamine 2011 "Ontario" filters 2011 ChemDiv 2011 Life_Chemicals 2011 Enoch 2012 Tetko, I.V.

Sushko et al, JCIM, 2012, 52(8):2310-6.

Analysis of virtual libraries

ToxAlerts: Screening results

The compounds that matched any alerts grouped by endpoints, publications and by alerts themselves

ENDPOINTS		x	View records for the fill	tered compounds 🛛 🚫 Tag the 1890 filtered molecules 🛛 💐 Export the screening results
			1 - 15 of 1890	15 items on page 1 of 126 > >>
 Functional groups 		1890 compounds		
PUBLICATIONS				Hydroxy compounds: alcohols of phenois (for Functional groups in 2005 Checkwol)
				Carboxylic acid derivatives (for Eurocional groups in 2005 CheckMol)
O 2005 CheckMol		1880 compounds	~	Carboxylic acid amides (for Functional groups in 2005 CheckMol)
O 2012 Salmina		1890 compounds		Carboxylic acid primary amides (for Functional groups in 2005 CheckMol)
				Aromatic compounds (for Functional groups in 2005 CheckMol)
DETECTED ALERTS			HO	Arenes (for Functional groups in 2005 CheckMol)
Hydroxy compounds: alcohols or phenols	2005 CheckMol	416 compounds		Nonmetals (for Functional groups in 2012 Salmina)
	2005 CheckMol	211 compounds	H ₂ N ⁻ O	Chalcogens (oxygen group) (for Functional groups in 2012 Salmina)
Carboxylic acid derivatives	2005 CheckMol	542 compounds	molecule profile	Priccogens (nitrogen group) (for Functional groups in 2012 Salmina)
	2005 CheckMol	101 compounds		
Carboxylic acid primary amides	2005 CheckMol	14 compounds		MoleculeID: M10446
	2005 CheckMol	993 compounds		Hydroxy compounds: alcohols or phonois (for Eurotional groups in 2005 ChockMol)
	2005 CheckMol	915 compounds		Phenols (for Functional arrours in 2005 CheckMol)
	2012 Salmina	1890 compounds		Nitriles (for Functional groups in 2005 CheckMol)
Chalcogens (oxygen group)	2012 Salmina	1433 compounds		Aromatic compounds (for Functional groups in 2005 CheckMol)
Phictogens (nitrogen group)	2012 Salmina	800 compounds		Arenes (for Functional groups in 2005 CheckMol)
Tetragens (carbon group)	2012 Salmina	1890 compounds	UN UN	Nonmetals (for Functional groups in 2012 Salmina)
O Nitriles	2005 CheckMol	55 compounds		Chalcogens (oxygen group) (for Functional groups in 2012 Salmina)
Carboxylic acid amidines	2005 CheckMol	4 compounds	N	Pricogens (arbogen group) (for Functional groups in 2012 Salmina)
Heterocyclic compounds	2005 CheckMol	309 compounds	molecule profile	
Five-membered heterocycles (LS)	2012 Salmina	143 compounds		MoleculeID: M1364
 Five-membered heterocycles with two heteroatoms (LS) 	2012 Salmina	61 compounds		Catherylic acid amidiaes (for Eurotional groups in 2005 CheckMol)
 Unsaturated five-membered heterocycles with two heteroatoms (LS) 	2012 Salmina	6 compounds		Aromatic compounds for Eurotional groups in 2005 CheckMol)
 Five-membered heterocycles (HS) 	2012 Salmina	69 compounds		Arenes (for Functional groups in 2005 CheckMol)
 Five-membered heterocycles with two heteroatoms (HS) 	2012 Salmina	31 compounds		Heterocyclic compounds (for Functional groups in 2005 CheckMol)
 Unsaturated five-membered heterocycles with two heteroatoms (HS) 	2012 Salmina	5 compounds		Nonmetals (for Functional groups in 2012 Salmina)
 Imidazolines (HS) 	2012 Salmina	2 compounds		Pnictogens (nitrogen group) (for Functional groups in 2012 Salmina)
O 2-Imidazolines (HS)	2012 Salmina	2 compounds		Tetragens (carbon group) (for Functional groups in 2012 Salmina)
Carbonyl compouns: aldehydes or ketones	2005 CheckMol	177 compounds	Ť	Five-membered heterocycles (LS) (for Functional groups in 2012 Salmina)
⊖ Aldehydes	2005 CheckMol	53 compounds	NH	Unsaturated five-membered heterocycles with two heteroatoms (LS) (for Functional groups in 2012 Salmina)
⊖ Thioethers	2005 CheckMol	14 compounds		Five-membered heterocycles (HS) (for Functional groups in 2012 Salmina)
O Dialkylthioethers	2005 CheckMol	5 compounds	molecule profile	Five-membered heterocycles with two heteroatoms (HS) (for Functional groups in 2012 Salmina)
Carbonyl compounds: aldehydes and ketones	2005 CheckMol	170 compounds		Unsaturated five-membered heterocycles with two heteroatoms (HS) (for Functional groups in 2012 Salmina)
⊖ Amines	2005 CheckMol	316 compounds		Imidazolines (HS) (for Functional groups in 2012 Salmina)
O Primary aliphatic amines	2005 CheckMol	59 compounds		2-Imidazolines (HS) (for Functional groups in 2012 Salmina)
Primary amines	2005 CheckMol	180 compounds		Molecula Dr. M24112

An easy fast way to identify possibly problematic compounds Comprehensive grouping of compounds according to a given feature

Functional groups

Online chemical d with modeling env	tabase	v2.4.45 Welcome, Guesti Logout
Home		A+ a-
ToxAlerts: Structural alerts browser Here you can browser structural alerts for vari	us toxicological endpoints	
FILTERS	😝 Upload new alerts 🛛 🔍 Screen compounds 🛛 🗸	
Article:	101 - 200 of 204 << < 100 0 items on page 2 of	3 > >>
Endpoint / Filter type: Functional groups	Thioxohetarenes R = H, alkyl, aryt; any heteroaromatic compound with a C=S structure	
Name / Alert ID:	SMARTS: [\$(c(=[SX1])[\$(n[#1.#6]).o.s]),\$(c(=[SX1])a[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaa[\$(n[#1.#6]).o.s]),\$(c(=[SX1])aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa	1]]aaa[\$(n[#1,#6]],o,s]]]
Show only approved alerts	CheckMol Security List of functional groups recognized by checkmol 2005;	= ©
	Alert ID: 7/4/207	13:32, 5 May 12 / 10:14, 7 Dec 12 SALMINA1987 gg
	N R2 N R1, R2 = H, alkyl, aryl; ; any heteroaromatic compound with a C=N structure SMARTS: [\$(c(=N)[\$(n[#1,#6]).o.s]),\$(c(=N)a[\$(n[#1,#6]).o.s]),\$(c(=N)aa[\$(n[#1,#6]).o.s]),\$(c(=N)aaaa[\$(n[#1,#6]).o.s]),\$(c(=N)aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa	s]),o,s])] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
	Alert ID: 7A7208	13:32, 5 May 12 / 16:14, 7 Dec 12 SALMINA1987 53
	Y Y Y Y R H, alkyl, aryl; Y = OH, alkoxy, aryloxy, (substituted) amino, etc. SMARTS: [CX4&is([CX4]([F,CI,Br,I])([F,CI,Br,I]))([#1,#6&is[(CX3]=[OX1,SX1,NX2])))([#1!#6])([!#1!#6])(6))[#11#6]
	Cathonylia asid atheasters	13.32, 0 M0912/16/14, 21 M0F13 SALMINA1987 88
	R1 = H, alkyl, aryl R1 = H, alkyl, aryl SMARTS: [#1,#6][CX4][(0X2][#6&I\$([CX3]=[0X1,SX1,NX2,C])])[[0X2][#6&I\$([CX3]=[0X1,SX1,NX2,C])])[0X2] Endpoint: Functional groups CheckMol Substrational groups recognized by checkmol 2005; Alert ID: 7A1270	[#6&I\$([CX3]=[OX1,SX1,NX2,C])]
	Carboxylic acid amide acetals R1, R2, R3 = H, alkyl, aryl; R4, R5 = alkyl, aryl	
	R4O OR5 SMARTS: [CH1,\$([CX4][#6])]([0X2][#6&I\$([CX3]=[0X1,SX1,NX2,C])])([0X2]]#6&I\$([CX3]=[0X1,SX1,NX2,C])] Endpoint: Functional groups	[(NX3)([#1,*])(#1,*] □

Screening Artefacts



Pan-Assay Interference Compounds (PAINS) filters by Baell and Holloway, 2010.

HIS tag frequent hitters, Schorpp, K. et al *J. Biomol. Screen.* 2014, 19(5):715-726. GST tag frequent hitters, Brenke, J.K. et al *J. Biomol. Screen.* 2016, 21(6):596-607.

MOA of AlphaScreenTM-HIS-FH



Schoorp et al, J Biomol Screen. 2014, 19(5):715-726.

SetCompare tool for data analysis

6	Onlin	e cher	nical data	abase ment	•			
Home -	Database -	Models -	Moderation -			1		
Ex	Welcome plore OC: rch chemical a sured publist	Create a Apply a r Create n Open pr	model model nultiple models redictor		; /			
our	users. You ca	Calculate SetCom	e descriptors	10003	SetComp	pare: Select the sets to compare	more two sets of molecules	hased on their structural features
Build	are QSAR mode erties. The me erimental data	els for pred odels can t published	ictions of chem be based on the in our database	ical e e.	Please, provi	the compounds in the first set:		2 Select the compounds in the second set:
					 (SDF/MC sheet) Provide a RN/SMIL 	a Name/CAS-	Durchsuchen)	(SDF/MOL2/SMILES/Excel Durchsuchen) (SoPi/MOL2/SMILES/Excel Durchsuchen) Provide a Name/CAS- RN/SMILES
					 Choose a set: 	a previously prepared		Chain indecute Cick on depiction to the right to draw) Choose a previously prepared set: []
					⊖ Select m Next >>	olecules by a tag: []		○ Select molecules by a tag: []

Examples of scaffolds analysis

SetCompare: Comparison results The comparison summary of the two selected sets

The following table shows the features (molecular descriptors) that were significantly overrepresented in one of the tw It includes appearance counts of the features in each set and the p-Value of such a distribution. Export results as a CSV file

1 - 15 of 253				1	5 ᅌ ite
Descriptor	In set 1 (13785 molecules)	In set 2 (228174 molecules)	Enrichment factor	p-Value	
R	3232 (23.4%)	26196 (11.5%)	2.0	1.33E-315	
Pnictogens Group 15: the nitrogen family N P As Sb Bi	13235 (96.0%)	202449 (88.7%)	1.1	1.42E-198	
R1 R2	3257 (23.6%)	79741 (34.9%)	1.5	-1.25E-172	
R - OH OH	280 (2.0%)	352 (0.2%)	13.2	4.21E-172	
R-NH ₂	2063 (15.0%)	18761 (8.2%)	1.8	2.23E-140	

Pyrolysis vs. melting point

SetCompare: Comparison results

The comparison summary of the two selected sets

The following table shows the features (molecular descriptors) that were significantly overrepresented in one It includes appearance counts of the features in each set and the p-Value of such a distribution. Export results as a CSV file



HIV Envelope glycoprotein GP120

Application of models



Modeling framework

Creation of QSAR/QSPR models

- It is easy to create new models from uploaded data
- More than 10 Machine learning algorithms
 - Neural networks, KNN, SVM, MLR, PLS, random forests, J48
- More than 20 descriptor packages
 - 0D to 3D descriptors, academic and commercial
 - New descriptor packages can be easily integrated
 - Packages can be easily combined (mix & match)
 - Dragon, CDK, ADRIANA.CODE, Chemaxon, E-state, MERA, ...

Models can be updated without complete rebuilt of model Estimation of accuracy for each prediction

Modeling iterative workflow

Select dataset

Over >1M measured values
Over 400 properties

Validate

Internal (N-Fold crossvalidation, Bagging) External validation Select descriptors (24 packages:0D, 1D, 2D 3D)

Build model

(MLR, ANN, KNN, Random Forest, SVM, FSMLR, WEKA-J48)

LogP of Pt(II) + Pt(IV) complexes



Tetko et al., J. Inorg. Biochem., 2016, 156(3), 1-13.

Model statistics

Predicted property: logPow Training method: Consensus

Data Set	#	R2	q2	RMSE	MAE
• Training set: LogPt final	187 records	0.93 ± 0.01	0.93 ± 0.01	0.41 ± 0.03	0.3 ± 0.02
• Test set: Hristo 💲 [x]	14 records	0.76 ± 0.1	0.76 ± 0.1	0.5 ± 0.06	0.43 ± 0.07



- Multiple statistical measures + confidence intervals
- Export for offline inspection
- Interactive scatter plot

Accuracy of prediction



Х

Outlying point on the Applicability Domain plot



- Estimation of applicability domain of models
- Identification of outliers

Accuracy of predictions for classification model

verview Applicability odel name: Ames levenb enchmarking of distance edicted property: AMES aining method: ANN	domain erg , publis to models fo	hed in or Ame	n Applicability es mutagenic	domains for ity set. publi	r classifica c identifie	ation problem r is 1	ms: Corre er	[OEstate] I. limit: 0.95 Variance threshold: 0.0, Maximum value: 999999, Levenberg, 1000 iterations, 3 neurons nsemble=100 additional param PARALLEL=10 5-fold cross-validation
Data Set			#	Accuracy	Balanced	d accuracy		
o Training set: Ames cha	allenge train	ning 4	4357 records (4359 selected)	78.1 ± 1.2	77.9	± 1.3		Calculated in 2402 seconds Size: 450 Kb
o Test set: Ames challer	nge test [x]	2	2181 records	79.9 ± 1.7	79.8	± 1.7		0.201 700 70
Real।/Predicted→	inactive	active	e	Real!/Predi	icted→	inactive	active	
inactive	1521	495		inactiv	e	802	207	
active	460	1883	3	active 232		940		
Training (C	Training (Original)			Test (Original)				

Overview Applicability domain



Solubility in DMSO classification



Example of a published solubility in DMSO model

✓ based on >163k compounds (provided by UCB and Enamine)
 ✓ nine folds decrease number of non-soluble in DMSO compounds

Tetko et al J. Chem. Inf. Model. 2013, 53(8):1990-2000.

300k Melting Point Datasets



Comprehensive modeling

Package name	Type of descriptors	Number of descriptors	Matrix size, billions	Non zero values, millions	Sparseness
Functional Groups	integer	595	0.18	3.1	33
QNPR	integer	1502	0.45	6.3	49
MolPrint	binary	688634	205	8.1	7200
Estate count	float	631	0.19	10	14
Inductive	float	54	0.02	11	1
ECFP4	binary	1024	0.31	12	25
Isida	integer	5886	1.75	18	37
ChemAxon	float	498	0.15	23	1.5
GSFrag	integer	1138	0.34	24	5.7
CDK	float	239	0.07	27	2
Adriana	float	200	0.06	32	1.3
Mera, Mersy	float	571	0.17	61	1.1
Dragon	float	1647	0.49	183	1.5

Comprehensive modeling

Predicted property: Melting Point

Training set: patents - decomposition (4 different versions detected)

Metrics RMSE - Root Mean Square Error C Training set	Validation: All validation protocols
	LibSVM
GSFrag	41.85
ChemaxonDescriptors (7.4)	39.63
InductiveDescriptors	49.54
MolPrint (Length 2)	61.2
StructuralAlerts	41.92
ECFP (4_1024)	61.46
OEstate	37.8
Fragmentor (Length 2 - 4)	38.54
QNPR (SMILES - length 1 - 3 threshold 5)	39.75
CDK (constitutional, topological, geometrical, electronic, hybrid)	38.45
Mera, Mersy	42.8
Dragon (blocks: 1-20)	43.04
Adriana	41.91
SIRMS (LABELING=ELM noH type=EXTENDED)	40.36
MolPrint (Length 3)	60.89
Consensus	
Misc. 36.01	

Outliers identified with applicability domain (AD) plot



Functional group analysis of pyrolysis and MP data

Descriptor	In set 1 (13785 molecules)	In set 2 (228174 molecules)	Enrichment factor	p-Value
R	3232 (23.4%)	26196 (11.5%)	2.0	1.33E-315
Pnictogens Group 15: the nitrogen family N P As Sb Bi	13235 (96.0%)	202449 (88.7%)	1.1	1.42E-198
R ₁ R ₂	3257 (23.6%)	79741 (34.9%)	1.5	-1.25E-172
	280 (2.0%)	352 (0.2%)	13.2	4.21E-172

Outliers have higher percentage of predicted decomposing compounds

20% of ~5k outliers are predicted as pyrolysis

VS.

17% of ~223k PATENTS compounds

Why do we need the Melting Point (MP)?

Use of MP for water solubility prediction

Yalkowsky equation:

logS = 0.5 - 0.01 (MP-25) - log Kow

Prediction of Huuskonen set using ALOGPS logP and MP based on 230k measurements

Predicted property: Aqueous Solubility modeled in log(mol/L) Training method: MLRA

Data Set	#	R2	q2	RMSE	MAE
• Training set: logS set	1311 records	0.842 ± 0.009	0.83 ± 0.01	0.84 ± 0.02	0.64 ± 0.02



Best models based on this set have RMSE ~ 0.6

Online chemical database Modeling of mixture deling environment



Mixtures' descriptors



Interpretation of models

Molecular Matched Pairs

A molecular matched pair (MMP) is a pair of molecules that have only a (minor) single-point difference. The typical way is to define a minor difference as a changed molecular fragment with less than 10 atoms.



MMPs for classification data (AMES mutagenicity)



Analysis of rules that were learnt by models



LogP of Pt(II) + Pt(IV) complexes



Model for LogP for Pt complexes

Predicted property: **logPow** Training method: Consensus

Data Set	#	R2	q 2	RMSE	MAE
• Training set: LogPt final	187 records	0.93 ± 0.01	0.93 ± 0.01	0.41 ± 0.03	0.3 ± 0.02
• Test set: Hristo 💲 [x]	14 records	0.76 ± 0.1	0.76 ± 0.1	0.5 ± 0.06	0.43 ± 0.07



Model for LogP for Pt complexes

Outlying point on the Applicability Domain plot

Prediction of logP for Pt complexes

Outliers of the model (functional groups descriptors)

Outliers of the model (Functional groups descriptors)

Storage of Chemical Reaction as Condensed Graph of Reaction (CRG)

Modeling of SN2 reactions using OCHEM

Predicted property: **SN2 reaction rate constant** modeled in log(L*mol^(-1)*s^(-1)) Training method: Consensus

Data Set	#	R2	q2	RMSE	MAE
• Training set: SN2 methanol 20-40	125 records	0.54 ± 0.08	0.54 ± 0.08	0.64 ± 0.08	0.44 ± 0.04

Consistent data: t = 20-40, solvent = methanol

based on 125 records for **SN2 reaction rate constant** (excluding 20 not indexed molecules) air similarity: Any +

SN2 vs. SN1 reactions

Evaluation in challenges

€FPA United States Environmental Protection Agency ALL EPA 💽 THIS AREA Advanced Search LEARN THE ISSUES | SCIENCE & TECHNOLOGY | LAWS & REGULATIONS | ABOUT EPA SEARCH Computational Toxicology Research Contact Us You are here: EPA Home » Research & Development » CompTox » Chemical Data Challenges & Release Key Links CompTox Home Research Publications Staff Profiles **Research Projects Basic Information** Chemical Databases Scientific Reviews CompTox Partners Organization ToxCast Stakeholder Events Communities of Practice Jobs and Opportunities EPA Exposure Research EPA Chemical Safety Research ToxCast Data Challenges

ToxCast Chemical Data Challenges and Release

EPA's high-throughput screening data on 1,800 chemicals is accessible through the interactive Chemical Safety for Sustainability Dashboards (iCSS dashboard). The iCSS dashboard provides user-friendly and customizable access to toxicity data from ToxCast and Tox21 high-throughput chemical screening technologies.

Using the **TopCoder** and **InnoCentive** crowd-sourcing platform, EPA invited the science and technology community to work with the data and provide solutions for how the new toxicity data can be used to predict potential health effects. The ToxCast data challenges focused on using this data and other publicly available data to predict the lowest effect level from traditional toxicity studies using laboratory animals. Challenge winners received awards for solving this challenge.

Key Links

- Lowest Effect Level Challenge Results (PDF, 497KB, 18pp)
- Chemical Safety for Sustainability Dashboards
- Complete ToxCast Phase II Data & Files
- TopCoder Challenge
- InnoCentive Challenge
- Stakeholder Workshops

National Center for Advancing Translational Sciences

Tox21 Data Challenge 2014

Contact Us

» Home

About Registration Data/Resources

Submissions

Discussion

Leaderboard

Survey NEW

About the Data 🚯

The Challenge

The 2014 Tox21 data challenge is designed to help scientists understand the potential of the chemicals and compounds being tested through the Toxicology in the 21st Century initiative to disrupt biological pathways in ways that may result in toxic effects.

The goal of the challenge is to "crowdsource"

All challenge winners will receive the opportunity to submit a paper for publication in a special thematic issue of Frontiers in Environmental Science

and recognition on the NCATS website and via social media.

Modeling capabilities

 Top-I rank submission model (May 2014) - entry by Sergii Novotarskyi*

 Two Top-I rank individual sub-challenges and overall best balanced accuracy for all targets (January 2015) – entry by Ahmed Abdelaziz**

> *Novotarskyi, S. et al. Chem. Res. Toxicol. 2016, 29, 768-75. **Abdelaziz, A. et al. Front. Environ. Sci. 2016, 4, 2.

Conclusions

• On-line tools are important for dissemination

- OCHEM
 - Can used for chemical data searching
 - ToxAlerts for understanding of hemical data
 - SetCompare can be useful for comparison of sets of molecules
 - Has a powerful modeling framework
 - Easily handles millions of compounds
 - Contribute highly predictive models
 - MMP analysis can complement modeling efforts

Acknowledgements

Dr. Y. Sushko Dr. S. Novotarskyi Mr. R. Körner Mr. A. Abdelaziz Mrs. E. Salmina

Prof. M. Sattler (HMGU) Prof. G. Wess (HMGU)

Dr. K. Hadian (TOX, HMGU) Dr. A. Williams (USA) Dr. D. Lowe (UK) Dr. H. Varbanov (Switzerland) Prof. N. Haider (Austria)

Thank you for your attention!