

Machine-Learning Methods in Property Predictions: *Quo Vadis?*

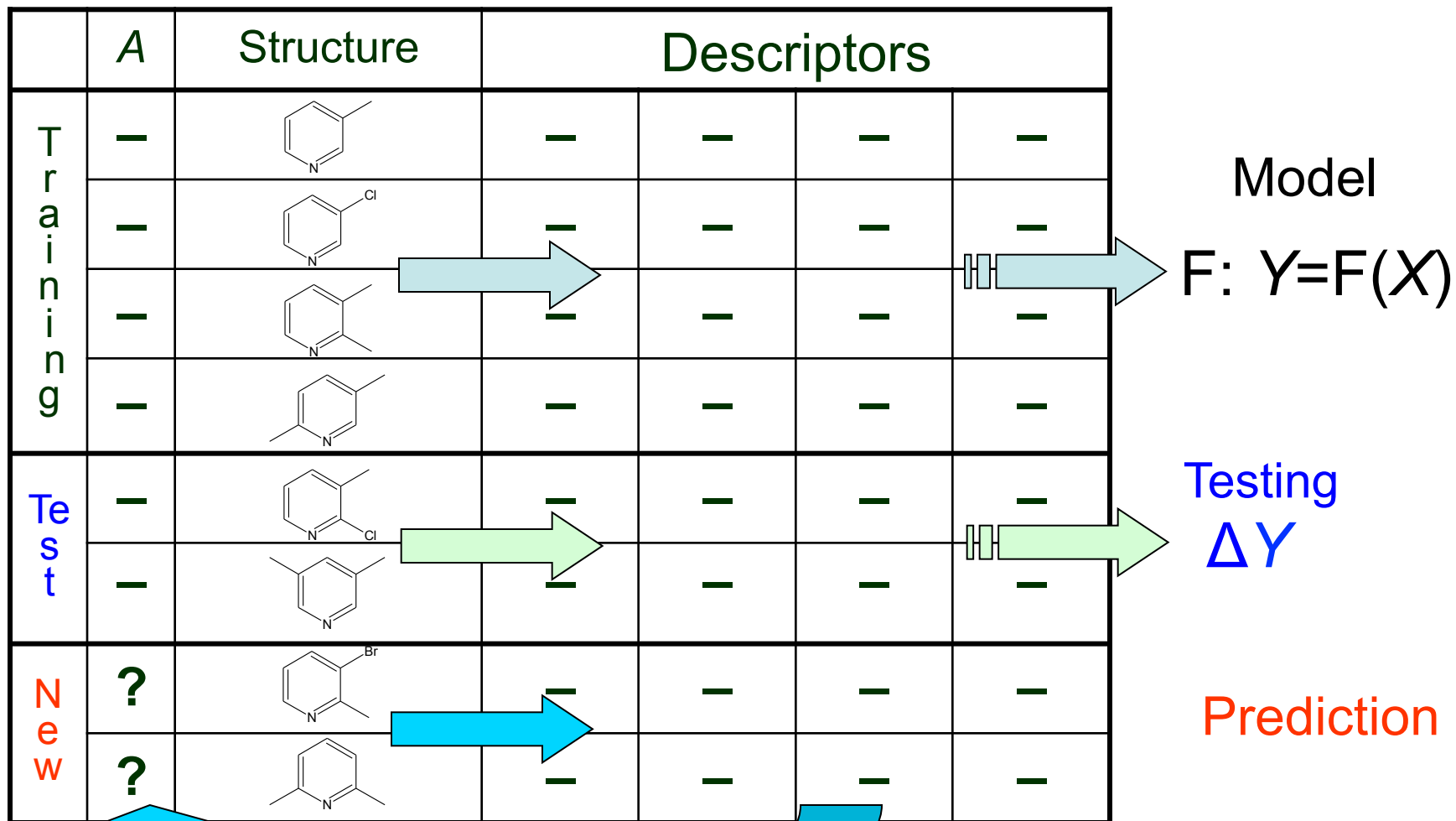
Igor I. Baskin

Lomonosov Moscow State University
RUSSIA

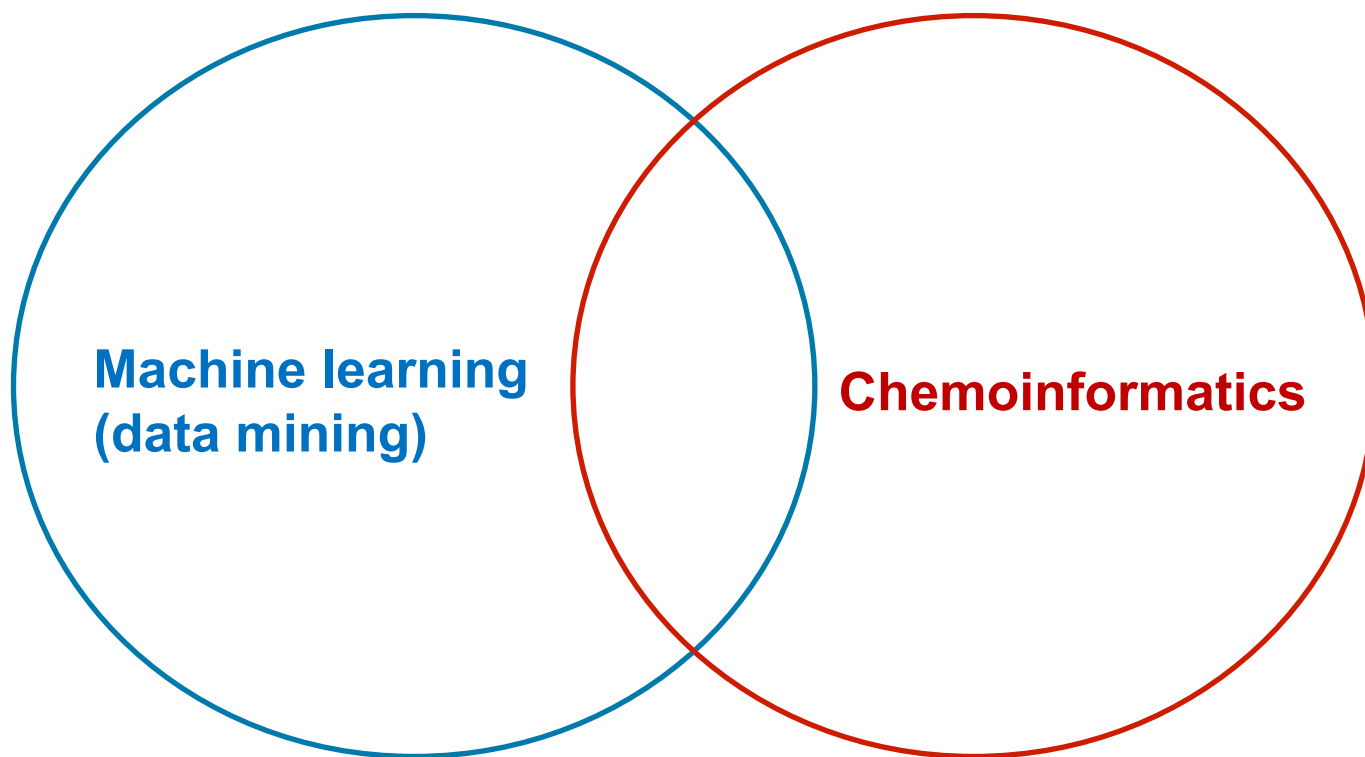


Физический факультет
Московский государственный университет имени М.В. Ломоносова

General Workflow for QSAR Modiling in Chemoinformatics



Machine Learning and Chemoinformatics: **different but overlapping fields**



Chemometrics

- Chemometrics is what chemometricians do.
- Chemometricians are people who drink beer and steal ideas from statisticians

Svante Wold



Chemoinformatics

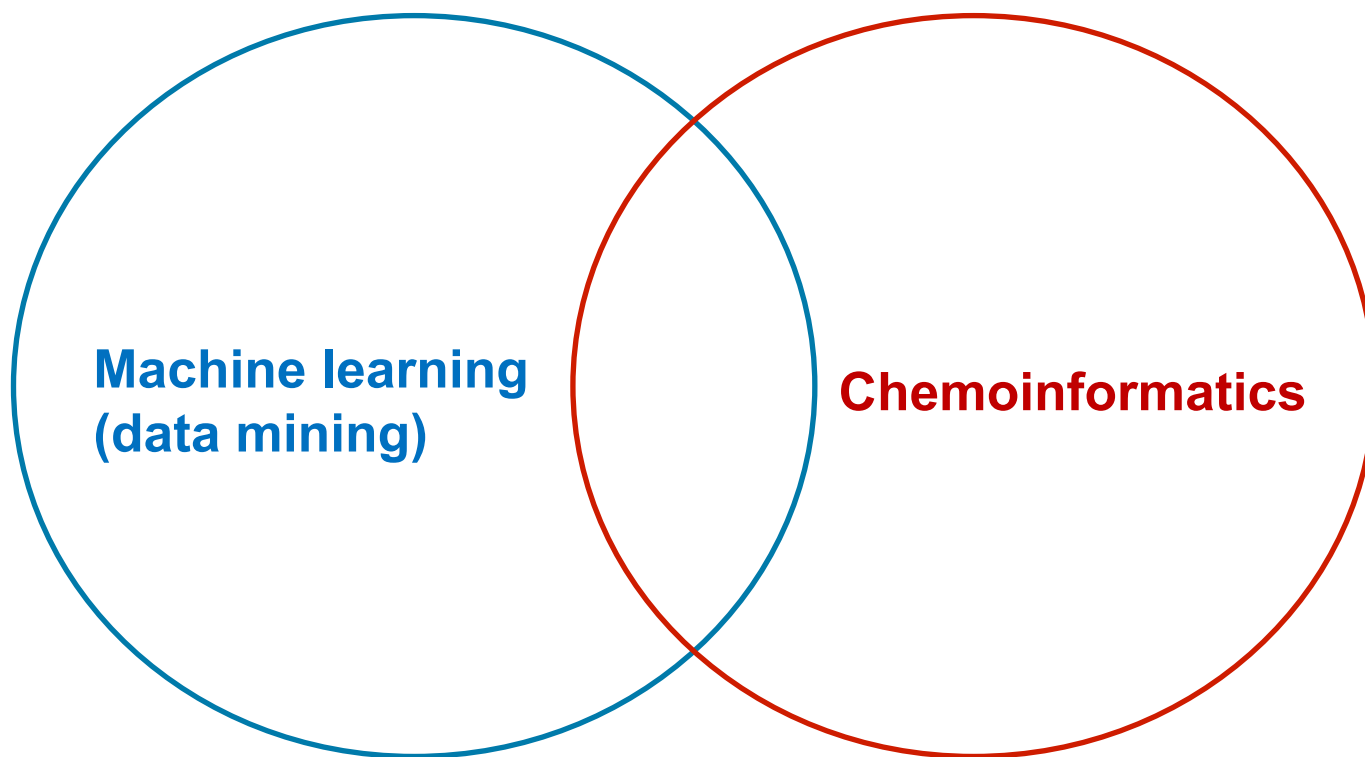
~~Chemometrics~~

- Chemoinformatics chemoinformaticians
- ~~Chemometrics~~ is what ~~chemometricians~~ do

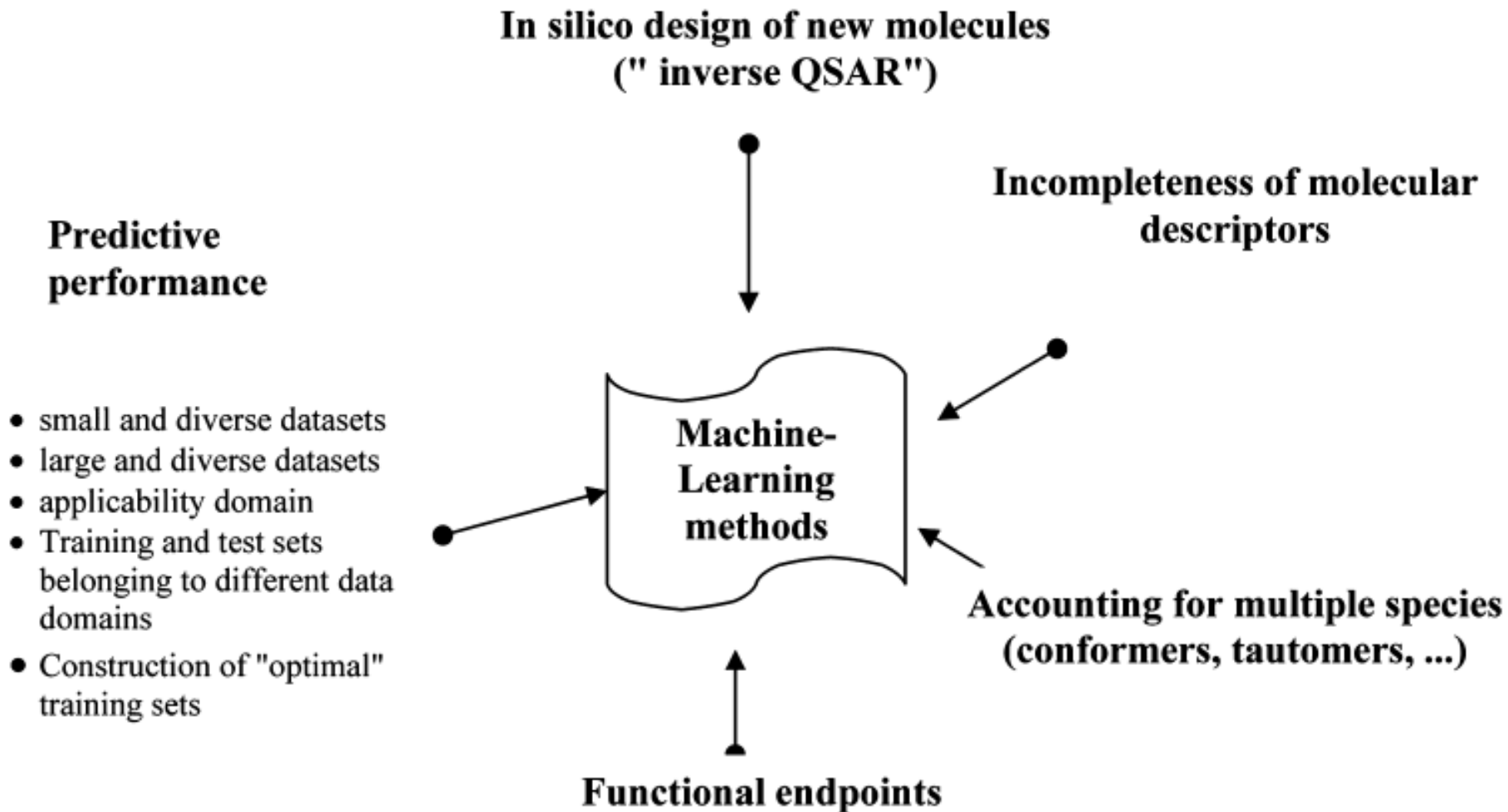
- Chemoinformaticians
- ~~Chemometricians~~ are people who drink beer (??)
 borrow machine-learners
and ~~steal~~ ideas from ~~statisticians~~



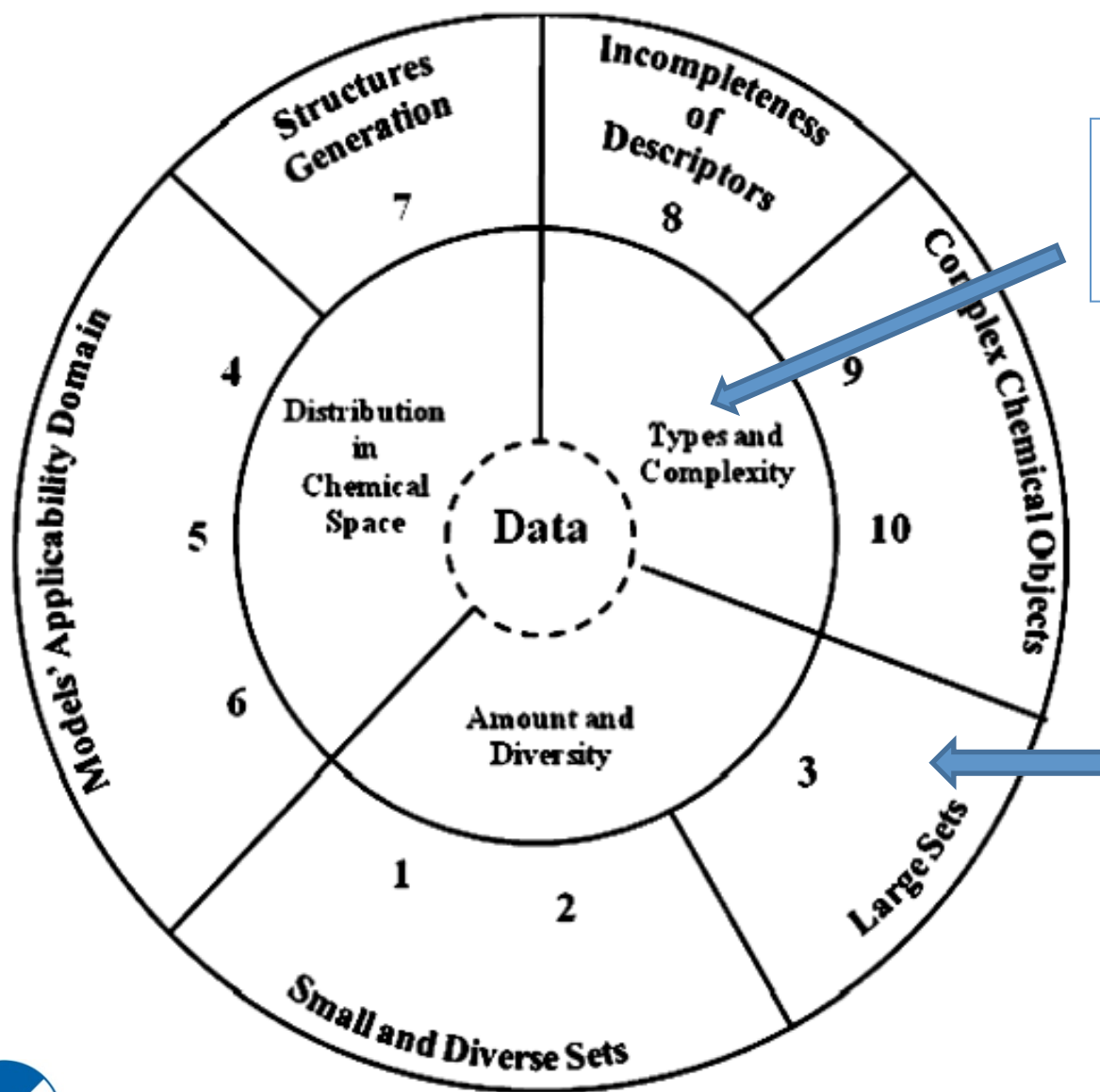
Machine Learning and Chemoinformatics: **different but overlapping fields**



Main Challenges of Machine-Learning Methods in Chemoinformatics



Guide to Choose Machine Learning Method to solve Chemical Problems

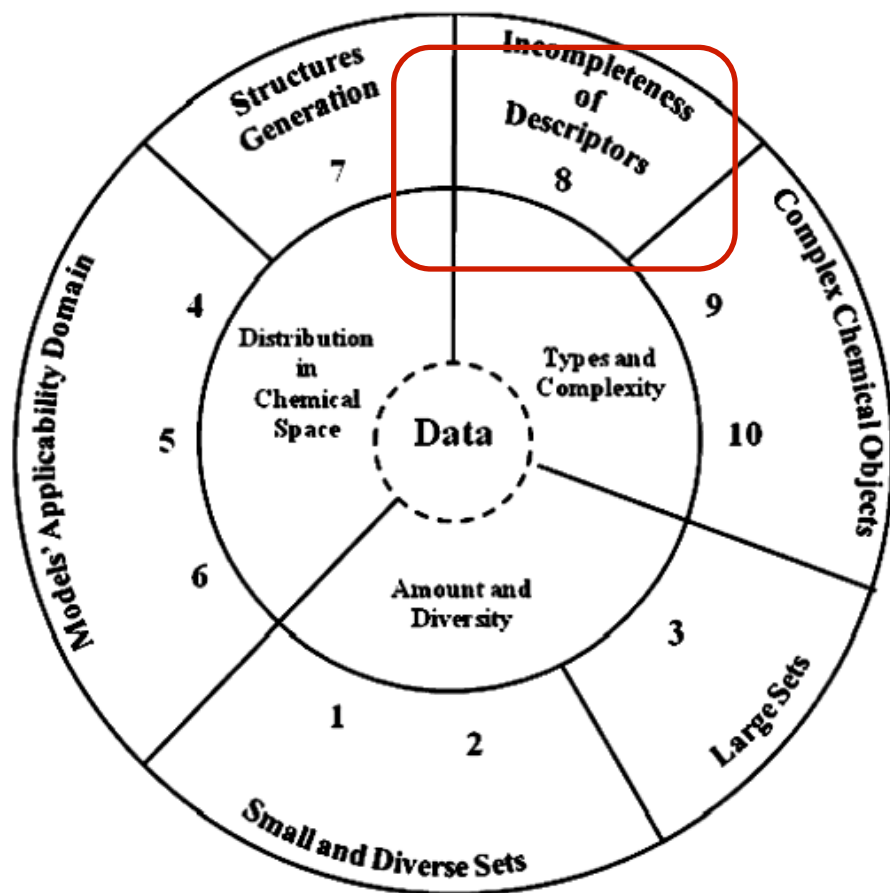


Different features of the data
(*inner circle*)

Challenges of
chemoinformatics
(*outer circle*)



Machine Learning on Molecular Graphs



Is it possible to build a model directly on molecular graphs instead of using fixed-sized vectors of descriptors?

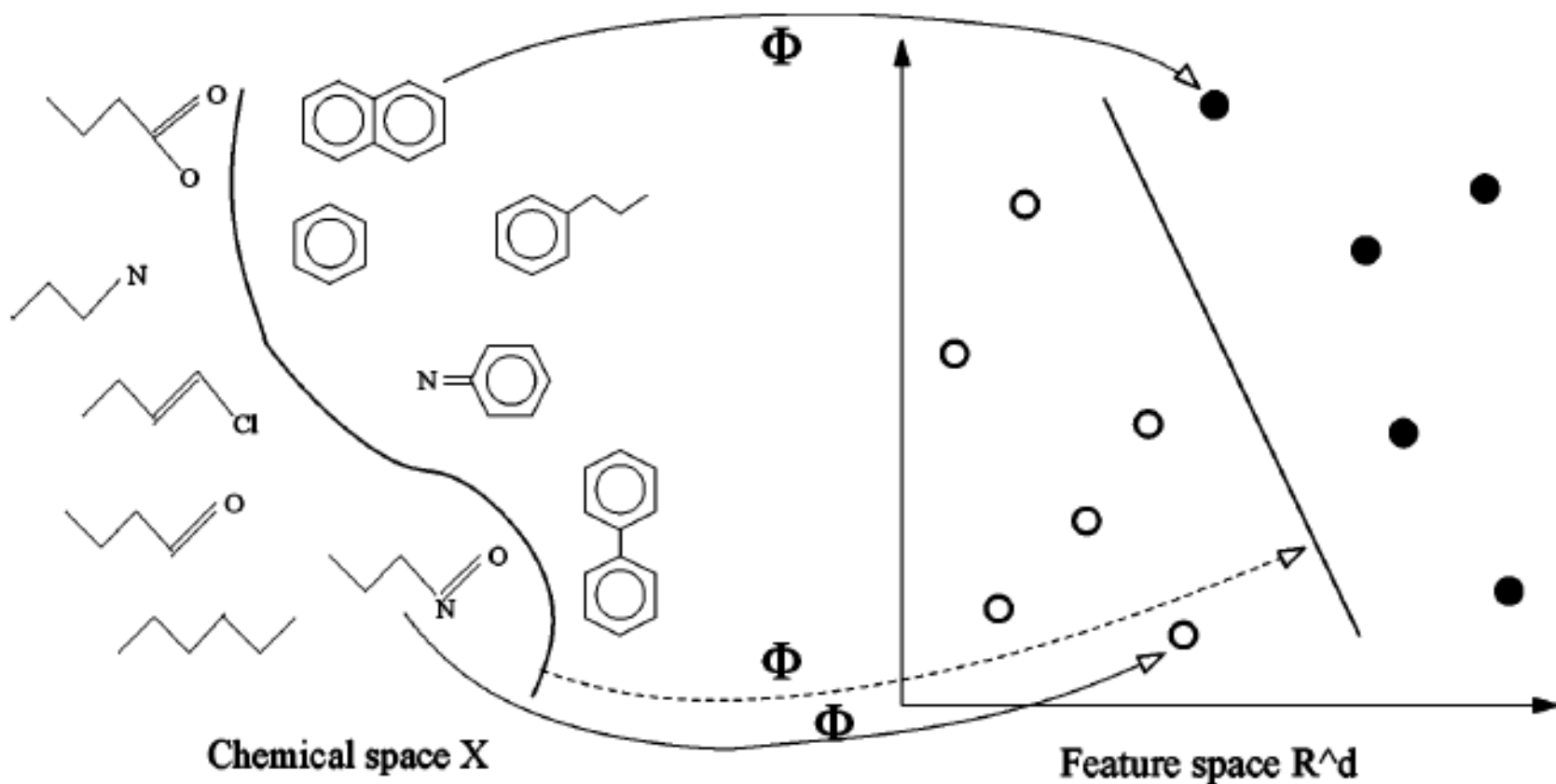
Graph \longrightarrow **Model** \longrightarrow Property

- Graph mining with special architectures of neural networks
- (Sub)Graph mining
- Graph kernels
- Inductive learning programming
- Symmetry-invariant machine learning with local features
- Energy-based learning
- etc

- G.Bakir, T.Hofmann, B.Schoelkopf, A.J.Smola, B.Taskar, S.V.N.Vishwanathan. Predicting Structured Data; The MIT Press:Cambridge, MA, 2007.
- D.J.Cook, L.B.Holder. Mining Graph Data; Wiley-Interscience: Hoboken, NJ, 2007.



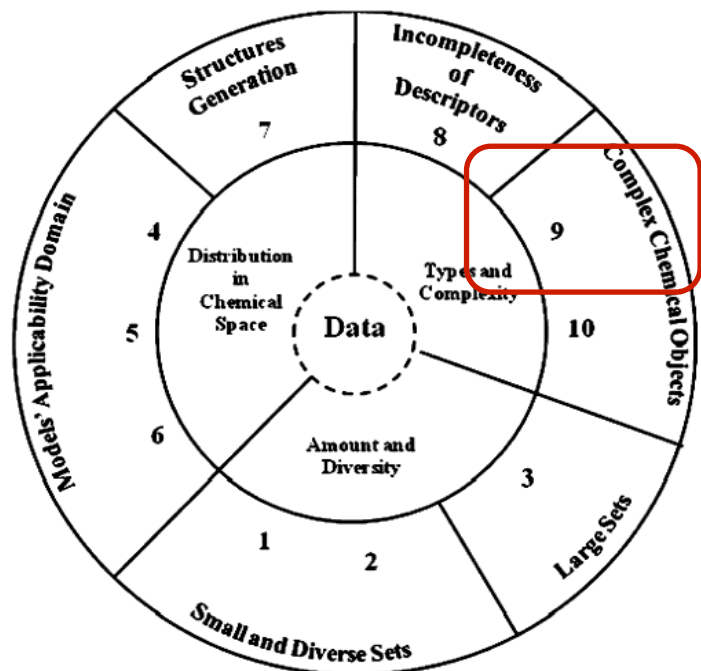
Machine Learning on Graph Kernels



$$\langle \Phi(x), \Phi(x') \rangle = K(x, x')$$

- M.Rupp, G.Schneider. *Mol. Inf.* **2010**, 29 (4), 266–273

Multi-Instance Learning

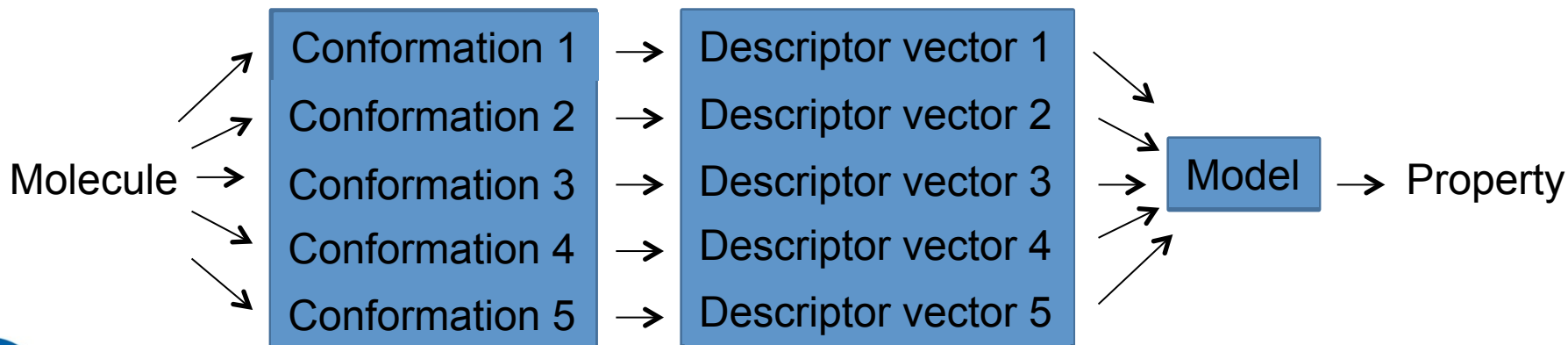


Representing molecule as a number of conformers, tautomers and ionization forms, ...

Every object represents an ensemble (so-called bag) of instances, each of which is described by a fixed-sized vector of descriptors.

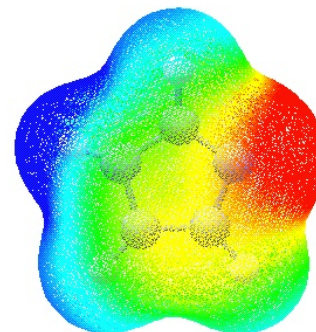
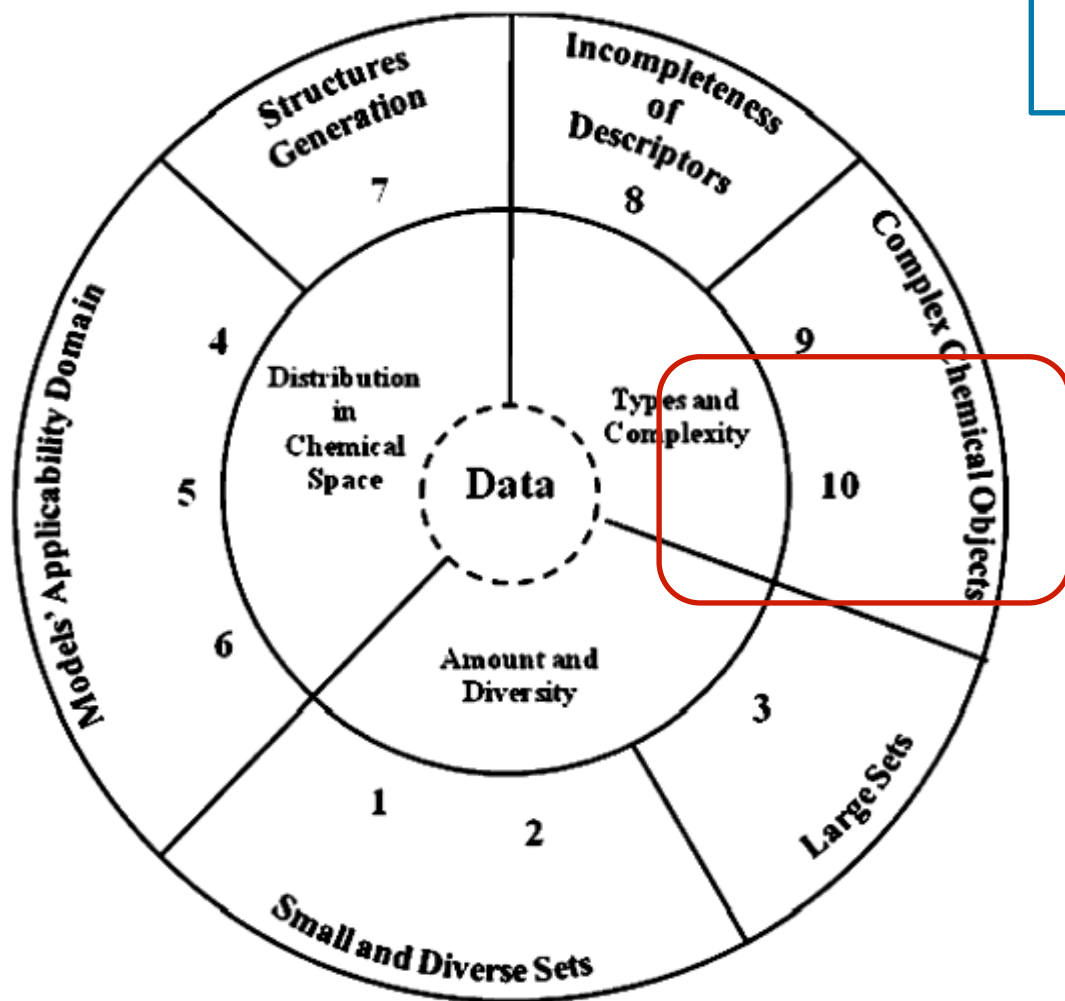
Instances
(conformations,
tautomers, etc)

Bag of feature
vectors (descriptor
vectors)

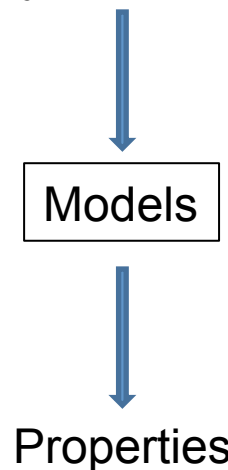


Functional Data Analysis

FDA allows one to build models for molecules represented by functions?



Objects represented by functions



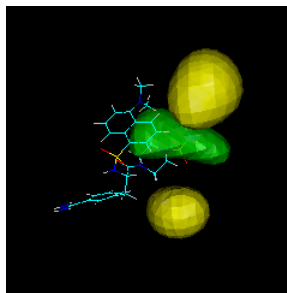
Continuous Molecular Fields (CMF)

Continuous Molecular Fields approach describes molecules by ensemble of continuous functions (*molecular fields*), instead of finite sets of molecular descriptors. CMF is kernel-based method.

traditional QSAR $Activity = F(X) = \sum c_i x_i$

CMF $Activity = F[X(\mathbf{r})] = \int C(\mathbf{r})X(\mathbf{r})d\mathbf{r}$

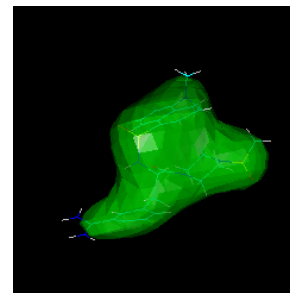
$$Activity = \int$$



C(r)

Calculated using special kernels
for molecular fields

.



X(r)

Gaussian functions approximation
of molecular fields

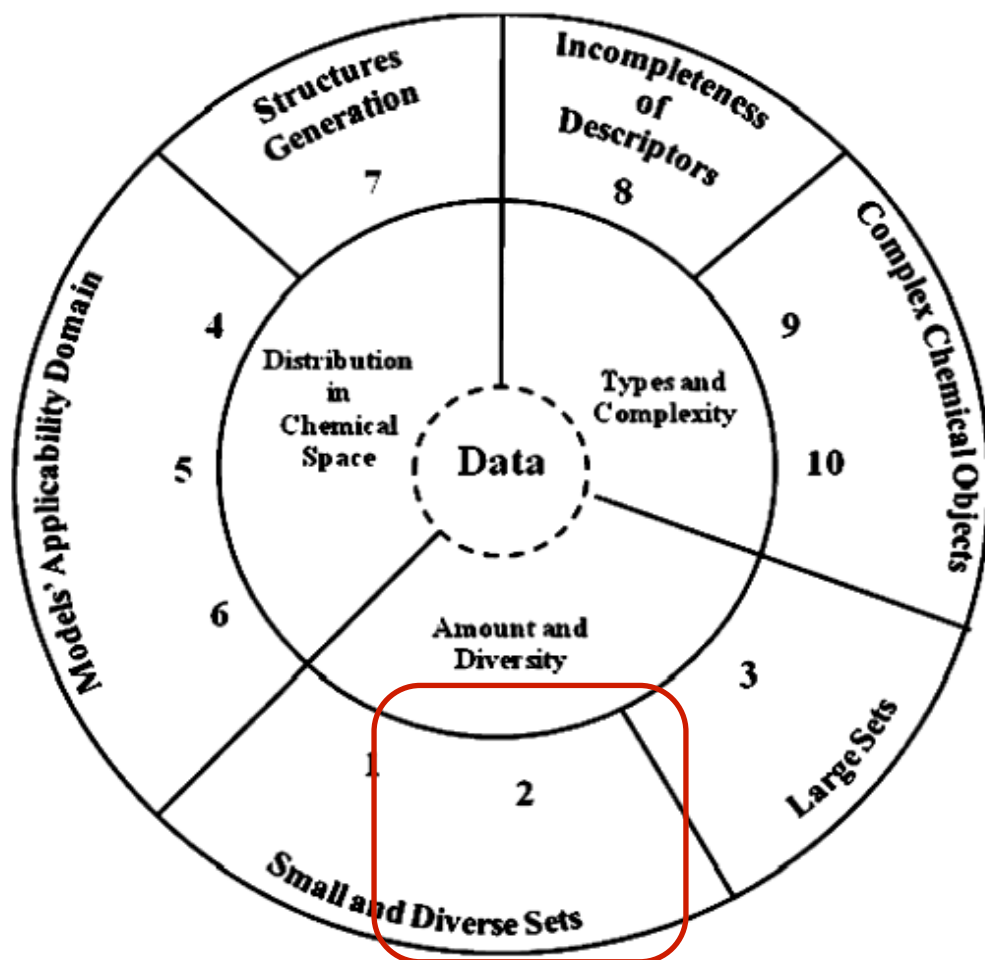
d \mathbf{r}

<http://sites.google.com/site/conmolfields/>



Inductive Knowledge Transfer

(inductive bias, lifelong learning, learning to learn, collaborative filtering, multi-task learning etc)

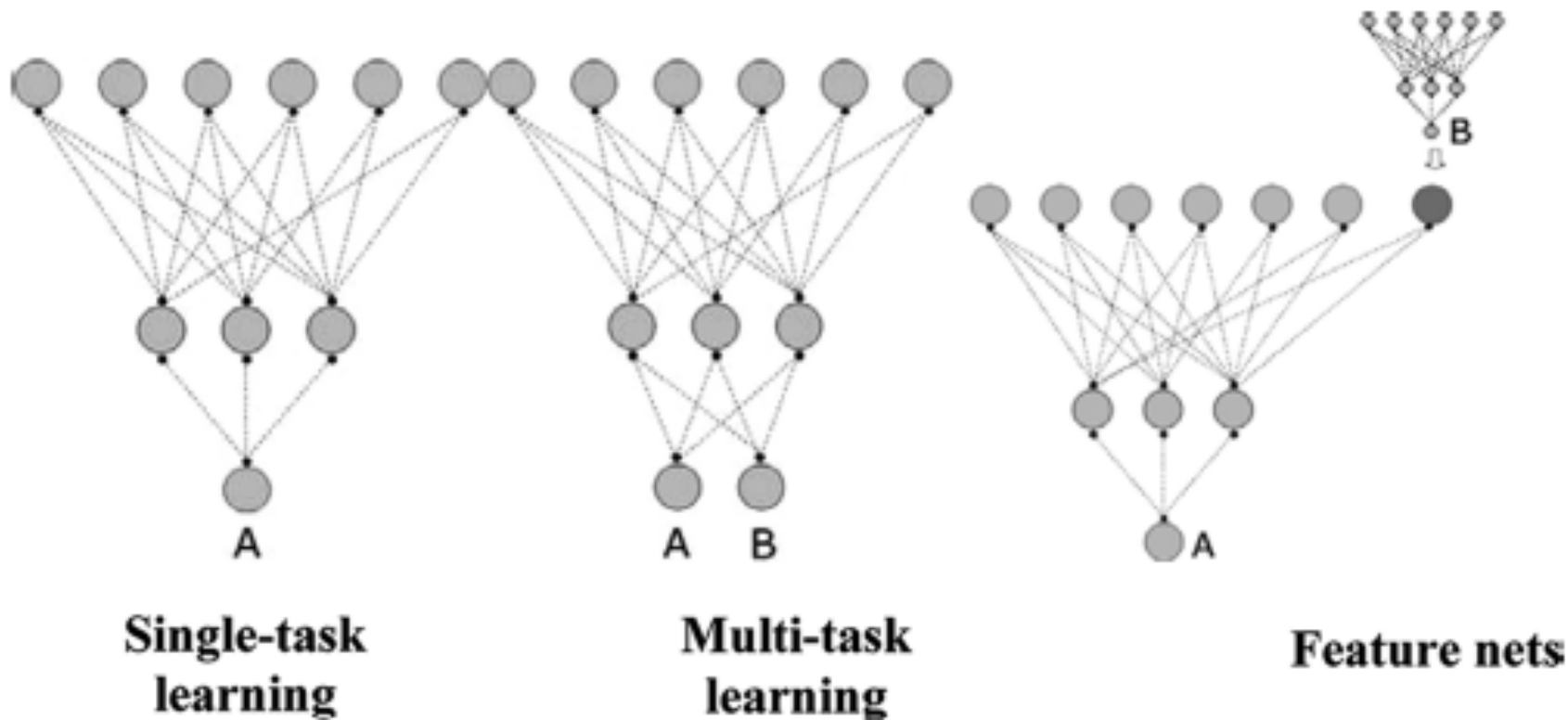


Transfer of information from one model, usually trained on sufficiently large dataset, to another model trained on small dataset

• *Learning to Learn*; S. Thrun, L.Y. Pratt, Eds.; Kluwer Academic Publishers: Boston, MA, 1998



Interference of Models (Inductive Knowledge Transfer)

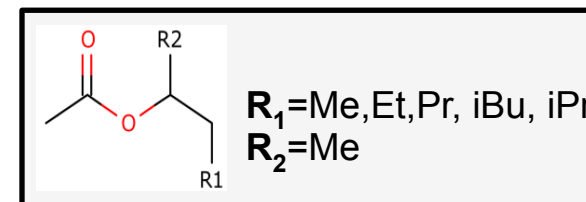
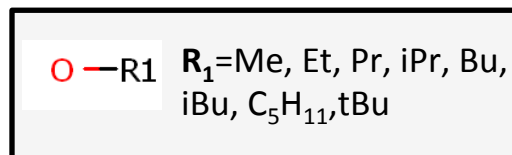
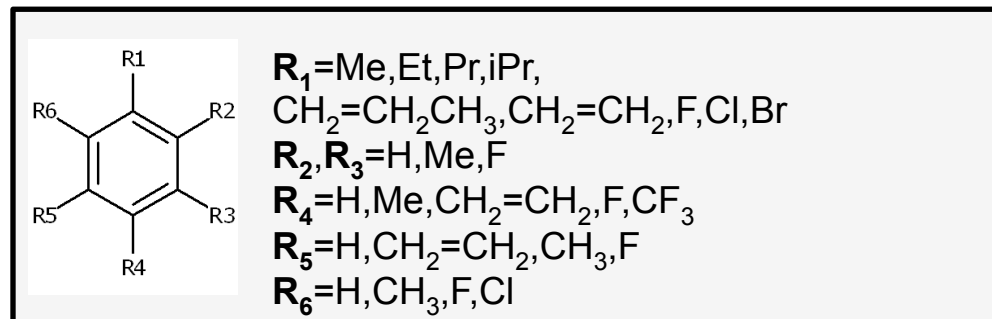


Partition coefficients air-tissue

The blood:air partition coefficient (PC) is an important determinant of the distribution of volatile organic chemicals (VOCs).

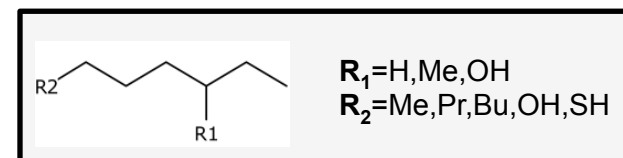
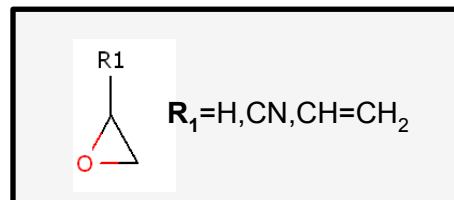
Human

Tissue	N
blood	139
fat	42
brain	36
liver	34
muscle	39
kidney	34

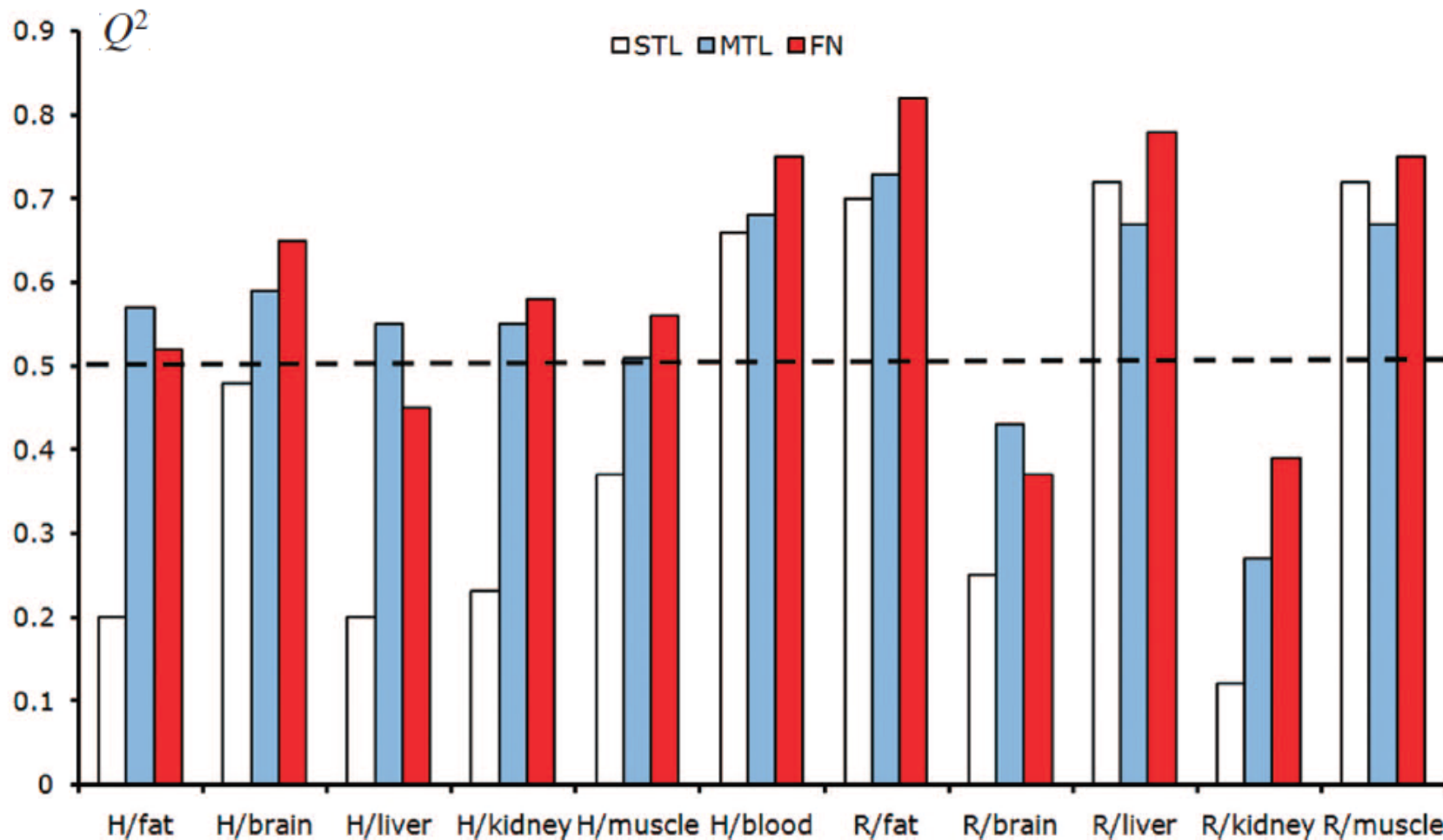


Rat

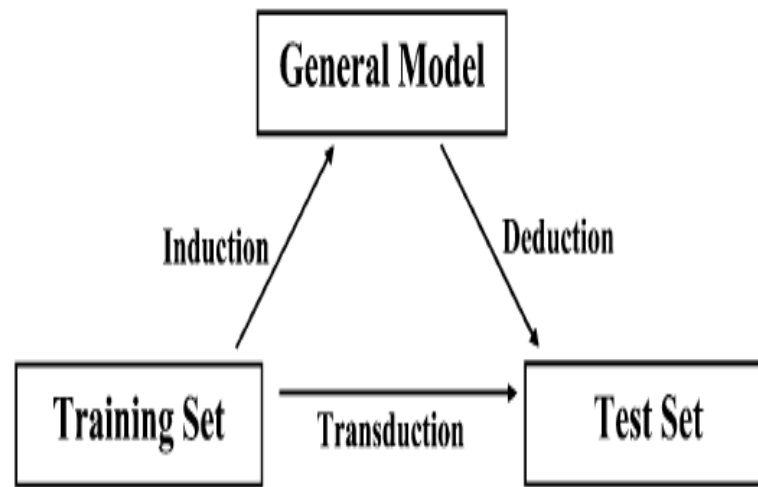
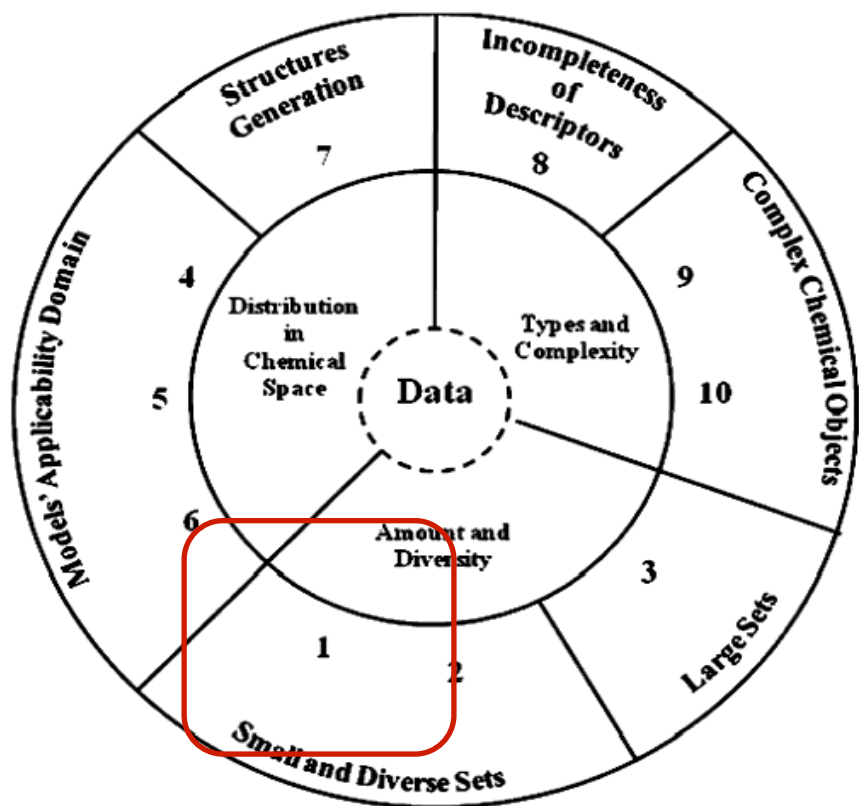
fat	99
brain	59
liver	100
muscle	97
kidney	27



Inductive Knowledge Transfer (Modeling Tissue-Air Partition Coefficients)



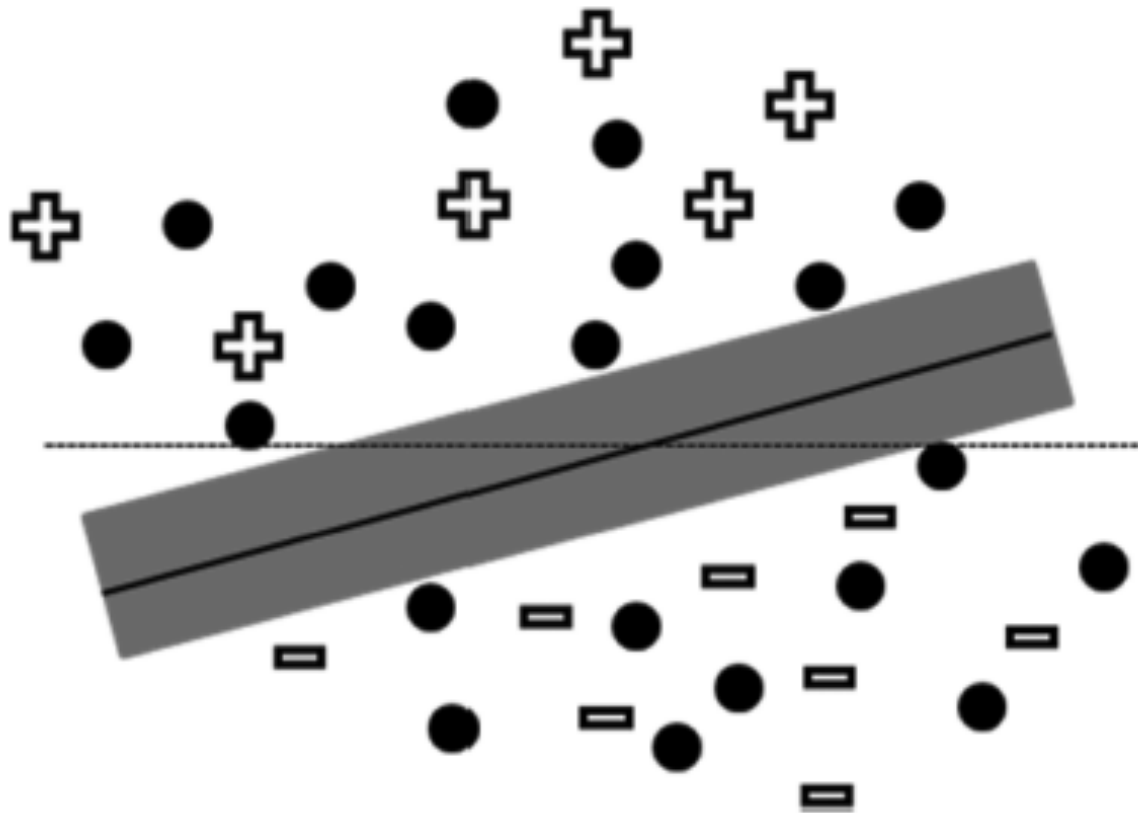
Transductive (Semi-Supervised) Machine Learning



Transductive modeling is used to build the models specifically oriented toward the best prediction performance on a particular test set instead of developing general models to be applied to any test set



Object Separation in SVM and TSVM

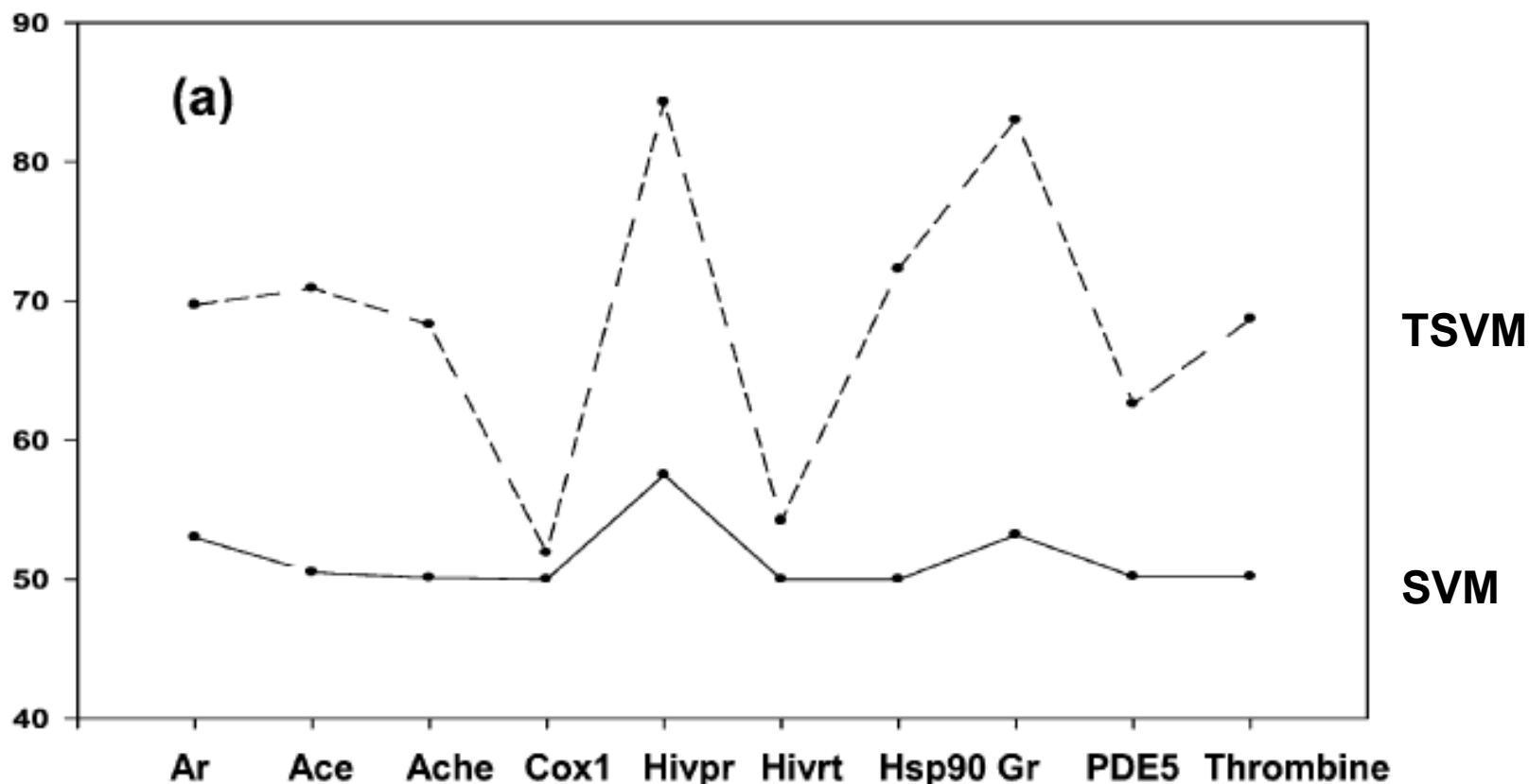


Labeled training set examples are depicted as signs - and +,. Unlabeled test set examples are shown as bold dots.



Prediction Performance (Balanced Accuracy) of SVM vs TSVM Models

(Training sets consist of 5 active and 50 inactive compounds)

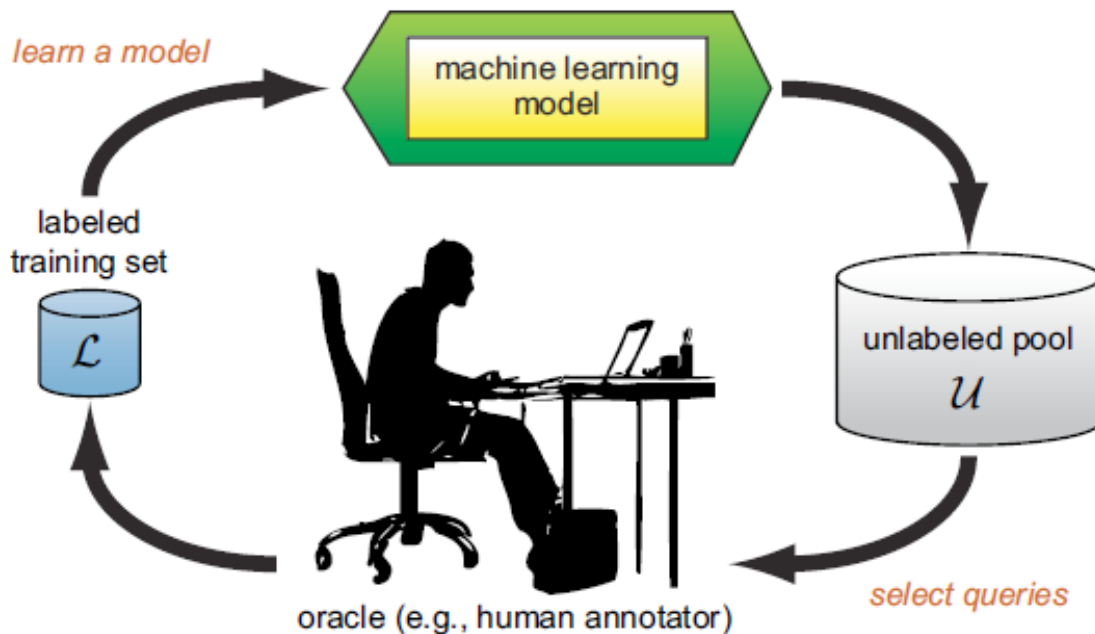
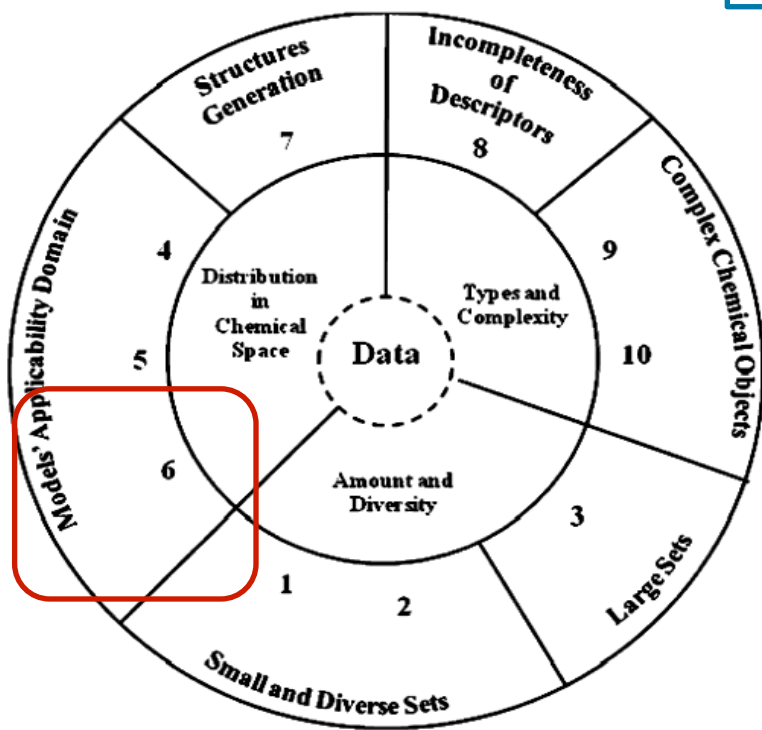


Transductive effect is the difference in prediction performance between transductive and inductive models



Active Learning

Active learning helps to form “optimal” training sets

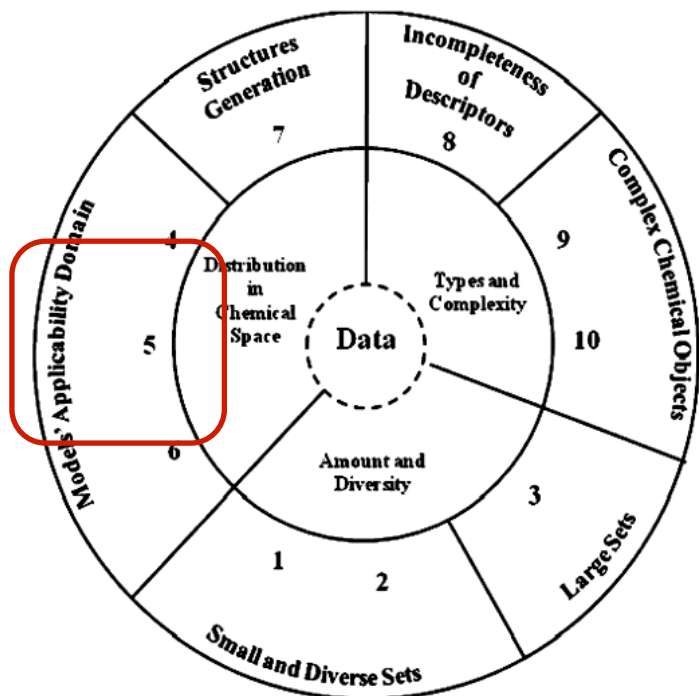


In each learning iteration, the most “useful” compound is selected from a pool, studied in experiment and added to the training set followed by model rebuilding

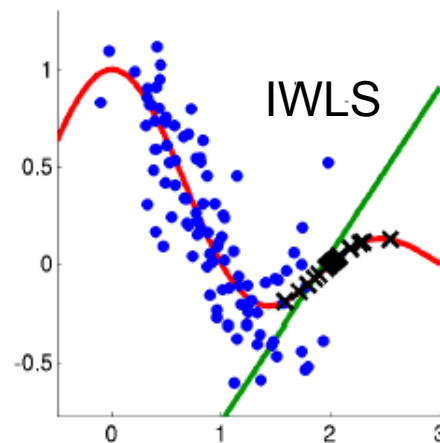
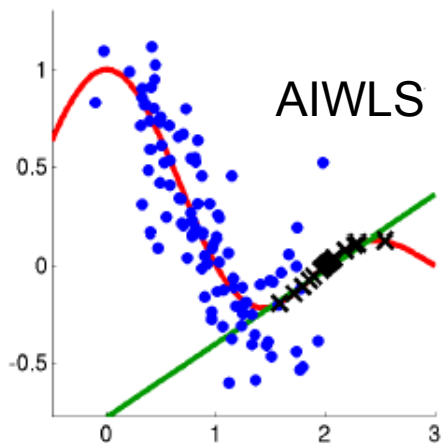
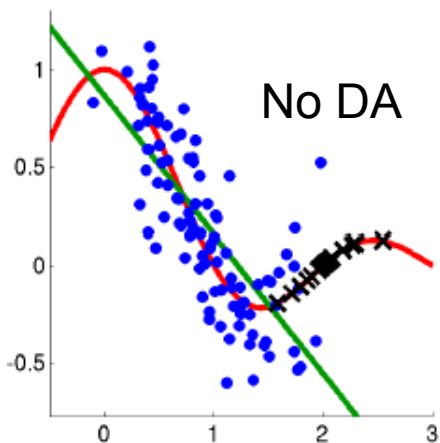
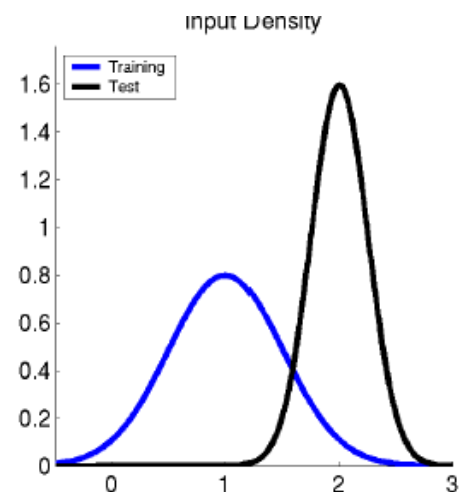
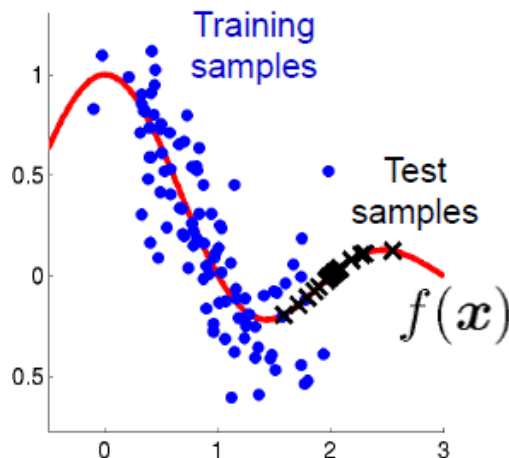


- Burr Settles. Active Learning Literature Survey. Computer Sciences Technical Report 1648, University of Wisconsin-Madison. 2009 (<http://active-learning.net>)
- Y.Fujiwara, Y.Yamashita, T.Osada et al. *J. Chem. Inf. Model.* **2008**, 48 (4), 930–940

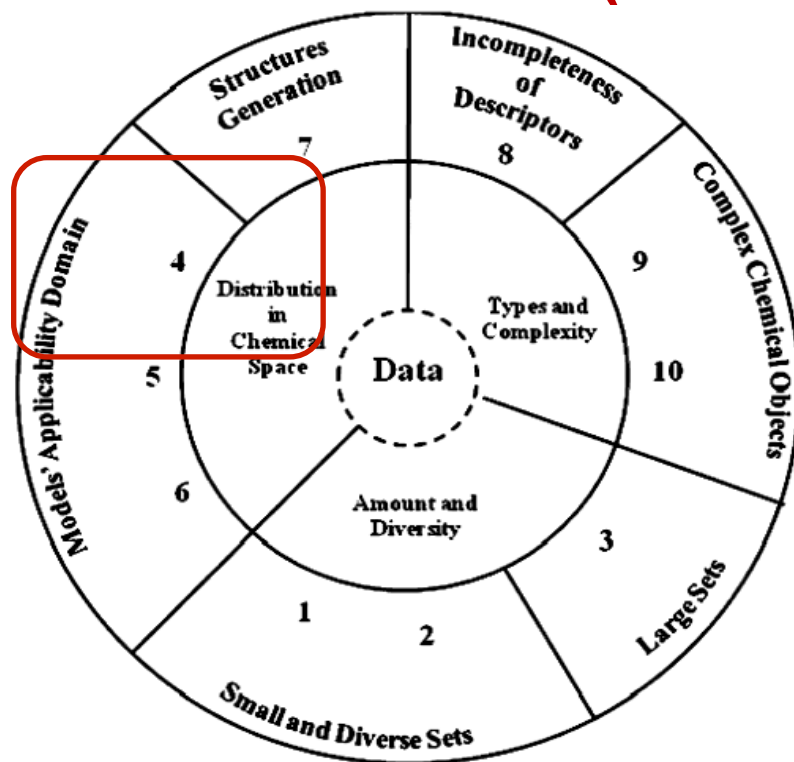
Domain Adaptation



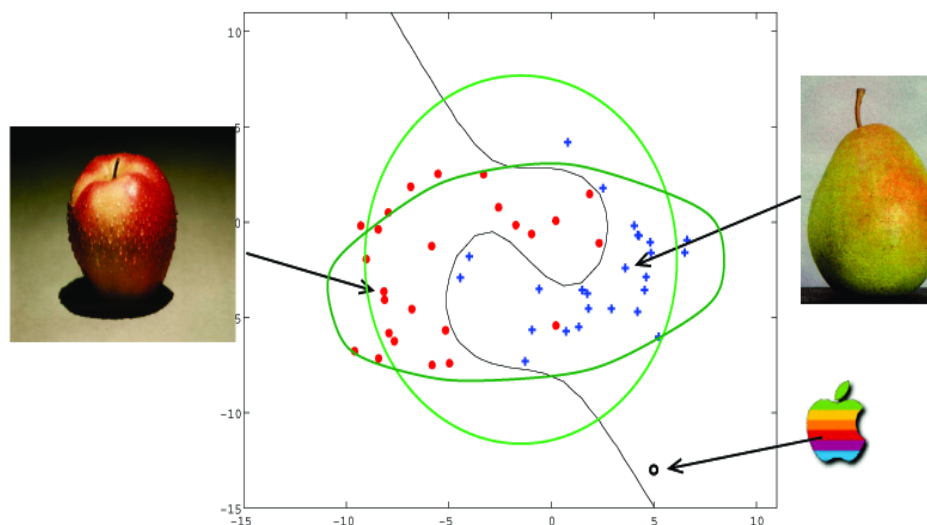
What to do if the training and the test sets are drawn from different distributions?



One-Class Classification (Novelty Detection)



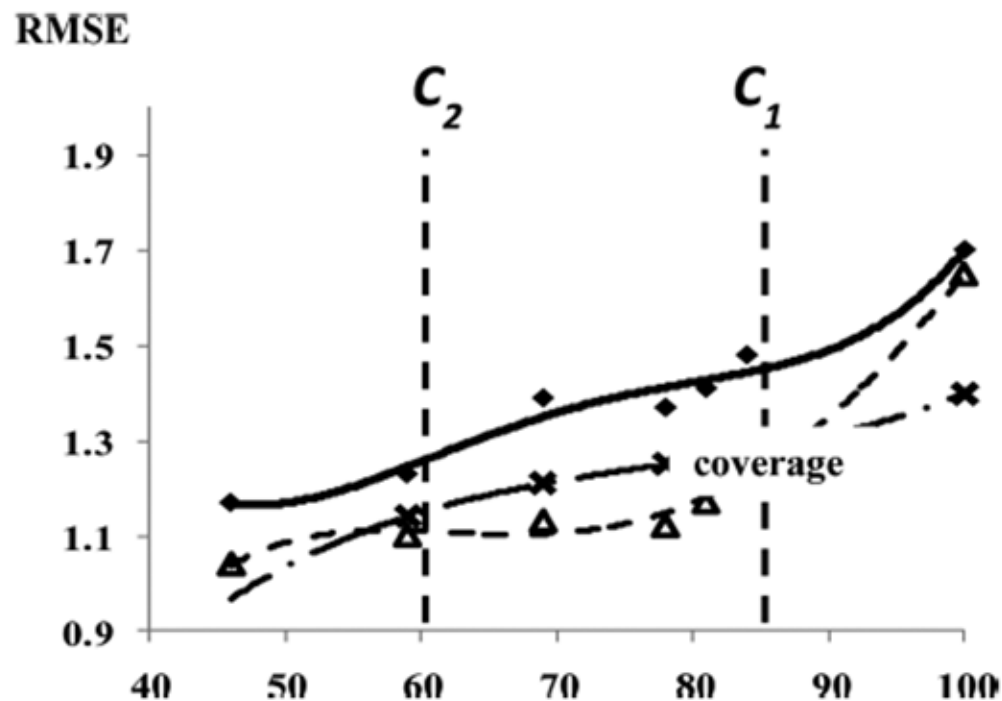
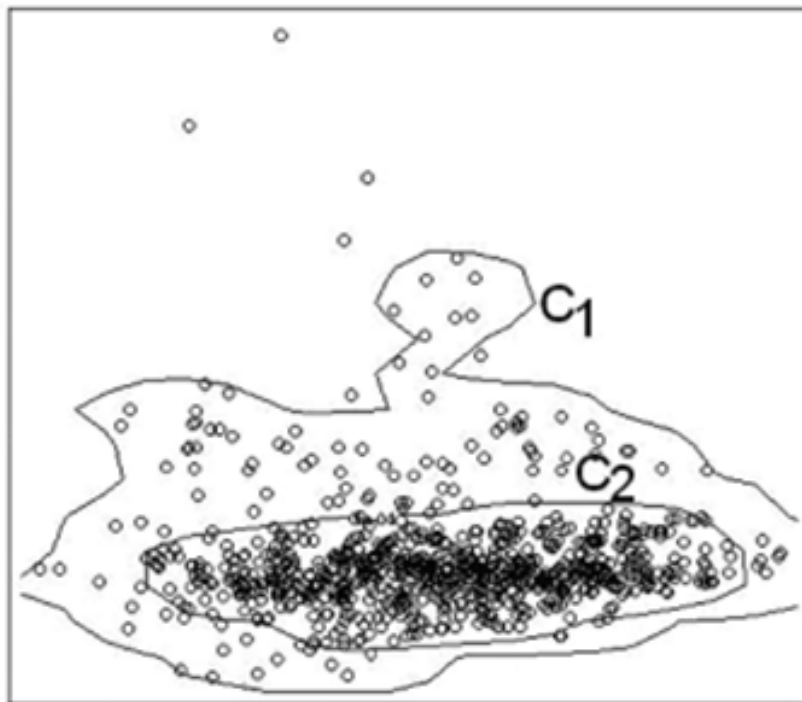
How to build classification models without counterexamples?



One-class classification (or novelty detection) methods allows one to build classification models without counterexamples. In contrast to conventional (two-class) classification, one-class classification tends to describe one single class of objects (*target class objects*), and distinguish it from all other objects (*outliers*).



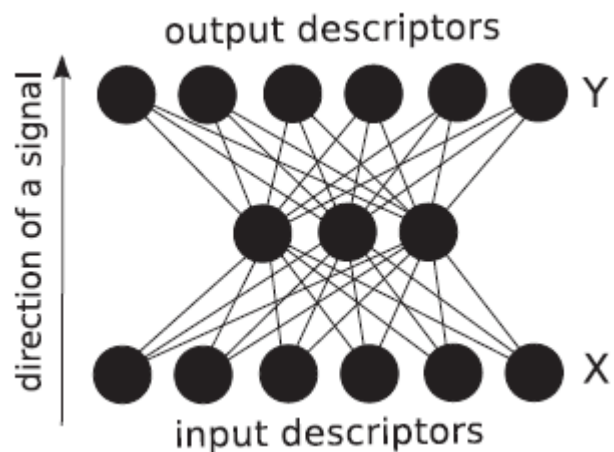
One-Class Classification (OCC) Approach to Defining Model Applicability Domain (AD)



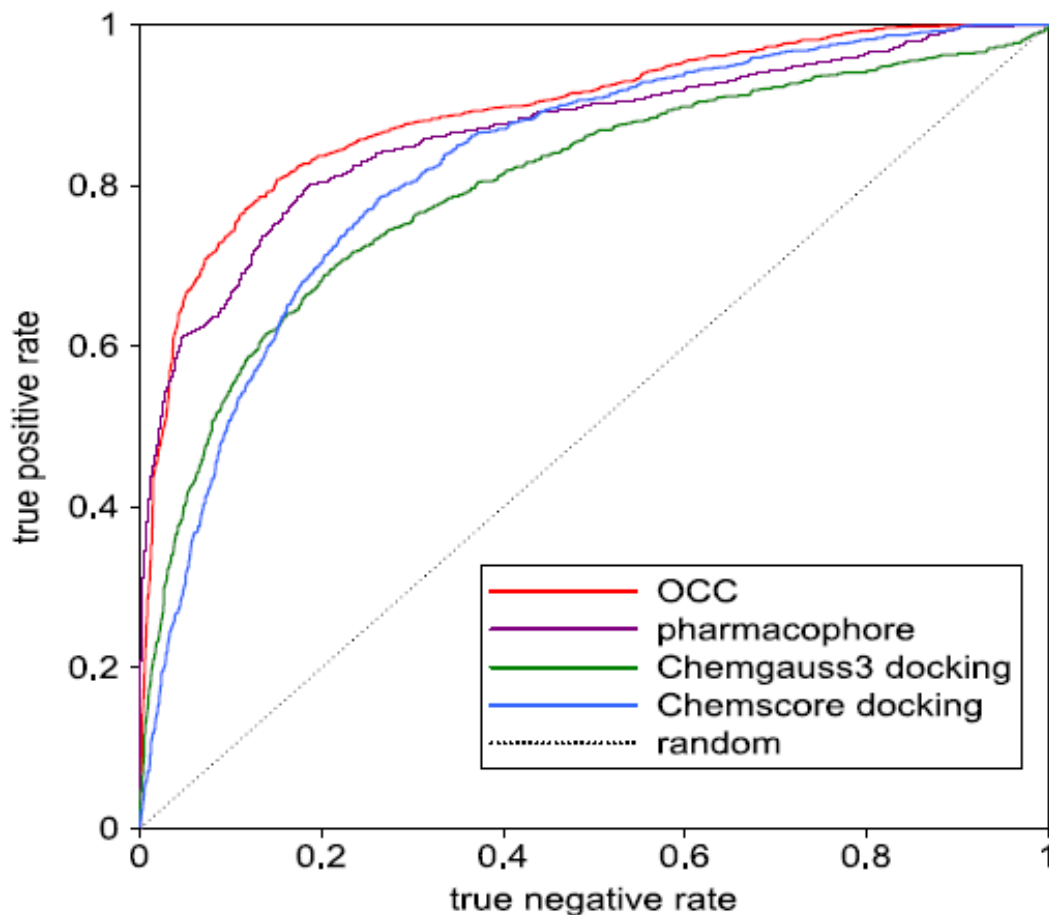
QSPR modeling of stability constants for of Ca^{2+} , Sr^{2+} and Ba^{2+} with organic ligands



Virtual Screening Based on One-Class Classification Using Auto-Encoder Neural Network



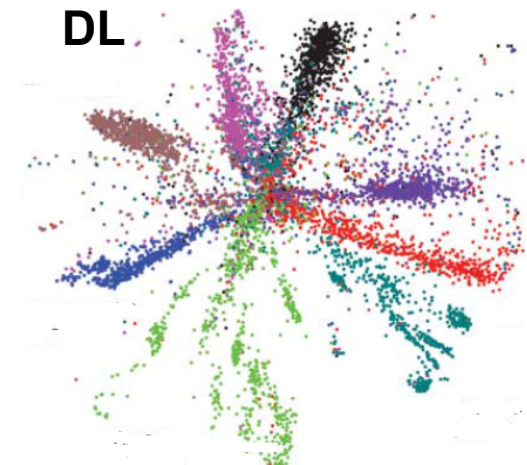
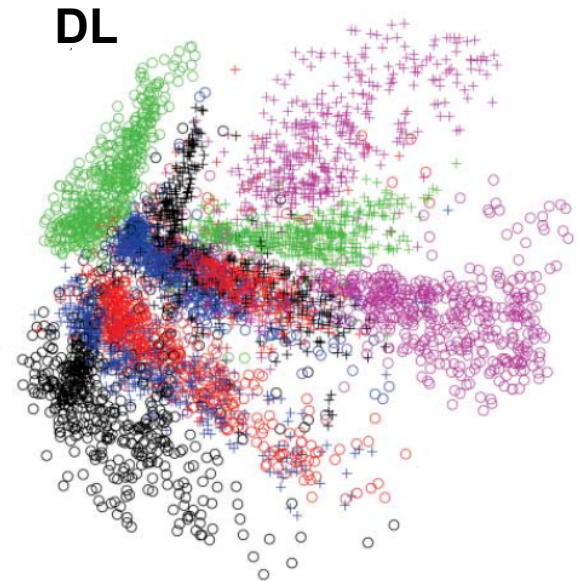
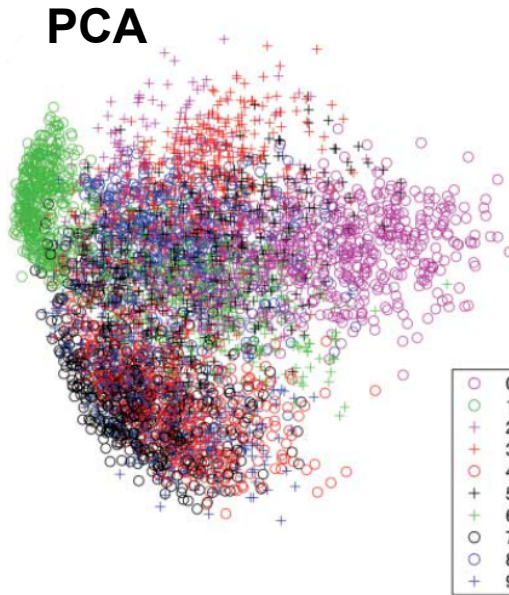
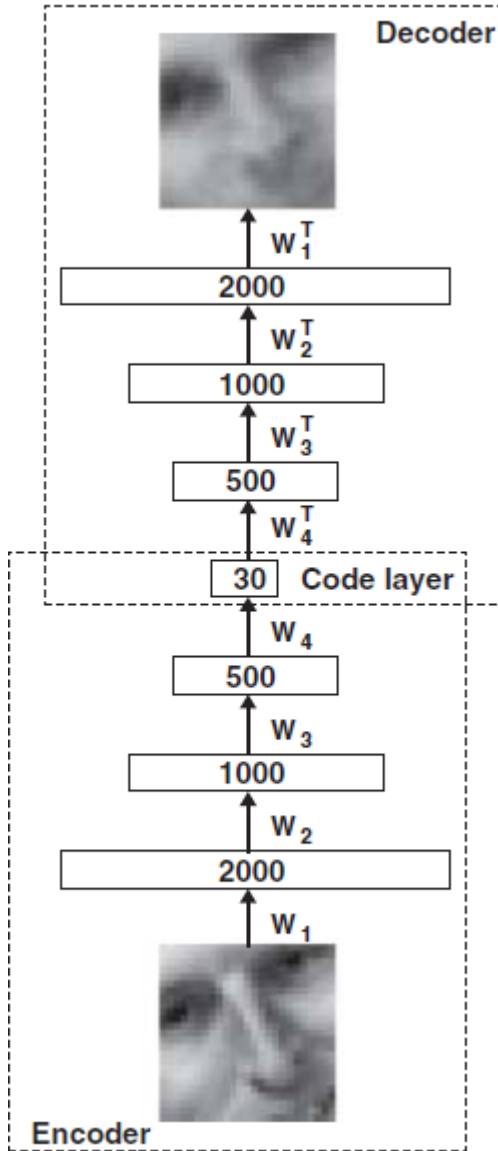
Test compounds with lower reconstruction error are supposed to have more chances to belong to the same activity class as the training compounds



glycogen synthase kinase 3 β inhibitors



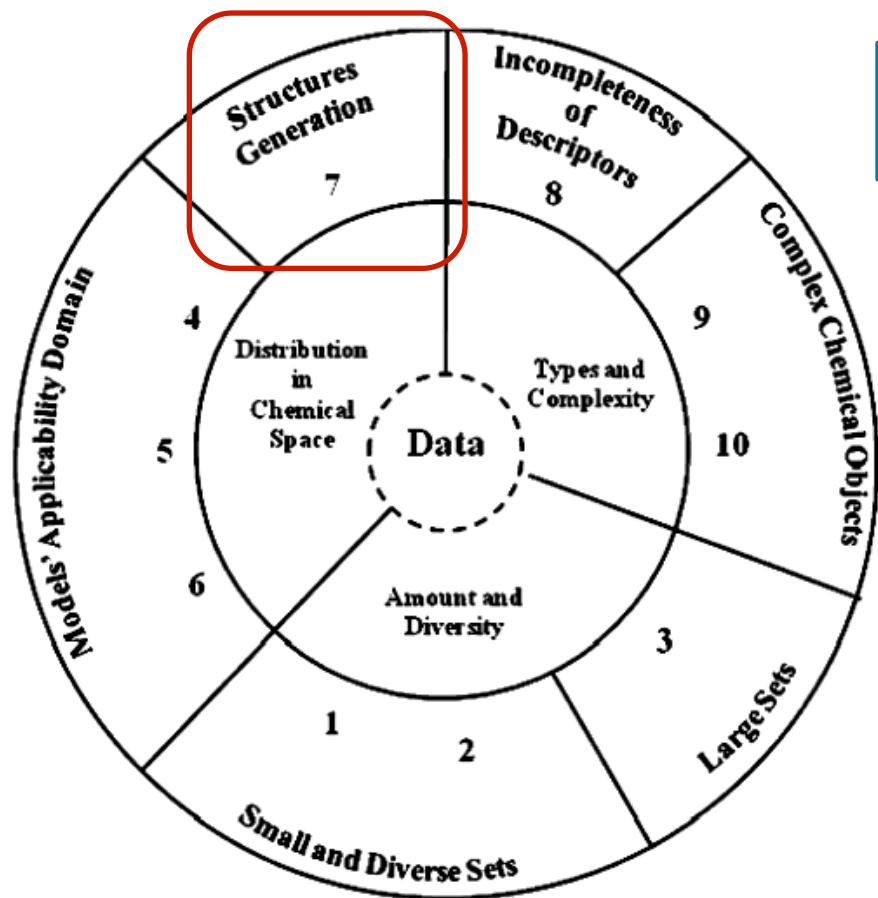
Deep Learning



- G.E.Hinton, R.R.Salakhutdinov, R. R. Science **2006**, 313 (5786), 504-507
- Y.Bengio. *Foundations and Trends in Machine Learning* **2009**, 2 (1), 1-127



Inverse QSAR



How to generate new chemical structures possessing desired properties?

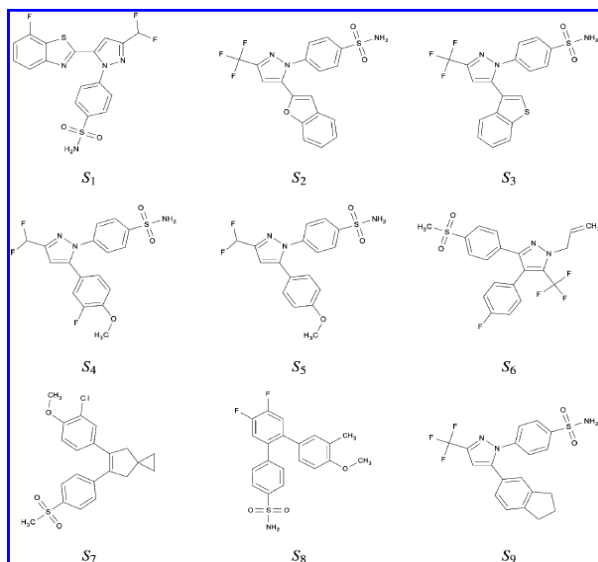
- Structure generation with filtering through QSAR models
- Combinatorial stochastic optimization utilizing QSAR models
- Solving pre-image problem for kernel-based QSAR models
- Building generative models for graphs

- I.I.Baskin et al. *Dokl. Akad. Nauk SSSR* **1989**, 307 (3), 613–617
- Churchwell et al. *J. Mol. Graphics Modell.* **2004**, 22 (4), 263–273
- W.Wong, F.A.Burkowski. *J. Cheminf.* **2009**, 1 (1), 4.
- D.White, R.C.Wilson. *J. Chem. Inf. Model.* **2010**, 50 (7), 1257–1274

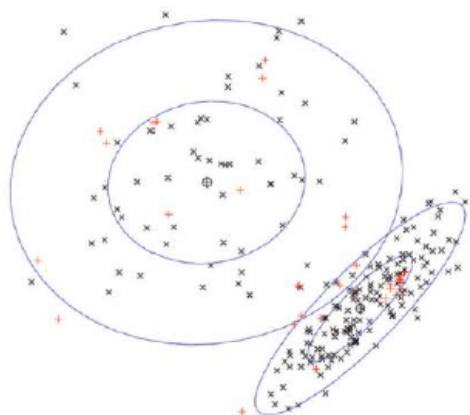


Generative Models for Chemical Graphs

Structures for training



COX2 inhibitors

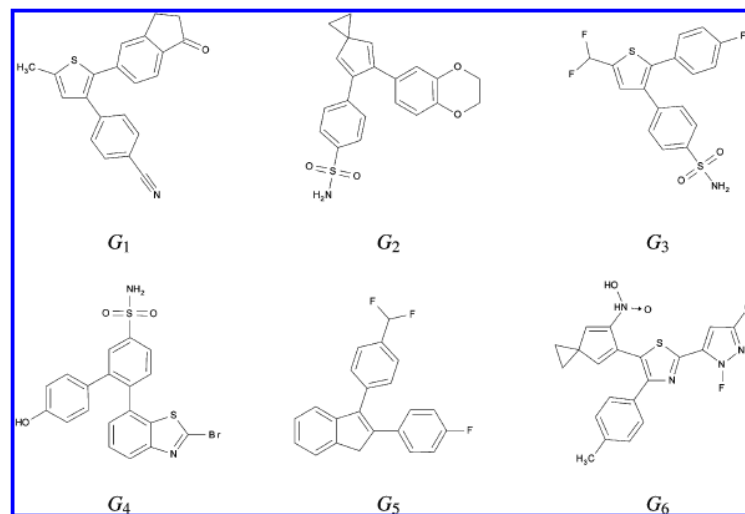


GMM model for $P(X|Y)$

sampling

Generative models are specified by either joint distribution $P(X,Y)$ or conditional distribution $P(X|Y)$

$$P(X|Y) = P(X,Y) / P(Y)$$



Generated structures



Machine Learning Methods for Property Prediction in Chemoinformatics: *Quo Vadis?*

Alexandre Varnek^{*,†} and Igor Baskin^{†,‡}

[†]Laboratoire d'Infochimie, UMR 7177 CNRS, Université de Strasbourg, 4, rue B. Pascal, Strasbourg 67000, France

[‡]Department of Chemistry, Moscow State University, Moscow 119991, Russia

Review of existing mathematical approaches potentially useful but rarely or never used in chemoinformatics



Chemoinformatics Tools and the Appropriate Machine Learning Concepts and Methods

Chemoinformatics problem	Machine learning concept	Machine learning method	Implementation in freely available software
1 Increase of the predictive performance of models built on small and diverse data sets	Ensemble learning ²⁹¹	Different methods of combining classifiers ²⁹² Bagging ⁷⁹ Boosting (classification) ⁸⁸ Boosting (regression) ⁹¹ Stacking ⁸⁶ Random subspace ⁸⁵ Random forest ⁸⁹	meta/ Vote (W) meta/ Bagging (W), adabag (R) meta/ AdaBoostM1 (W), ada, adabag (R) meta/ AdditiveRegression (W) GAM-Boost, mboost (R) meta/ Stacking (W) meta/ RandomSubSpace (W) trees/ RandomForest (W) randomForest (R) SVMlight ²⁹⁶ SGTlight ²⁹⁷ SemL ²⁹⁸
	Semisupervised and transductive learning ^{96,293}	TSVM (transductive SVM) ^{97,294,295} SGT (Spectral Graph Transducer) SemL (Semisupervised Learning) ²³⁰ LapSVM (Laplacian SVM) ²⁹⁹ Semisupervised learning based on one-class classification ³⁰⁰ and ensemble learning ³⁰¹	SVMLight ²⁹⁶ SGTlight ²⁹⁷ SemL ²⁹⁸
	Inductive knowledge transfer, ²⁸¹ multitask learning, ^{153,154} collaborative filtering ²²⁶	Multitask learning using backpropagation neural networks ¹⁵⁴ Multitask learning using multitask kernel, ¹⁵⁵ Bayesian multitask learning, ^{303,304} multitask learning using Partial Least Squares (PLS) method, ³⁰⁵ online multitask learning, ³⁰⁶ multitask learning with data editing ³⁰⁷ semisupervised multitask learning using Dirichlet process, ³⁰⁸ conic programming for multitask learning, ³⁰⁹ multitask learning by multiple kernel learning ³¹⁰	RSNNS, AMORE, neuralnet, nnet (R); SNNS ³⁰²
	L ₁ -Regularized methods	Support Vector Machines(SVM) ^{665,311,312} Ridge regression ³¹⁴	functions/SMO (W); kernlab (R); LibSVM ³¹³ SVMlight ²⁹⁶ functions/LinearRegression (W); Penalized, RXshrink (R)



Acknowledgements

- Alexandre Varnek
- Gilles Marcou
- Dragos Horvath
- Nathalie Kireeva

- Nelly Zhokhova
- Pavel Karpov
- Dmitry Osolodkin



Strasbourg University



Lomonosov Moscow State University



Helmholtz Zentrum München

- Igor Tetko

