

## Drug Design with Machine Learning

1 OVIDIU IVANCIUC  
2 Department of Biochemistry and Molecular Biology,  
3 University of Texas, Medical Branch, Galveston, USA

### 4 Article Outline

5 Glossary  
6 Definition of the Subject  
7 Introduction  
8 Decision Trees  
9 Lazy Learning and  $k$ -Nearest Neighbors  
10 Bayesian Methods  
11 Support Vector Machines  
12 Comparative Studies  
13 Future Directions  
14 Bibliography

### 15 Glossary

16 **Bayesian classifier** Bayes' theorem of conditional proba-  
17 bility is a method of statistical inference that represents  
18 the basis of several classification ML used in drug de-  
19 sign and chemoinformatics to classify libraries of com-  
20 pounds into active and inactive. A Bayesian classifier  
21 considers each structural feature or descriptor inde-  
22 pendent of the other descriptors, and the probability  
23 that a compound is active is proportional to the ra-  
24 tio of active to inactive compounds that have the same  
25 structural feature or have the same value for that de-  
26 scriptor. The final probability that a compound is ac-  
27 tive is a product of all descriptor-based probabilities.  
28 Structural descriptors that are real numbers are usu-  
29 ally binned prior to their evaluation with a Bayesian  
30 classifier.

31 **Decision tree** A decision tree is a sequence of rules ap-  
32 plied to selected structural descriptors. The training  
33 phase comprises the selection of the structural descrip-  
34 tors that are evaluated, the order in which the rules  
35 are applied, and the decision taken at each leaf. Us-  
36 ually, each rules evaluates a descriptor ( $\geq$  or  $<$  than  
37 a threshold) and splits the objects into two or more  
38 populations. Then each population is selected and the  
39 splitting procedure is performed with a new rule, until  
40 a stopping condition is met (for example, when all ob-  
41 jects in the population belong to the same class). The  
42 prediction phase starts from the root node and eval-  
43 uates each rule on a pathway determined by the out-  
44 come (true or false) of the previous rule. When a leaf  
45 is reached the algorithm predicts the class of the object

(classification trees) or the numerical value of a prop-  
erty (regression trees).

**$k$ -nearest neighbors**  $k$ -nearest neighbors ( $k$ -NN) is a su-  
pervised learning algorithm that predicts the property  
of an object based on a local interpolation model. In  
classification, the class of a new object is predicted  
based on the majority vote of its  $k$  nearest neighbors. In  
regression, the property value for a new object is pre-  
dicted as an average value of the property values for its  
 $k$  nearest neighbors.

**Lazy Learning** Lazy learning is a memory based local  
learning that defers the computation until a prediction  
is requested for an object. The first step is to insert the  
query object into the space of the training objects, and  
to identify the training objects located in a set neigh-  
borhood. The predicted property of the query object is  
based on an interpolation of the properties of the ob-  
jects situated in the neighborhood.

**Machine learning** Machine learning is an important field  
of artificial intelligence, and includes a diversity of  
methods and algorithms that extract rules and func-  
tions from large datasets, such as decision trees,  
lazy learning,  $k$ -nearest neighbors, Bayesian methods,  
Gaussian processes, support vector machines, and ker-  
nel algorithms. Machine learning algorithms extract  
information from experimental data by computational  
and statistical methods and generate a set of rules,  
functions or procedures that allow them to predict the  
properties of novel objects that are not included in the  
learning set.

### Quantitative structure-activity relationships

Quantitative structure-activity relationships (QSAR)  
represent regression models that define quantita-  
tive correlations between the chemical structure of  
molecules and their physical properties (boiling point,  
melting point, aqueous solubility), chemical properties  
and reactivities (chromatographic retention, reaction  
rate), or biological activities (cell growth inhibition,  
enzyme inhibition, lethal dose). The fundamental  
hypotheses of QSAR is that similar chemicals have  
similar properties, and small structural changes result  
in small changes in property values. The general form  
of a QSAR equation is  $P(i) = f(\mathbf{SD}_i)$ , where  $P(i)$  is  
a physical, chemical, or biological property of com-  
pound  $i$ ,  $\mathbf{SD}_i$  is a vector of structural descriptors of  $i$ ,  
and  $f$  is a mathematical function such as linear regres-  
sion, partial least squares, artificial neural networks, or  
support vector machines. A QSAR model for a prop-  
erty  $P$  is based on a dataset of chemical compounds  
with known values for the property  $P$ , and a matrix of  
structural descriptors computed for all chemicals. The

Please note that the pagination is not final; in the print version an entry will in general not start on a new page.

learning (training) of the QSAR model is the process of determining the optimum parameters of the regression function  $f$ . After the training phase, a QSAR model may be used to predict the property  $P$  for novel compounds that are not present in the learning set of molecules.

**Support vector machines** Support vector machines (SVM) are a class of supervised machine learning methods based on the structural risk minimization and the statistical learning theory of Vapnik. SVM may be applied to data classification and regression, using selected objects (support vectors) to generate the SVM model. Nonlinear classification problems are transformed into linear classification problems by using kernel functions that combine the input space into a higher-dimensional feature space in which a hyperplane may discriminate the classes. An SVM classification model computes a maximum margin hyperplane that separates the classes in the feature space. The maximal margin hyperplane maximizes the distance to the hyperplane of the closest patterns from the two classes. An SVM regression model builds a regression tube with the property that all objects inside the tube do not contribute to the overall error of the model. The shape of the regression tube is determined by selected objects (support vectors) situated outside the tube.

**Structural descriptor** A structural descriptor (SD) is a numerical value computed from the chemical structure of a molecule, which is invariant to the numbering of the atoms in the molecule. Structural descriptors may be classified as constitutional (counts of molecular fragments, such as rings, functional groups, or atom pairs), topological indices (computed from the molecular graph), geometrical (volume, surface, charged-surface), quantum (atomic charges, energies of molecular orbitals), and molecular field (such as those used in CoMFA, CoMSIA, or CoRSA).

**Structure-activity relationships** Structure-activity relationships (SAR) represent classification models that can discriminate between sets of chemicals that belong to different classes of biological activities, usually active/inactive towards a certain biological receptor. The general form of a SAR equation is  $C(i) = f(\mathbf{SD}_i)$ , where  $C(i)$  is the activity class of compound  $i$  (active/inactive, inhibitor/non-inhibitor, ligand/non-ligand),  $\mathbf{SD}_i$  is a vector of structural descriptors of  $i$ , and  $f$  is a classification function such as  $k$ -nearest neighbors, linear discriminant analysis, random trees, random forests, Bayesian networks, artificial neural networks, or support vector machines.

## Definition of the Subject

The process of drug discovery has the goal to identify lead chemicals that have a significant activity against a selected biological target. A disease state may be the result of changes in the structure and function of cell-signaling receptors, enzymes, hormone receptors, or other functional proteins. The drug target is a protein whose activity is modulated by its interaction with a chemical compound, and thus may control a disease. The lead compounds identified in the drug discovery step are optimized in the drug development phase that results in a small number of chemicals that are evaluated in human clinical trials. The first priority in drug development is to increase the biological activity of a lead compound while preserving its drug-like properties. The lead compound is expanded into a chemical library that conserves the structure responsible for the biological activity (pharmacophore) and adds chemical groups that might improve its activity. Then the chemicals are synthesized and tested in biological assays against the target, which may result in the identification of more active compounds. The cycle “library design – chemical synthesis – biological assay” is repeated several times until a number of good candidates are identified for clinical trials. The design of effective drug candidates is a multi-objective optimization problem, because simultaneous with a good biological activity the chemicals must pass several other important tests, including pharmacokinetics, pharmacodynamics, toxicity, mutagenicity, metabolism, and excretion.

Computer-assisted drug design (CADD) uses computational chemistry to increase the chances of finding valuable drug candidates. CADD methods include machine learning, structure-activity relationships (SAR), quantitative structure-activity relationships (QSAR), molecular mechanics, quantum mechanics, molecular dynamics, and drug-protein docking. SAR and QSAR are based on the theory that chemical structure determines all physical, chemical and biological properties of a molecule. To obtain structure-activity models, the chemical structure is characterized with structural descriptors and then ML models are used to identify a statistical relationship between the descriptor space and a molecular property. Other considerations for the application of SAR and QSAR in drug design are that similar molecules have similar properties, and that a small modification in a chemical structure result in a small modification in its properties.

Machine learning (ML) procedures are applied in drug discovery in an iterative way, in which experimental activity data are used to train ML models which in turn offer predictions for novel molecules with improved bi-

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197 ological properties. Then the molecules designed with  
198 ML are tested in biological assays, and the best can-  
199 didates are selected for further cycles of optimization.  
200 ML models may greatly improve the chances of finding  
201 lead molecules by screening a larger diversity of chem-  
202 ical topologies. In the classical approach, a number of  
203 chemicals in the range  $10^3$ – $10^4$  is synthesized and assayed  
204 to identify lead compounds, whereas the introduction of  
205 combinatorial chemistry and high-throughput screening  
206 (HTS) increased these numbers to  $10^5$ – $10^6$ . Although the  
207 robots offer great speed and reliability in HTS, the chem-  
208 ical compounds still have to be synthesized before screen-  
209 ing, thus limiting the diversity of structures evaluated in  
210 HTS. This is why ML models are used to enhance the lead  
211 identification process with a virtual HTS (vHTS) screening  
212 of  $10^7$ – $10^8$  molecules. Based on experimental results from  
213 a HTS campaign, ML models are computed and then used  
214 in vHTS experiments. These ML are mainly classification  
215 models that predict molecules that have a high probabili-  
216 ty of interacting with the selected target. vHTS has sev-  
217 eral obvious advantages compared to HTS, mainly coming  
218 from the fact that the chemical compounds are generated  
219 only in a computer, which opens the possibility to explore  
220 a larger diversity of chemical skeletons and a much higher  
221 number of molecules.

222 Classification (SAR) and regression (QSAR) ML mod-  
223 els are applied during the drug development cycles to op-  
224 timize the biological activity, target selectivity, and other  
225 physico-chemical and biological properties of selected  
226 chemicals. ML models are essential tools in increasing the  
227 chances of bringing a drug to market, but only when in-  
228 tegrated with the usual chemical, biological, pharmaco-  
229 logical and clinical procedures used by the pharmaceuti-  
230 cal industry. Boid enumerated several successful CADD  
231 applications in which computational methods had a deci-  
232 sive contribution to the discovery of a drug, namely  
233 norfloxacin (Merck, antibacterial), losartan (Merck, an-  
234 tihypertensive), dorzolamide (Merck, antiglaucoma), ri-  
235 tonavir (Abbott, antiviral), indinavir (Merck, antiviral),  
236 donepezil (Esai, anti-Alzheimer's disease), zolmitriptan  
237 (AstraZeneca, antimigraine), nelfinavir (Pfizer, antiviral),  
238 amprenavir (GlaxoSmithKline, antiviral), zanamivir  
239 (GlaxoSmithKline, antiviral), oseltamivir (Roche, antiviral),  
240 lopinavir (Abbott, antiviral), imatinib (Novartis, anti-  
241 neoplastic), erlotinib (OSI, antineoplastic) [17]. This is an  
242 impressive list that adds convincing evidence to advocate  
243 the integration of machine learning and other chemoin-  
244 formatics methods in drug discovery and development.

## Introduction

245

246 Machine learning is an important field of artificial intel-  
247 ligence, and includes a diversity of methods and algo-  
248 rithms that extract rules and functions from large datasets,  
249 such as linear discriminant analysis (LDA), artificial neu-  
250 ral networks (ANN), decision trees, lazy learning,  $k$ -near-  
251 est neighbors, Bayesian methods, Gaussian processes,  
252 support vector machines (SVM), and kernel algorithms.  
253 Several influential books are recommended for a detailed  
254 overview of ML algorithms: *Pattern Recognition and Neu-  
255 ral Networks* by Ripley [110], *Machine Learning and Data  
256 Mining* by Kononenko and Kukar [83], *Pattern Recogni-  
257 tion and Machine Learning* by Bishop [15], *Data Mining:  
258 Practical Machine Learning Tools and Techniques* by Wit-  
259 ten and Frank [151], *Introduction to Machine Learning*  
260 by Alpaydin [4], *The Elements of Statistical Learning* by  
261 Hastie, Tibshirani, and Friedman [52], *Pattern Classifica-  
262 tion* by Duda, Hart, and Stork [35], *Machine Learning* by  
263 Mitchell [96], and *Neural Networks for Pattern Recognition*  
264 by Bishop [16].

265 ML algorithms extract information from experimen-  
266 tal data by computational and statistical methods and gen-  
267 erate a set of rules, functions or procedures that allow  
268 them to predict the properties of novel objects that are  
269 not included in the learning set. In drug design the task  
270 is to learn how chemical structure determines some im-  
271 portant drug property, such as physico-chemical prop-  
272 erties (aqueous solubility, skin penetration, hydropho-  
273 bicity, intestinal absorption, blood-brain barrier penetra-  
274 tion), biological activity (enzyme inhibition), metabolism,  
275 toxicity, mutagenicity, or excretion. Each molecule is as-  
276 sociated with a set of structural descriptors and an ex-  
277 perimental property. The structural descriptors are used  
278 as input for an ML algorithm and the property is the  
279 target (output) of the model. A structural descriptor is  
280 a structural feature of the chemical structure, and can  
281 be a list of substructures or substituents, graph descrip-  
282 tors [19,64,136], topological indices [8,18,65], connectiv-  
283 ity indices [74,75] electrotopological indices [76], or in-  
284 dices derived from the molecular geometry and quan-  
285 tum calculations [73,133]. The experimental property  
286 that is modeled may be a class label (+1/–1), such as  
287 soluble/non-soluble, inhibitor/non-inhibitor, ligand/non-  
288 ligand, toxic/non-toxic, mutagen/non-mutagen, carcino-  
289 gen/non-carcinogen. For this type of data one obtains  
290 a classification model or SAR. Classification models are  
291 mainly used to filter large chemical libraries and reduce  
292 them a small number of chemicals with a high probability  
293 of having a desired property (ligand for a biological target,  
294 inhibitor for an enzyme, non-toxic, non-mutagen, or non-

carcinogen). The experimental property may be also a continuous value, such as inhibition constant to an enzyme, binding constant to a target, hydrophobicity, aqueous solubility, skin penetration, or intestinal absorption, in which case one computes a regression model or QSAR. Regression models, which may be linear or non-linear, are used to optimize the drug-like properties of a chemical compound. A drug-related property may be considered categorical variable or a continuous variable, depending on the drug development phase in which is applied. For example, in developing CNS (central nervous system) drugs, the blood-brain barrier penetration is used as a categorical variable in the early stages of drug discovery, to identify and eliminate compounds that do not pass the blood-brain barrier. In later stages of drug development this property is used as a continuous variable in QSAR models to optimize the brain concentration of a small number of selected chemicals that show promising pharmacological properties. Both SAR and QSAR models belong to the class of supervised learning algorithms, because the target (output) value is provided together with the input chemical structures and descriptors. Supervised learning algorithms have two distinct phases. The first one is learning or training, in which a ML method is used to learn the relationships between input (structural descriptors) and output (an experimental property provided for each molecule). The result is a model that can be a statistical function or a set of rules. The second phase of a ML algorithm is the prediction, when the trained model is used to predict the property for novel chemicals that were not present in the training set or that are not even synthesized. The predictions obtained are then used as a guide in synthesizing and testing novel chemicals.

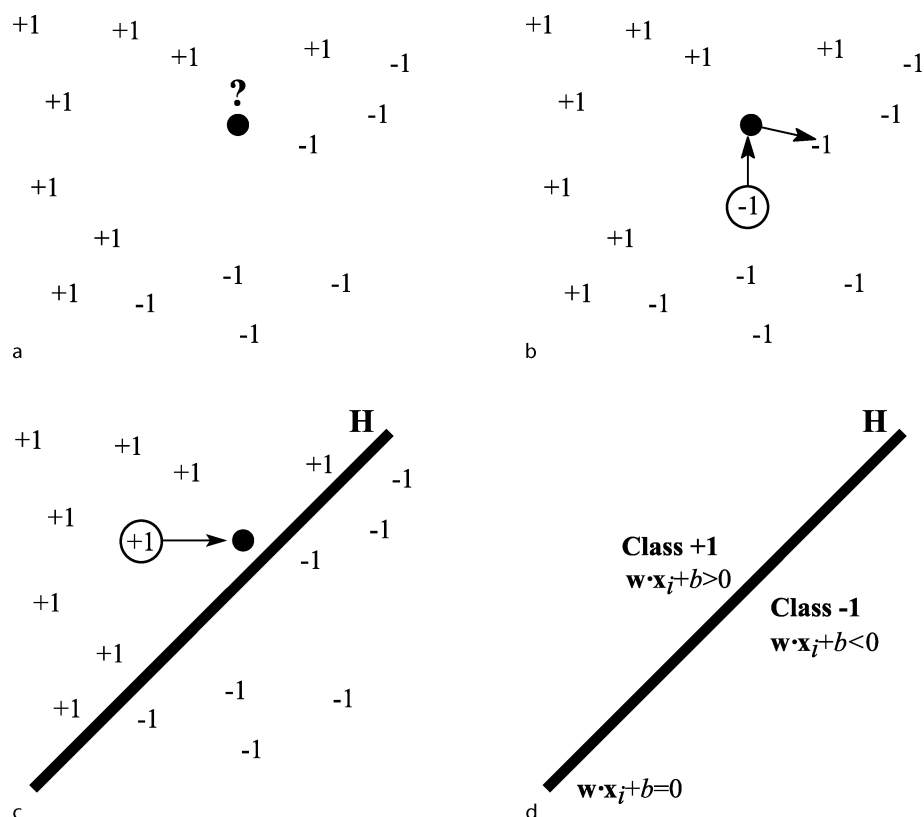
In unsupervised learning the dataset contains only chemical structures and structural descriptors, without output values. The ML objective is to identify how molecules form clusters based on their structural similarity. Obviously, starting from a particular set of structural descriptors, different similarity indices and different clustering algorithms will result in different clusters. Clustering algorithms are applied in drug design to identify similar molecules in chemical libraries. For example, starting from a lead compound one can find similar compounds in catalogs and databases, and thus purchase chemicals that might be useful in the drug discovery process. Another application is to identify which chemical structures are under-represented or missing from a company collection, and to guide compound synthesis and acquisition.

To illustrate two basic approaches in data classification we consider the objects from Fig. 1a. The two classes of objects may represent two populations of molecules,

for example enzyme inhibitors form class +1, whereas non-inhibitors form class -1. From the distribution of the two classes one can identify a cluster for class +1 and a distinct cluster for class -1. Such a clear-cut situation is not the norm in drug design application, because classes of molecules may overlap which makes more difficult the property prediction for novel chemicals. The problem considered in Fig. 1a is to use the information provided by the two classes of objects to predict the class of the unknown object marked with “?”. The first ML tested is  $k$ -nearest neighbors ( $k$ -NN), that belongs to the class of lazy learning algorithms, and performs a local approximation of the model that is computed only in the moment of the prediction. An object is classified in the most populated class among its  $k$  nearest neighbors, where for two-class problems  $k$  must be odd to avoid undecided situations. We test here the simplest variant, 1-nearest neighbor, which considers only the closest neighbor. An inspection of the distances between the prediction object • and all other objects shows that the closest object belongs to class -1, and thus the prediction object is assigned to class -1 (Fig. 1b). The predicted class may change if one considers a larger number of nearest neighbors, and the optimum value for  $k$  is usually determined through cross-validation.

An efficient procedure to classify linearly separable classes is represented by a hyperplane  $H$  (Fig. 1c) that separates the descriptor space into a region for class +1 and another region for class -1. A new object that is located in the space region that belongs to class +1 will be assigned to class +1 irrespective of its distance to the separating hyperplane. Similarly, if the new object is situated in the space region that belongs to class -1, then the object is predicted in class -1. The unknown object from Fig. 1a is assigned to class +1 by the hyperplane  $H$  (Fig. 1c). For the same classification problem we obtained two different predictions for the unknown object, namely class -1 with 1-NN and class +1 with a hyperplane. Such divergent predictions obtained with different ML methods are quite common, and it is not unusual to obtain different predictions for the same object just by changing the ML parameters. These practical aspects of ML should receive proper consideration, and it is a good practice to evaluate a large diversity of ML methods to obtain consensus predictions.

The two ML demonstrated here are very different in their formulation and their properties.  $k$ -NN can be used for multi-class discrimination and it is a nonlinear classifier that is useful for classes that cannot be separated with a linear hypersurface. Also,  $k$ -NN does not produce a statistical function to replace the data, and the learning phase is missing. The second classifier replaces the data with the equation of a hyperplane (Fig. 1d). This function may be



Drug Design with Machine Learning, Figure 1

Classification with machine learning: a a class prediction (+1/−1) is sought for the object •; b a  $k$ -nearest neighbor classifier ( $k = 1$ ) predicts the object • in class −1; c a linear classifier defined by the hyperplane  $H$  predicts the object • in class +1; d the classification hyperplane  $H$  separates the space in region for class +1 and another region for class −1

397 used only for two-class problems in which the objects may  
 398 be separated by linear hypersurface. A hyperplane  $H$  has  
 399 the formula  $\mathbf{w} \cdot \mathbf{x} + b = 0$ , where  $\mathbf{w}$  is the normal vector  
 400 to  $H$ . A molecule  $i$  characterized by a vector of structural  
 401 descriptors  $\mathbf{x}_i$  belongs to class +1 if  $\mathbf{w} \cdot \mathbf{x}_i + b > 0$ , or to  
 402 class −1 if  $\mathbf{w} \cdot \mathbf{x}_i + b < 0$ . The rules for predicting the class  
 403 for an unknown object  $k$  are:

$$404 \quad \text{class}(k) = \begin{cases} +1 & \text{if } \mathbf{w} \cdot \mathbf{x}_k + b > 0 \\ -1 & \text{if } \mathbf{w} \cdot \mathbf{x}_k + b < 0. \end{cases} \quad (1)$$

405 The hyperplane classifier may be extended with kernel  
 406 functions to nonlinear classification, as it will be shown in  
 407 Sect. “Comparative Studies” in which we give an overview  
 408 of support vector machines and their applications in structure-  
 409 activity models.

410 A decision tree performs a mapping of the structural  
 411 descriptors to a conclusion about a property of the  
 412 object. The learning process establishes a series of successive  
 413 rules that are based on the numerical values of a subset of  
 414 descriptors. Rules are applied in a set se-

415 quence and at each branching point a rule is evaluated  
 416 and its outcome (true or false) determines which branch  
 417 is selected in the next step. Starting from the root node,  
 418 the rules are evaluated and depending on the actual numerical  
 419 values of the descriptors the pathway ends on a leaf that  
 420 assigns a class label or a value to the investigated property.  
 421 The class of decision tree algorithms includes the C4.5 decision  
 422 tree with naïve Bayes classifiers at the leaves [82], the alternating  
 423 decision tree ADTree [42], random trees and random forests  
 424 [22]. Decision trees are computationally very efficient, and  
 425 their predictions are usually better than those obtained with  
 426 more complicated ML algorithms. Drug design applications  
 427 of decision trees include diverse applications, such as prediction  
 428 of drugs that cross the blood-brain barrier [30], structure-  
 429 activity relationships for estrogen receptor ligands [135],  
 430 prediction of  $P$ -glycoprotein substrates [91], evaluation of  
 431 the cytochrome P450 metabolism [154], identification of drug-  
 432 like chemicals [114], design of target chemical libraries [31],  
 433 prediction of protein-protein interaction [98], and modeling the

415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435

oral absorption of drugs [57]. A review of these studies is presented in Sect. “Decision Trees”.

Lazy learning is a memory-based local learning method that stores the training objects and performs all computations only when a prediction request is received by the system [5,6]. The prediction is performed by inserting the unknown object into the prediction space and then identifying all objects situated in its neighborhood. Each prediction is made by computing a local model from the neighboring objects. Lazy learning may be used both for classification and regression, with predictions based on local interpolations of selected objects according to a distance measure. Local learning may perform better than global learning, as shown in a number of SAR and QSAR studies that use various techniques of lazy learning [47,85,159].

The  $k$ -nearest neighbors ( $k$ -NN) algorithm is a lazy learning method that can be used for classification and regression [1]. The training set of objects is stored in the memory, and then each prediction is made by computing a local interpolation model. The prediction object is added to the descriptor space in which the training objects are placed and then its  $k$  nearest neighbors are identified. For classification, the class of the prediction object is assigned to class of the majority of its neighbors. Ties are avoided by selected an odd value for  $k$ , and the best value for  $k$  is usually determined by cross-validation. For regression, the predicted property for the query object is an average value of the property values for its  $k$  nearest neighbors. A distance weighting may be added in such a way to increase the contribution of closer objects and to decrease the influence of more distant objects. Another lazy learning algorithm is  $K^*$  which uses an entropy-based distance function [26]. Basak was an early proponent of the  $k$ -NN regression for physico-chemical and biological properties [10,11,49]. The method was further refined by Tropsha by optimizing the descriptor selection,  $k$ , and the distance weighing function [60,160,161]. Other  $k$ -NN applications were proposed for three-dimensional QSAR [2] and for ligand-based virtual screening of chemical libraries [54]. In Sect. “Lazy Learning and  $k$ -Nearest Neighbors” we review SAR and QSAR models computed with lazy learning and  $k$ -nearest neighbors.

Bayes’ theorem provides mathematical tools that explain how the probability that a theory is true is affected by a new piece of evidence [12]. Several Bayesian classifiers [71,151] were proposed to estimate the probability that a chemical is active based on a set of descriptors and a chemical library with known class attributes (active/inactive). Bayesian classifiers consider that each descriptor is statistically independent of all other descrip-

tors. Binary descriptors, structural counts, fingerprints and other descriptors that have integer values are directly evaluated with Bayesian classifiers, whereas real number descriptors are discretized and transformed into an array of bins. The probability that a compound is active is proportional to the ratio of active to inactive compounds that have the same structural feature or have the same value for that descriptor. A final probability is computed as the product of all descriptor-based probabilities. Bayesian classifiers have numerous applications in chemoinformatics and structure-activity models, such as for the prediction of multidrug resistance reversal [125], to estimate the phospholipidosis inducing potential [105], to identify chemical similar to natural products [39], to improve high-throughput docking [78,79], and to screen virtual chemical libraries [77]. Drug design applications of Bayesian methods are presented in Sect. “Bayesian Methods”.

Support vector machines belong to a heterogeneous group of machine learning methods that use kernels to solve efficiently high-dimensional nonlinear problems. SVM extend the generalized portrait algorithm developed by Vapnik [144] by using elements of statistical learning theory [143] that describe the ML properties that guarantee dependable predictions. Vapnik elaborated further the statistical learning theory in three more recent books, *Estimation of Dependencies Based on Empirical Data* [138], *The Nature of Statistical Learning Theory* [139], and *Statistical Learning Theory* [140]. Vapnik and co-workers developed the current formulation of the SVM algorithm at AT&T Bell Laboratories [20,25,28,33,50,116,141,142].

SVM models have several interesting and appealing properties, namely the maximum margin classification, the kernel transformation of the input space into a feature space where a hyperplane may separate the classes, and a unique solution. The SVM introduction generated an enormous interest, comparable only with that produced by the development of artificial neural networks. Many applications were investigated in a short period of time, simultaneous with developments of novel SVM algorithms, such as least squares SVM [126,127] and other kernel algorithms [21,58,119]. The theory and applications of SVM are presented in a number of books, including *Learning with Kernels* by Schölkopf and Smola [115], *Learning Kernel Classifiers* by Herbrich [53], *An Introduction to Support Vector Machines* by Cristianini and Shawe-Taylor [29], and *Advances in Kernel Methods: Support Vector Learning* by Schölkopf, Burges, and Smola [117]. SVM have numerous applications in drug design, from screening of chemical libraries to QSAR [70], with accurate predictions that frequently surpass those obtained with more established

538 algorithms. In Sect. “Support Vector Machines” we review  
539 the most important aspects of SVM models and we present  
540 several applications to drug design, SAR and QSAR.

541 With an ever-increasing number of machine learn-  
542 ing algorithms available, it is difficult to select those  
543 methods that have better chances of solving a particu-  
544 lar SAR or QSAR problem. To address this topic, we re-  
545 view in Sect. “Comparative Studies” several studies that  
546 compare diverse ML models for drug design and struc-  
547 ture-activity relationships. A comparison of 21 machine  
548 learning algorithms is presented for data selected from  
549 the National Cancer Institute 60-cell line screening panel  
550 (NCI-60) [69]. The structure-anticancer activity models  
551 are developed for three cell lines, namely for lung large cell  
552 carcinoma NCI-H460, glioma SF-268, and melanoma SK-  
553 MEL-5. Finally, we present an assessment of the results ob-  
554 tained in the machine learning competition CoEPrA 2006  
555 (Comparative Evaluation of Prediction Algorithms, <http://www.coepra.org/>).  
556

## 557 Decision Trees

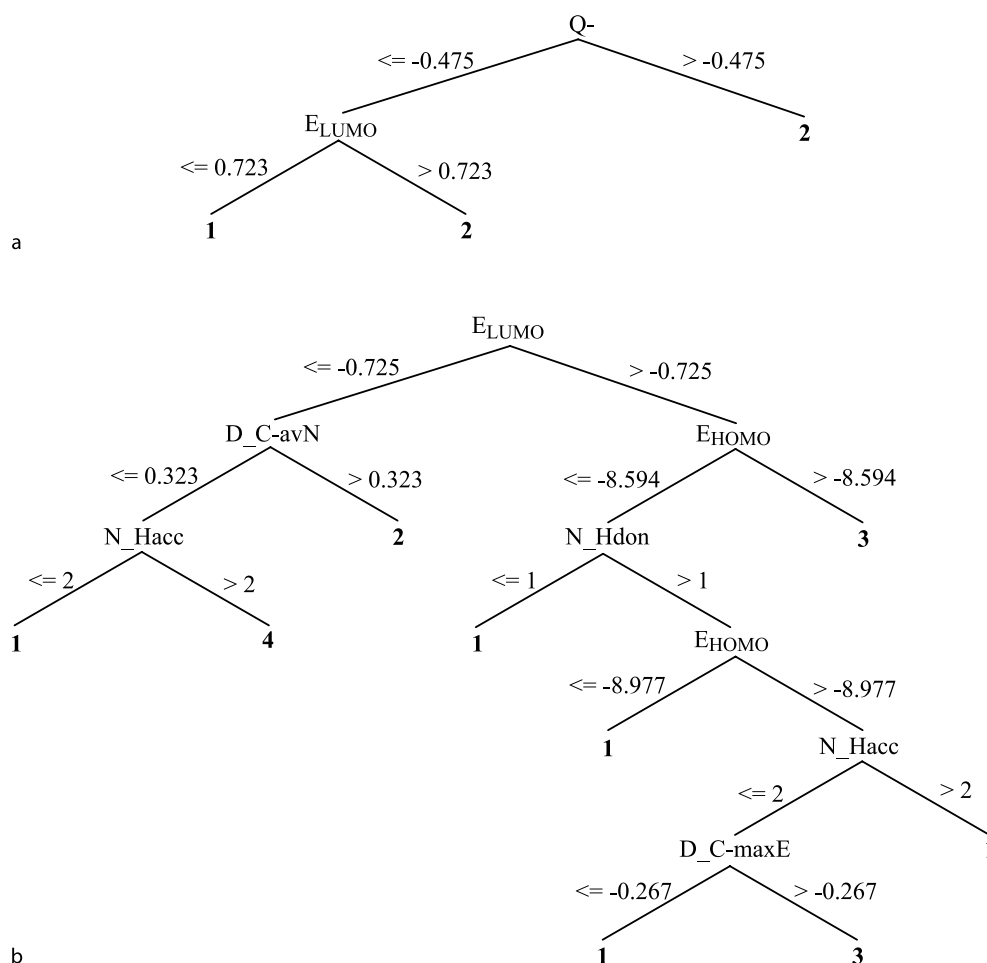
558 A decision tree represents a series of rules that perform  
559 a mapping of the structural descriptors to a prediction for  
560 a molecular property. The class of decision tree algorithms  
561 includes the C4.5 decision tree [108], the NBTree decision  
562 tree with naïve Bayes classifiers at the leaves [82], the alter-  
563 nating decision tree ADTree [42], random trees and ran-  
564 dom forests [22]. A typical decision tree algorithm, such  
565 as C4.5 and related methods, generates a sequence of rules  
566 based on splitting the objects into two subgroups based  
567 on a selected structural descriptor. A threshold of the de-  
568 scriptor separates the objects into a subgroup of objects  
569 that have a descriptor value lower than the threshold, and  
570 another subgroup of objects that have a descriptor value  
571 higher than the threshold. The descriptor and the thresh-  
572 old are selected such as to maximize the difference in en-  
573 tropy, which maximizes the separation of objects based on  
574 their class, i. e., one subgroup contains mainly object from  
575 one class whereas the second subgroup contains mainly  
576 objects from the other class. The process is repeated un-  
577 til a proper separation in classes is obtained for all objects.

578 The structure of a C4.5 decision tree is demon-  
579 strated for the classification of narcotic pollutants  
580 into polar chemicals (class 1) and nonpolar chemicals  
581 (class 2) [66,109,137]. The dataset consists of 190 com-  
582 pounds, from which 76 are polar and 114 are nonpo-  
583 lar. The structural descriptors are the octanol-water par-  
584 tition coefficient  $\log K_{ow}$ , the energy of the highest oc-  
585 cupied molecular orbital  $E_{HOMO}$ , the energy of the low-  
586 est unoccupied molecular orbital  $E_{LUMO}$ , the most neg-

587 ative partial charge on any non-hydrogen atom in the  
588 molecule  $Q-$ , and the most positive partial charge on  
589 a hydrogen atom  $Q+$ . The J48 (C4.5) decision trees com-  
590 puted with Weka [41,151] shows that only two descrip-  
591 tors, namely  $Q-$  and  $E_{LUMO}$ , can separate the leaning  
592 dataset (Fig. 2a). The first rule of this decision tree sep-  
593 arates molecules based on their  $Q-$  value, namely those  
594 with  $Q- \leq -0.475$  are sent to the second rule whereas  
595 those with  $Q- > -0.475$  are classified in class 2. The sec-  
596 ond rule separates chemicals according to their  $E_{LUMO}$   
597 value, namely those with  $E_{LUMO} \leq 0.723$  are classified in  
598 class 1 whereas those with  $E_{LUMO} > 0.723$  are classified in  
599 class 2.

600 The second example considers a J48 (C4.5) decision  
601 trees with four classes (Fig. 2b) representing four modes  
602 of toxic action of phenols in the *Tetrahymena pyriformis*  
603 assay, namely polar narcotics (class 1), oxidative uncou-  
604 plers (class 2), proelectrophiles (class 3), and soft elec-  
605 trophilic (class 4). The dataset consists of 220 chemicals  
606 and seven structural descriptors, namely  $\log K_{ow}$ ,  $E_{HOMO}$ ,  
607  $E_{LUMO}$ , the maximum donor (electrophilic) delocalizabil-  
608 ity for C atoms  $D\_C\text{-maxE}$ , the average acceptor (nucleo-  
609 philic) delocalizability for C atoms  $D\_C\text{-avN}$ , the hydro-  
610 gen bond donor count  $N\_Hdon$ , and the hydrogen bond  
611 acceptor count  $N\_Hacc$ . Out of the seven descriptors, only  
612  $\log K_{ow}$  is not selected to form a splitting rule, whereas  
613  $E_{HOMO}$  and  $N\_Hacc$  are selected in two rules each. Deci-  
614 sion trees are efficient ML algorithms that provide accurate  
615 and fast predictions, with the added benefit of performing  
616 also a selection of the most useful descriptors.

617 Combinatorial chemistry and high-throughput  
618 screening accelerate the drug discovery process by pro-  
619 ducing vast quantities of biological assay results. The data  
620 mining of all these chemical and biological experiments  
621 requires faster structure-activity algorithms, as shown  
622 by Rusinko et al for a library of monoamine oxidase  
623 inhibitors that was evaluated with recursive partition-  
624 ing [112]. The advantage of RP is that it scales linearly with  
625 the number of descriptors, and it can be used for datasets  
626 containing more than  $10^5$  chemicals and  $2 \times 10^6$  struc-  
627 tural descriptors. To generate chemical libraries focused  
628 on a specific target, Deng et al developed the structural  
629 interaction fingerprints, a class of descriptors that encode  
630 ligand-target binding interactions [31]. These fingerprints  
631 were used in a decision tree model to identify ligands for  
632 MAP kinase p38, showing that the method can be used  
633 with success for the structure-based screening of combi-  
634 natorial chemical libraries. Statistical studies of FDA ap-  
635 proved drugs show that there are certain physico-chemical  
636 properties that define drug-like structural requirements,  
637 independent of the drug targets. Drug-like filters are used

**Drug Design with Machine Learning, Figure 2**

Classification with J48 (C4.5) decision trees: **a** polar (class 1) and nonpolar (class 2) narcotic pollutants; **b** toxicity modes of phenols against the *Tetrahymena pyriformis*, namely polar narcotics (class 1), oxidative uncouplers (class 2), proelectrophiles (class 3), and soft electrophiles (class 4)

638 to eliminate at early stages those compounds that do not  
 639 have the general structural features of a drug. Schneider  
 640 et al applied decision trees to the classification of drug-like  
 641 compounds based on drug-like indices, hydrophobicity,  
 642 and molar refractivity [114]. This simple filter correctly  
 643 identifies 83% of drugs and 39% of the non-drugs. The  
 644 classification trees suggest several properties that sepa-  
 645 rate drugs from non-drugs, such as a molecular weight  
 646 higher than 230, a molar refractivity higher than 40, as  
 647 well as the presence of structural elements such as rings  
 648 and functional groups. Natural products are frequently  
 649 screened to identify active compounds for selected tar-  
 650 gets. A database of 240 Chinese herbs containing 8264  
 651 compounds was screened with a random forest algorithm  
 652 to identify inhibitors for several targets, including cy-  
 653 clooxygenases, lipoxygenases, aldose reductase, and three

HIV targets [36]. The screening results show that random  
 forests give dependable predictions even for unbalanced  
 libraries in which inactive compounds are more numerous  
 than active compounds. A literature search confirmed the  
 computational predictions for 83 inhibitors identified in  
 Chinese herbs.

Torsade de pointes (TdP) is a polymorphic ventricular  
 arrhythmia that may be caused by drugs that induce the  
 prolongation of the QT interval by inhibiting the heart  
 potassium channel hERG (human ether-á-go-go). Gepp  
 and Hutter investigated the TdP potential of 339 drugs,  
 and found that a decision tree has a success rate of up to  
 80% in predicting the correct class [46]. Ekins et al com-  
 bined a recursive partitioning (RP) tree with Kohonen and  
 Sammon maps in order to identify hERG inhibitors [38].  
 RP is used as a fast and accurate filter to screen and prior-

654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669



670 itize drug databases for an in-depth assessment with Ko-  
671 honen and Sammon maps. The human intestinal absorp-  
672 tion of potential drugs may be investigated with in silico  
673 methods that assess the intestinal permeability of a chem-  
674 ical compound before synthesis. Hou et al used recursive  
675 partitioning to classify the intestinal absorption (poor or  
676 good) of chemical compounds [57]. The SAR model based  
677 on decision tree classification has good results for a train-  
678 ing set of 481 compounds (95.9% success rate for the poor  
679 absorption class and 96.1% for the good absorption class)  
680 and a prediction set of 98 chemicals (100% success rate for  
681 the poor absorption class and 96.8% for the good absorp-  
682 tion class).

683 All drugs that have a CNS (central nervous system) ac-  
684 tion must penetrate the blood-brain barrier (BBB), and  
685 BBB permeability models are used to filter compounds  
686 that cannot pass the barrier. Andres and Hutter found that  
687 a decision tree provides fast and accurate results, with an  
688 accuracy of 96% for 186 training compounds and 84% for  
689 38 test chemicals [3]. The drugs that pass the BBB are pre-  
690 dicted with a higher accuracy (94%) than those that do not  
691 pass the BBB (89%). Deconinck et al predicted BBB per-  
692 meability with classification and regression trees (CART)  
693 alone and aggregated in a boosting approach [30]. The  
694 training was performed for 147 drugs, and the structural  
695 descriptors were computed with Dragon [133]. The per-  
696 centage of correctly classified drugs in cross-validation is  
697 83.6% for a single tree, and increases to 94% for a boosting  
698 model that contains 150 trees. The single tree model may  
699 be used for a fast screening of large libraries, but for more  
700 reliable predictions a boosting model should be preferred.

701 The drug affinity for the cytochrome P450 2C9 was  
702 predicted with a consensus method based on four ML al-  
703 gorithms [59]. Two of these algorithms are variants of  
704 recursive partitioning, the third is a local interpolation  
705 method, and the fourth is a subgraph search method.  
706 These four models were trained with a set of 276 com-  
707 pounds, and then the SAR models are tested for a pre-  
708 diction set of 50 chemicals. The consensus predictions  
709 have an accuracy of 94%, demonstrating that the com-  
710 putational methods may predict which drugs are metabo-  
711 lized by P450 2C9. The metabolic stability of 161 drugs  
712 against six isoforms of the cytochrome P450 (1A2, 2C9,  
713 2C19, 2D6, 2E1, and 3A4) was evaluated with a clas-  
714 sification tree [154]. The SAR models several structural  
715 characteristics that determine the P450 specificity, namely  
716 3A4 substrates are larger compounds, 2E1 are smaller  
717 molecules, 2C9 substrates are anionic, and 2D6 substrates  
718 are cationic.

719 *P*-glycoprotein (Pgp) mediated drug efflux is responsi-  
720 ble for the low cellular accumulation of anticancer drugs,

721 for reduced oral absorption, and for low blood-brain bar-  
722 rier penetration. Pgp affects also the hepatic, renal, or  
723 intestinal elimination of drugs. Li et al. proposed a  
724 decision tree model to differentiate Pgp substrates from  
725 non-substrates [91]. Four-point 3D pharmacophores ob-  
726 tained from the three-dimensional structure of 163 chem-  
727 ical compounds were subsequently used as descriptors  
728 for the classification tree. Nine pharmacophores were se-  
729 lected in the decision tree, giving an accuracy of 87.7% for  
730 the training set and 87.6% for the prediction set. These  
731 pharmacophores highlight the structural features of Pgp  
732 substrates, and can be used as a fast filter for chemical  
733 libraries. The renal clearance of 130 drugs was investi-  
734 gated with Molconn-Z topological indices [76] and re-  
735 cursive partitioning [32]. RP results for the separation  
736 of high-clearance compounds from low-clearance com-  
737 pounds show 88% correct predictions for the training set  
738 of 130 compounds and 75% correct predictions for a pre-  
739 diction set of 20 compounds.

740 The prediction power of classification and regression  
741 algorithms can be significantly improved by using an en-  
742 semble of models that are combined into a final predic-  
743 tion by an averaging or a voting procedure. Such ML al-  
744 gorithms are called ensemble, committee, or jury meth-  
745 ods. Tong et al proposed an ensemble method, decision  
746 forest (DF), which is based on classification and regres-  
747 sion trees [134]. A DF is obtained in four steps: (a) gen-  
748 erate a CART; (b) generate a second CART based only  
749 on structural descriptors that were not used in the first  
750 CART; (c) repeat the first two steps until no more trees can  
751 be generated; (d) predict a property of a chemical com-  
752 pound based on all trees generated. DF was used to pre-  
753 dict the estrogen receptor binding activity of organic com-  
754 pounds from Molconn-Z topological indices [135]. After  
755 2000 simulations of a 10-fold cross-validation, the pre-  
756 dicted accuracy for a dataset of 232 compounds was 81.9%,  
757 and the accuracy for another dataset of 1092 chemicals was  
758 79.7%. The DF models can be used to identify potential  
759 endocrine disruptors, which are chemicals that affect the  
760 endocrine functions. A similar algorithm is the random  
761 forest (RF) [22], which was extensively tested for several  
762 structure-activity studies [128]. RF has three properties  
763 that are important in SAR, namely the algorithm evalu-  
764 ates the importance of descriptors, identifies similar com-  
765 pounds, and measures the prediction performance. The  
766 skin sensitization activity of 131 chemicals was modeled  
767 with RF, using Dragon and Molconn-Z indices [90]. Com-  
768 pared with single tree models, RF improve significantly the  
769 ability to identify chemical allergens.

### 770 Lazy Learning and $k$ -Nearest Neighbors

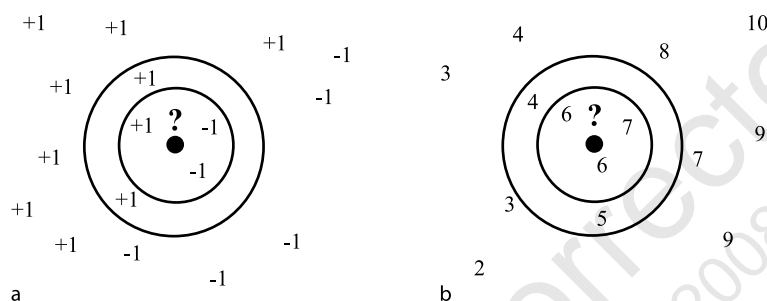
771 Lazy learning is a local learning method that keeps in the  
772 memory the entire set of training objects, and computed  
773 an interpolation model based on the neighborhood on the  
774 prediction object Atkeson, 1997 #1955 [6]. The prevalent  
775 lazy learning algorithm in SAR and QSAR is the  $k$ -nearest  
776 neighbors [1]. A  $k$ -NN prediction is based on a descriptor  
777 space, which may be the entire set of structural descrip-  
778 tors evaluated in a study or a subset that gives better pre-  
779 dictions, a distance function, and a neighborhood  $k$  of the  
780 prediction object. The  $k$ -NN classification may be applied  
781 to two or more classes (Fig. 3a). The prediction object,  
782 shown here as a black circle with an unknown class “?”,  
783 is inserted into the descriptor space and then its  $k$  nearest  
784 neighbors are identified. The class of the prediction object  
785 is the class of the majority of its neighbors. In the example  
786 from Fig. 3a, the predicted class is  $-1$  if  $k = 3$ , and  $+1$  if  
787  $k = 5$ . The prediction statistics of the  $k$ -NN classifier de-  
788 pend on the value of  $k$ , and its optimum value is usually  
789 determined by cross-validation.

790 A similar procedure is used for  $k$ -NN regression  
791 (Fig. 3b). For regression, each training object has a real  
792 number as property value, and the predicted property for  
793 the query object is an average value of the property values  
794 for its  $k$  nearest neighbors. For the situation depicted in  
795 Fig. 3b, the predicted value for  $k = 3$  is  $(6+6+7)/3 = 6.3$ ,  
796 whereas the prediction for  $k = 5$  is 5.6. A distance weight  
797 is usually added to increase the contribution of closer ob-  
798 jects and to decrease the influence of more distant objects.  
799 For a distance function  $d$ , a possible weighting of the prop-  
800 erty values is the inverse distance  $1/d$ .

801 Lazy learning applications in drug design are still rare,  
802 mainly because the community of practitioners only re-  
803 cently discovered the rich field of machine learning algo-  
804 rithms. Kumar et al used locally linear embedding for non-  
805 linear dimensionality reduction coupled with lazy learn-  
806 ing for two QSAR benchmarks, the Selwood dataset and

807 the steroids dataset originally analyzed with CoMFA [85].  
808 Guha et al applied local lazy regression to three QSAR  
809 models, namely to 179 artemisinin analogues with anti-  
810 malarial activity, to 79 platelet-derived growth factor in-  
811 hibitors, and to 756 inhibitors of dihydrofolate reduc-  
812 tase [47]. In all cases, the local models give better pre-  
813 dictions compared to global regression models. An au-  
814 tomated lazy learning quantitative structure-activity rela-  
815 tionship (ALL-QSAR) approach was developed based on  
816 locally weighted linear regression models computed from  
817 the compounds in training set that are chemically most  
818 similar to a test compound [159]. ALL-QSAR was applied  
819 with good results for three datasets, namely 48 anticonvul-  
820 sant compounds, 48 dopamine receptor antagonists, and  
821 to model the toxicity of 250 phenols against *Tetrahymena*  
822 *pyriformis*. Sommer and Kramer reported several appli-  
823 cations of lazy learning to chemical datasets with a high  
824 structural diversity [123]. Collections of noncongeneric  
825 compounds are usually difficult to model, but their ap-  
826 proach obtained good results for several datasets, namely  
827 739 mutagenic chemicals from the Carcinogenic Potency  
828 Database, 342 chemicals tested for biodegradability, and  
829 1481 from the AIDS Antiviral Screen of the NCI Develop-  
830 mental Therapeutics Program.

831 The fundamental assumptions in  $k$ -NN applications  
832 to SAR and QSAR are that molecules with similar struc-  
833 tures are close in the descriptor space, and that similar  
834 molecules have similar properties. Basak advocated the  
835 use of  $k$ -NN in QSAR as an alternative to the multiple  
836 linear regression [10,11]. In this approach, the chemical  
837 compounds are represented with a collection of topologi-  
838 cal indices, and the Euclidean distance is used to identify  
839 the structures most similar with the test molecule. This  
840  $k$ -NN method was applied with success to model the hy-  
841 drophobicity of 4067 compounds [10], the mutagenicity  
842 of a diverse set of 520 chemicals [11], the boiling tem-  
843 perature of 139 hydrocarbons and the mutagenicity of 95



Drug Design with Machine Learning, Figure 3

Applications of the  $k$ -NN algorithm to a classification and b regression

aromatic and heteroaromatic amines [9]. Good models are obtained with optimized  $k$  values over a range between 1 and 25, as shown in QSAR model for vapor pressure [48]. The descriptors selected in the structural space are essential in establishing a suitable similarity measure. General collections of descriptors usually give good predictions, but property tailored structural spaces can increase significantly the predictivity of the model. Tailored structural spaces include only those descriptors that contribute significantly to the model, as demonstrated for hydrophobicity and boiling temperature [49]. A structural space consisting of electrotopological state indices gives reliable predictions for hydrophobicity, Henry's law constant, vapor pressure, and OH-radical bimolecular rate constant [24].

Tropsha applied the  $k$ -NN to numerous QSAR models, demonstrating its value in comparisons with more established algorithms. For a dataset of 29 dopamine antagonists,  $k$ -NN predictions are much better than those obtained with CoMFA, which represents the standard 3D QSAR model [55].  $k$ -NN may be used also as a robust method to select those structural descriptors that are important in modeling a particular biological property, as shown for 58 estrogen receptor ligands [161] and for 48 amino acid derivatives with anticonvulsant activity [120]. In a more elaborate implementation,  $k$ -NN is used to select descriptors, to find the best value for  $k$ , and to optimize the function that weights the contribution of its neighbors [60]. The anticancer activity of 157 epipodophylotoxin derivatives was modeled with molecular connectivity indices and  $k$ -NN [152]. In a comparison with a CoMFA model obtained for the same dataset it was found that  $k$ -NN gives better predictions in a cross-validation test. The metabolic stability of 631 drug candidates was evaluated with  $k$ -NN and a collection of structural descriptors that included molecular connectivity indices and atom pairs [121]. The prediction in a test set of compounds had accuracy higher than 85%, whereas an external validation set of 107 chemicals was predicted with 83% success rate.

$k$ -NN is also used as a component in more elaborate drug design algorithms. Ajmani et al combined molecular field analysis and  $k$ -NN to obtain a robust QSAR model that was tested for three datasets, namely the CoMFA steroids, cyclooxygenase-2 inhibitors, and anticancer compounds [2]. Another QSAR model that incorporates  $k$ -NN is Complimentary Ligands Based on Receptor Information (CoLiBRI) [103]. CoLiBRI represents both ligands and receptor binding sites in the same space of universal chemical descriptors. In a test consisting of 800 X-ray protein-ligand receptors from PDB, CoLiBRI ranked the correct ligands in the top 1% chemicals selected, and was able to quickly eliminate improbable lig-

ands. The skin permeability coefficients of 110 compounds were predicted with an ensemble consisting of a collection of  $k$ -NN and ridge regression models [99].

The melting temperature of organic compounds is particularly difficult to model, due mainly to a lack of good descriptors for the interactions in the solid and liquid states. Nigsch et al used several types of molecular descriptors to model the melting temperature of 4119 diverse organic molecules with  $k$ -NN [101]. The contribution of selected neighbors was evaluated with four weighting functions (arithmetic and geometric average, inverse distance weighting, and exponential weighting), with better results obtained by using the exponential weighting scheme. The optimized model has an average error of 42.2 °C. In a comparison between global and local ( $k$ -NN) regression models for the blood-brain distribution of 328 chemicals, it was found that  $k$ -NN consistently gives better predictions as measured in cross-validation [84].

### Bayesian Methods

Although Bayesian classifiers were only recently applied to drug design and structure-activity studies, they have several characteristics that make them particularly useful in screening large chemical libraries. Bayesian classifiers give dependable predictions, can tolerate noise and errors in experimental data, and can be computed fast compared to other ML algorithms, which is particularly important in screening large collections of molecules. Bayes' theorem describes the relationship between the prior and conditional probabilities of two events  $A$  and  $B$ , when  $B$  has a non-zero probability [12]:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)},$$

where the terms have the following meaning:  $P(A)$  is the prior or marginal probability of  $A$ , and does not consider any information regarding  $B$ ;  $P(B)$  is the prior (marginal) probability of  $B$ ;  $P(A|B)$  is the conditional (posterior) probability of  $A$  given the occurrence of event  $B$ ;  $P(B|A)$  is the conditional (posterior) probability of  $B$  given the occurrence of event  $A$ . Bayesian classifiers [71,151] evaluate a set of descriptors to predict the activity of a chemical based on the known descriptors and activities of a training set of chemicals. Each descriptor is considered to be statistically independent of all other descriptors. Based on a selected descriptor, the probability that a compound is active is proportional to the ratio of active to inactive compounds that have the same value for that descriptor. The final prediction is obtained by multiplying the descriptor-based probabilities.

941 A variety of drug discovery and development problems can be solved efficiently with Bayesian classifiers. 942  
943 The virtual screening of large chemical libraries represents a major application of Bayesian classifiers, mainly 944  
945 due to their tolerance to noise in experimental data [146]. 946  
947 Klön and Diller proposed a method to prioritize chemicals for synthesis and screening based on physico-chemical 948  
949 properties [77]. It was found that the Bayesian classifier is superior to the usual pairwise Tanimoto similarity 950  
951 filter. High-throughput docking is an essential tool for structure-based drug design, but the current scoring functions 952  
953 do not provide an accurate ranking of the molecules. Various improvements to the scoring functions were proposed, 954  
955 including Bayesian classifiers [78,79]. Simulations with two biological targets showed that a Bayesian classifier 956  
957 increases the enrichment in active compounds. Finding the potential targets of chemicals is increasingly relevant 958  
959 in drug discovery. Nidhi et al developed a multi-class Bayesian model in which chemicals represented with extended- 960  
961 connectivity fingerprints are evaluated against 964 target classes [100]. The practical utility of the approach 962  
963 is demonstrated by its high success rate, with 77% correct target identification. Adverse drug reactions are investigated 964  
965 during early phases of drug discovery in preclinical safety pharmacology tests for various targets. Multi-class 966  
967 Bayesian classifiers developed for 70 such targets can identify 93% of ligands with a 94% correct classification 968  
969 rate [14]. In a test using chemicals from World Drug Index, the classifier identifies 90% of the adverse drug 970  
971 reactions with a 92% correct classification rate. Adverse reactions of several drugs withdrawn from the market were 972  
973 also predicted with good accuracy.

974 Watson used 2D pharmacophore feature triplet vectors and Bayesian classifiers to find active compounds in 975  
976 library screening [150]. Then good enrichment in active compounds obtained for several chemogenomic databases 977  
978 demonstrates the practical utility of this approach. Natural products represent an excellent source of starting 979  
980 structures for the design of novel drugs. Ertl trained a Bayesian classifier to provide the natural product-likeness 981  
982 score for a chemical [39]. This score measures the similarity between a molecule and the structural space covered by 983  
984 natural products, and can separate between natural products and synthetic molecules. Usual similarity search in 985  
986 chemical libraries is based on a set of structural descriptors and a similarity (or distance) measure. Bender et al 987  
988 showed that the retrieval rate of active compounds increases considerably by using Bayes affinity fingerprints that 989  
990 describe the ligand bioactivity space [13]. In this space, Bayes scores for ligands from about 1000 activity classes 991  
992 describe the chemical structures. This approach was able to

992 improve with 24% the results from a standard similarity screening. 993

994 Bayesian methods are also used in various structure-activity models. Klön et al demonstrated several applica- 995  
996 tions of modified Bayesian classifiers that model continuous numerical data with a Gaussian distribution [80]. 997  
998 Several models developed for absorption, distribution, metabolism and excretion property prediction show that 999  
1000 the new classifiers have better performance compared to classical Bayesian classifiers. Multidrug resistance represents 1001  
1002 the ability of cancer cells to become simultaneously resistant to several drugs. A possible solution to this problem 1003  
1004 is the use of multidrug resistance reversal (MDDR) agents. Sun found that a Bayesian classifier based on atom 1005  
1006 types could be used to identify MDDR agents in a set of 424 training compounds [125]. The prediction of MDDR 1007  
1008 agents in a group of 185 test compounds had a success rate of 82.2%. Phospholipidosis is an adverse effect that is 1009  
1010 investigated during preclinical studies in animals, and may stop the development of an otherwise promising drug 1011  
1012 candidate. Pelletier et al proposed an *in silico* Bayesian classifier to identify phospholipidosis as a viable alternative to 1013  
1014 more expensive and time-consuming *in vivo* tests [105]. Based on two simple descriptors,  $pK_a$  and ClogP, the 1015  
1016 model has a success rate of 83% for a dataset of 201 compounds. UGT is an enzyme family that catalyzes the 1017  
1018 reaction between glucuronic acid and a chemical that has a nucleophilic group. Drug glucuronidation represents an 1019  
1020 important mechanism for their elimination, but the selectivity of different UGT isoforms is difficult to predict. 1021  
1022 The Bayesian classifier suggested by Sorich et al predicts the glucuronidation site for eight UGT isoforms based on 1023  
1024 the partial charge and Fukui function of the nucleophilic atom and the presence of an aromatic ring attached to the 1025  
1026 nucleophilic atom [124]. The models have a good predictive ability, with a sensitivity and specificity in the range 1027  
1028 75–80%. The screening test to identify bitter compounds presented by Rodgers et al uses selected substructural 1029  
1030 features selected from 649 bitter chemicals and 13,530 random drug-like compounds [111]. A Bayesian classifier 1031  
1032 trained with these data can predict correctly 72.1% of the bitter compounds. The bitter taste is determined mainly 1033  
1034 by substructures representing sugar moieties and highly branched carbon skeletons. 1035

### 1036 Support Vector Machines

1037 Support vector machines represent a major development for SAR and QSAR models, as suggested by the large 1038  
1039 number of publications that apply SVM and related kernel methods to drug design. To better understand the math- 1040

emathical basis of SVM and the parameters that influence their results, we start this section with a brief theoretical presentation of SVM for classification and regression. Several kernel functions are presented together with SVM plots that demonstrate the complex shapes that may be simulated with these functions. The influence of various kernels on QSAR predictions is shown for a dataset of benzodiazepine receptor ligands, and several applications in drug design are reviewed.

### Hard Margin Linear SVM Classification

The most simple case of SVM is the classification of two classes of objects that may be separated by a hyperplane (Fig. 4a). It is obvious from this figure that there are many hyperplanes that can separate the two classes of objects, but not all of them give good predictions for novel objects. Intuitively, a hyperplane that maximizes the separation of the two classes, like hyperplane H, is expected to offer the best predictions. The SVM algorithm determines a unique hyperplane H that has the maximum margin, i. e., the maximum distance between hyperplanes H<sub>1</sub> and H<sub>2</sub>. Hyperplane H<sub>1</sub>, which is determined by three objects from class +1 that are shown inside circles, defines the border with class +1. Hyperplane H<sub>2</sub>, which is determined by two objects from class -1 that are shown also inside circles, defines the border with class -1. The objects represented inside circles are called support vectors and they determine the unique solution for the SVM model. By removing all other objects and keeping only the support vectors one obtains the same SVM model. The hyperplane H is defined by  $\mathbf{w} \cdot \mathbf{x} + b = 0$ , where  $\mathbf{w}$  is the normal vector to H, the hyperplane H<sub>1</sub> that borders class +1 is defined by  $\mathbf{w} \cdot \mathbf{x} + b = +1$ , and the hyperplane H<sub>2</sub> that borders class

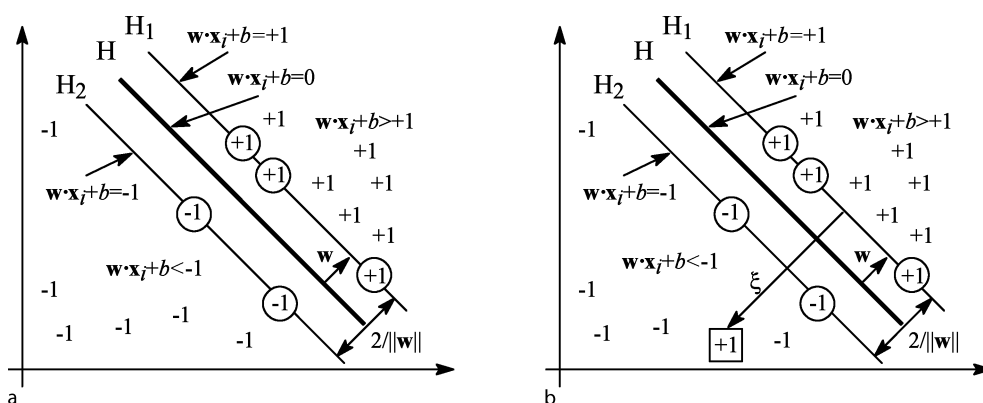
-1 is defined by  $\mathbf{w} \cdot \mathbf{x} + b = -1$ . The distance between the origin and the hyperplane H is  $|b|/\|\mathbf{w}\|$ , and the SVM margin (the distance between hyperplanes H<sub>1</sub> and H<sub>2</sub>) is  $2/\|\mathbf{w}\|$ . A molecule  $i$  characterized by a vector of structural descriptors  $\mathbf{x}_i$  belongs to class +1 if  $\mathbf{w} \cdot \mathbf{x}_i + b \geq +1$ , or to class -1 if  $\mathbf{w} \cdot \mathbf{x}_i + b \leq -1$ . Obviously, support vectors from class +1 are determined by  $\mathbf{w} \cdot \mathbf{x}_i + b = +1$ , whereas support vectors from class -1 are determined by  $\mathbf{w} \cdot \mathbf{x}_i + b = -1$ .

Based on the formula of the SVM margin, it follows that the maximum separation hyperplane is obtained by maximizing  $2/\|\mathbf{w}\|$ , which is equivalent to minimizing  $\|\mathbf{w}\|^2/2$ . The linear SVM model (Fig. 4a) is formulated as:

$$\begin{aligned} \text{minimize } f(\mathbf{x}) &= \frac{\|\mathbf{w}\|^2}{2} \\ \text{with the constraints } g_i(\mathbf{x}) &= y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0, \\ & i = 1, \dots, m \end{aligned} \quad (2)$$

where  $y_i$  is the class (+1 or -1) of the object  $i$ . This optimization equation is a quadratic programming that is further transformed with a Lagrangian function into its dual formulation. All SVM formulations are solved for the dual formulation, because the solution obtained for this simple case is easily extended to more complex situations. The minimization problem from Eq. (2) is expressed with a Lagrangian function:

$$\begin{aligned} L_P(\mathbf{w}, b, \mathbf{A}) &= f(\mathbf{x}) + \sum_{i=0}^m \alpha_i g_i(\mathbf{x}) \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) \end{aligned}$$



Drug Design with Machine Learning, Figure 4

SVM classification models: **a** linearly separable data; **b** linearly non-separable data

$$\begin{aligned}
 &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i y_i (\mathbf{w} \cdot \mathbf{x}_i + b) + \sum_{i=1}^m \alpha_i \\
 &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i y_i \mathbf{w} \cdot \mathbf{x}_i \\
 &\quad - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i, \quad (3)
 \end{aligned}$$

where  $\mathbf{A} = (\alpha_1, \alpha_2, \dots, \alpha_m)$  is the set of Lagrange multipliers for the objects that have  $\alpha_i \geq 0$ , and  $P$  in  $L_P$  indicates the primal formulation of the problem. The Lagrangian function  $L_P$  must be minimized with respect to  $\mathbf{w}$  and  $b$ , and maximized with respect to  $\alpha_i$ , subject to the constraints  $\alpha_i \geq 0$ . The dual optimization problem  $L_D$  is:

maximize

$$L_D(\mathbf{w}, b, \mathbf{A}) = \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

subject to  $\alpha_i \geq 0$ ,  $i = 1, \dots, m$

$$\text{and } \sum_{i=1}^m \alpha_i y_i = 0.$$

The optimization problem is usually solved with the sequential minimal optimization (SMO) proposed by Platt [106]. For a complete formulation of the SVM solution see the books and reviews mentioned in Introduction [70]. The training objects with  $\alpha > 0$  constitute the support vectors, whereas the objects with  $\alpha = 0$  may be removed from the learning set without affecting the SVM model. The vector  $\mathbf{w}$  is obtained as:

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

and  $b$  is computed as the average value for all support vectors. The SVM model is the optimum separation hyperplane  $(\mathbf{w}, b)$  that can now be used to predict the class membership for new objects. The class of a molecule  $k$  with the structural descriptors  $\mathbf{x}_k$  is:

$$\text{class}(k) = \begin{cases} +1 & \text{if } \mathbf{w} \cdot \mathbf{x}_k + b > 0 \\ -1 & \text{if } \mathbf{w} \cdot \mathbf{x}_k + b < 0. \end{cases} \quad (6)$$

The classifier is further simplified by substituting  $\mathbf{w}$  with its expression (5):

$$\text{class}(k) = \text{sign} \left( \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_k + b \right), \quad (7)$$

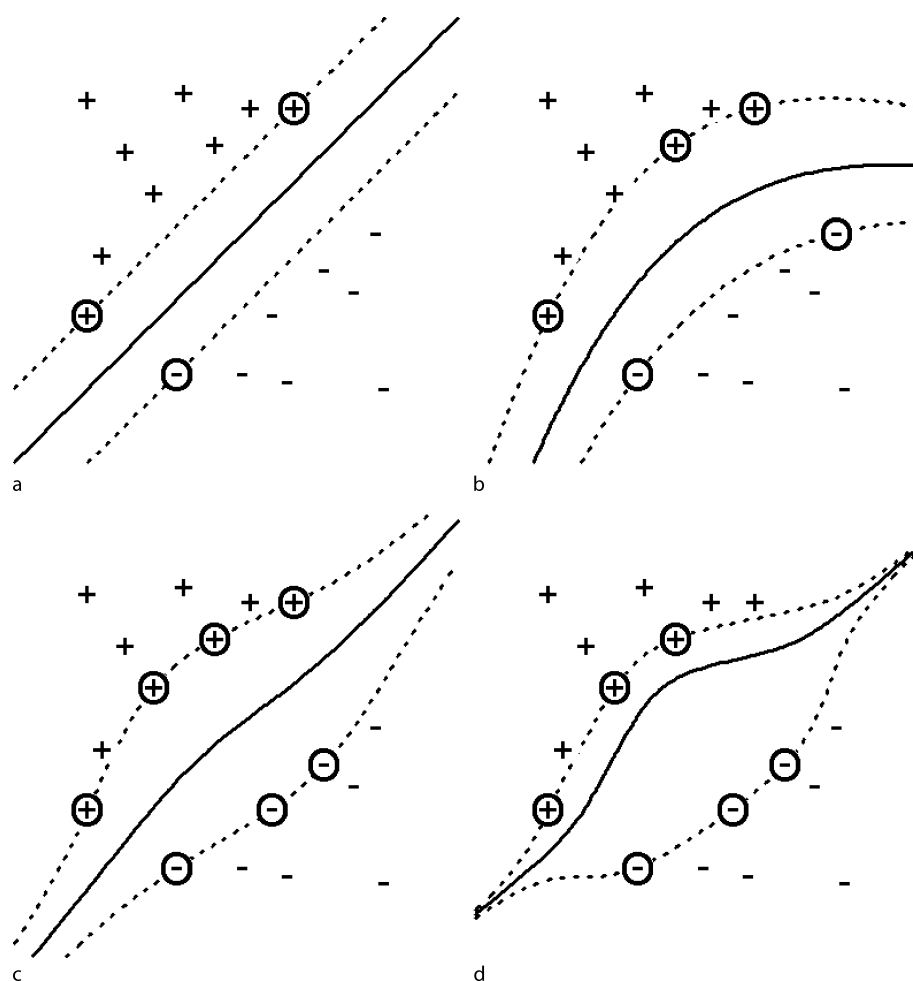
where the summation goes only over the support vectors. The SVM described in the above equations is called a “hard

margin” SVM because it does not allow for classification errors.

The class separation with a linear SVM is simple and intuitive. However, SVM are usually used with nonlinear kernels that have a much more complex shape, which might result in too complicated separation surfaces. To demonstrate the shape of the separation surface for linearly separable data we compare results obtained with a linear kernel and several nonlinear kernels that will be introduced later. All calculations were performed with *R* (<http://www.r-project.org/>) and the *kernlab* package. In all figures, class + 1 patterns are represented by “+” and class -1 patterns are represented by “-”. The SVM hyperplane is depicted with a continuous line, whereas the margins of the SVM hyperplane are shown with dotted lines. Support vectors from class + 1 are represented as “+” inside a circle, and support vectors from the class -1 are depicted as “-” inside a circle.

The linear kernel gives an SVM model that has an optimum separation of the two classes (Fig. 5a) with three support vectors, namely two from class + 1 and one from class -1. The hyperplane has the maximum width, and no object is situated inside the margins (represented with dotted lines). The classification of new objects is made by applying Eq. (7) to the three support vectors. The same dataset is modeled with a degree 2 polynomial kernel (Fig. 5b) in which case the SVM model has five support vectors, namely three for class + 1 and two for class -1. The margin width varies, being larger in the middle and smaller towards the extremes. The hyperplane topology is different from that obtained with a linear kernel, and obviously for some objects the two SVM models will give different predictions. The SVM classifier obtained with a degree 3 polynomial kernel (Fig. 5c) has four support vectors from class + 1 and three support vectors from class -1, and the margin becomes smaller towards the ends. Higher order polynomial kernels produce a complex separation hyperplane, as shown also for a degree 10 polynomial kernel (Fig. 5d), in which case the margin almost vanished towards the extremes.

The Gaussian radial basis function (RBF) kernel is the most common option for SVM, and in Fig. 6 we present four SVM models obtained for different  $\sigma$  values. For low  $\sigma$  values the RBF kernel (Fig. 6a,  $\sigma = 0.01$ ) approximates a linear kernel (Fig. 5a) and the two hyperplanes are very similar. As  $\sigma$  increases the nonlinear behavior of the kernel becomes apparent (Fig. 6b,  $\sigma = 0.1$ ), but the SVM model still has a good separation of the data. As  $\sigma$  increases to 1 (Fig. 6c) and then to 5 (Fig. 6d) the nonlinearity of the SVM model becomes more apparent, and the number of support vectors increases. In the last case all data points are



Drug Design with Machine Learning, Figure 5

SVM classification models for linearly separable data: a dot kernel (*linear*); b polynomial kernel, degree 2; c polynomial kernel, degree 3; d polynomial kernel, degree 10

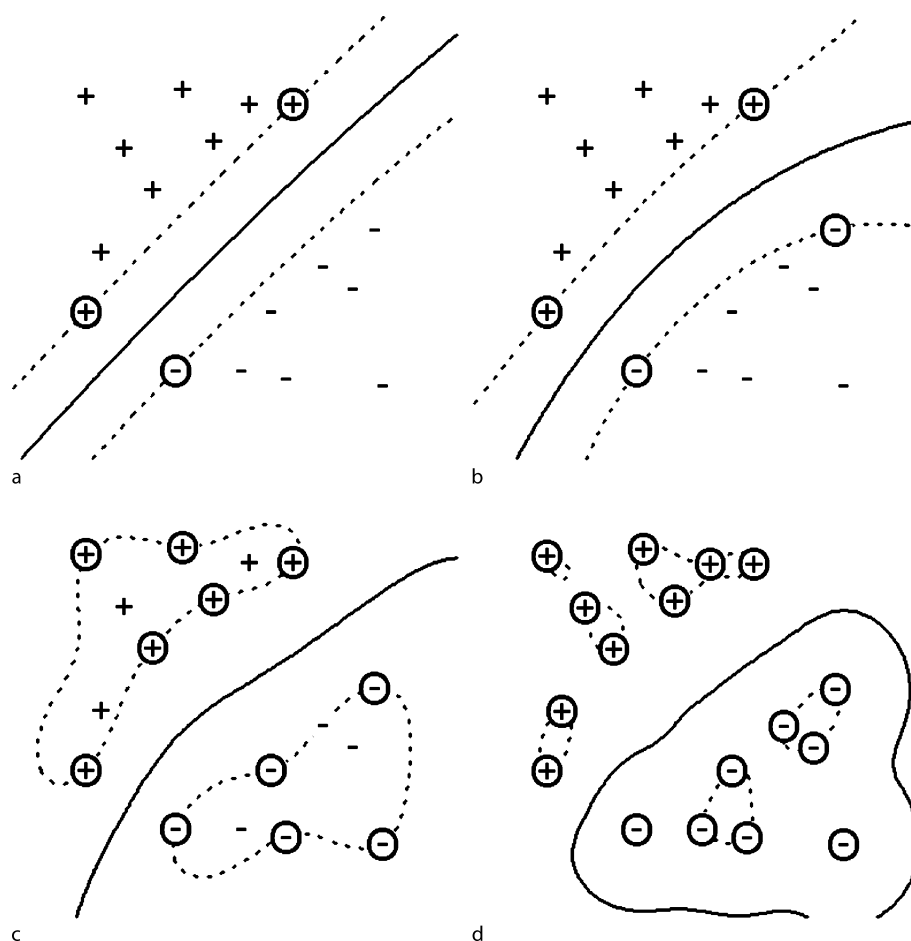
1178 support vectors, and the SVM model has a very complex  
1179 topology. The examples presented here show that the RBF  
1180 kernel may generate too complex hyperplanes that do not  
1181 properly reflect the structure of the data.

1182 Another important kernel is the hyperbolic tangent  
1183 ( $\tanh$ ), which is also used in artificial neural networks.  
1184 Some combinations of parameters result in SVM models  
1185 almost identical with those obtained with the linear kernel  
1186 (Fig. 7a) or similar (Fig. 7b). Higher values for the  
1187 two parameters result in SVM models with smaller margins  
1188 and nonlinear separation of the two classes (Fig. 7c,d).  
1189 The SVM models demonstrated here for linearly separable  
1190 data show that nonlinear kernels, such as RBF and  $\tanh$ ,  
1191 may give SVM models similar to those obtained with linear  
1192 kernels, for certain values of their parameters. However,  
1193 nonlinear kernels may also give SVM models that are

too complex and that do not follow the structure of the  
data. As a general rule, the SVM computed with nonlinear  
kernels should be compared with a linear kernel SVM.

#### Soft Margin Linear SVM Classification

The classification problem solved in the previous section  
considered a two-class situation with classes separated  
with a linear classifier. However, such cases are rare in  
drug design applications, and most often the classes of  
compounds have regions where they superpose. There are  
several reasons why classes of chemicals cannot be separated  
with a linear classifier. The most common problem is the  
identification of the structural descriptors that determine  
the property. There are thousands of structural descriptors  
proposed in the literature, but it is difficult to



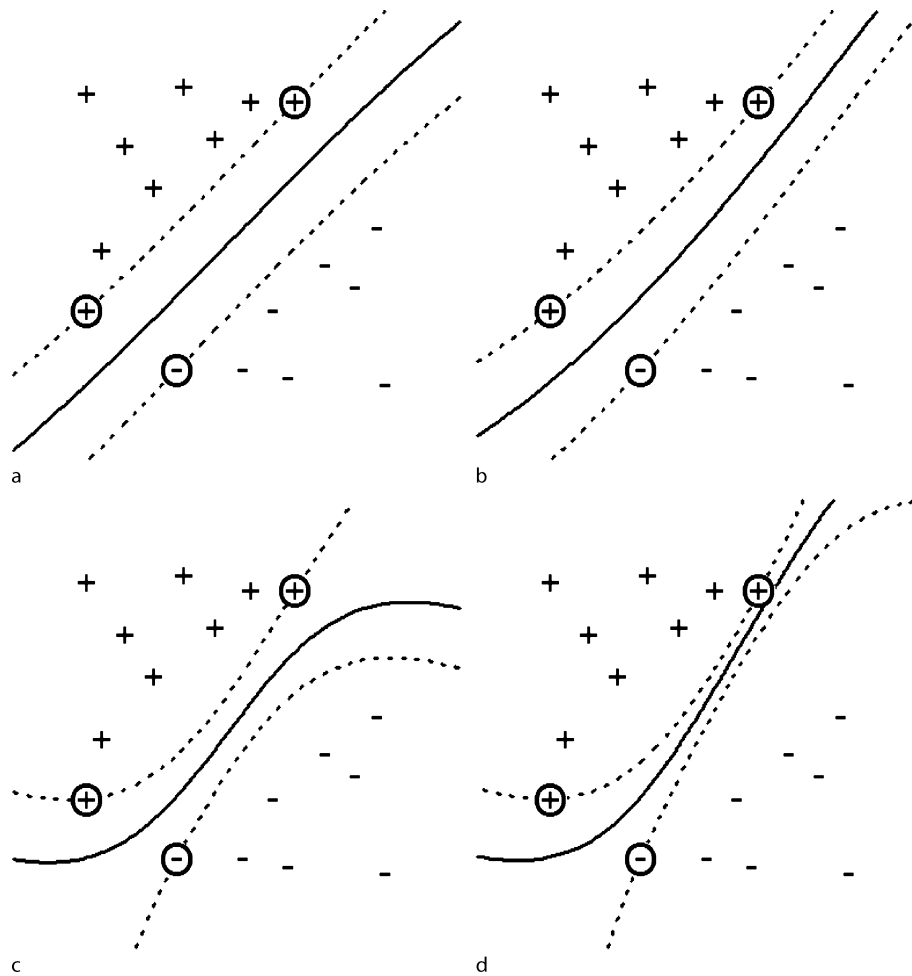
Drug Design with Machine Learning, Figure 6

SVM classification models for linearly separable data with the Gaussian RBF kernel: a  $\sigma = 0.01$ ; b  $\sigma = 0.1$ ; c  $\sigma = 1$ ; d  $\sigma = 5$ 

1208 try them all in a SAR or QSAR study. It is always a possi- 1225  
 1209 bility that the proper type of descriptors is not even discovered. 1226  
 1210 The second problem is to identify the subset of descriptors 1227  
 1211 that give the best predictions, which is a time-consuming 1228  
 1212 task that is usually solved with heuristic algorithms that give a 1229  
 1213 sub-optimal solution. The mapping between the descriptor space 1230  
 1214 and the property might be nonlinear, and obviously a linear classifier 1231  
 1215 fails to model the data. As a final point we have to mention that 1232  
 1216 experimental data may be affected by noise, measurement errors, 1233  
 1217 interference with other biological targets, or unaccounted factors 1234  
 1218 that determine the outcome of the experiments. Although the hard 1235  
 1219 margin linear SVM is limited in its abilities to model real data, the 1236  
 1220 formalism presented in the previous section is the basis for developing 1237  
 1221 more complex SVM models. The linear SVM with soft margin is such 1238  
 1222 an extension that allows classification errors. 1239  
 1223 1240  
 1224 1241

The classification problem shown in Fig. 4b is very similar with the linear separable case from Fig. 4a, because all but one object can be separated with a linear hyperplane  $H$ . The exception is a object from class +1 that is situated in the  $-1$  space region (shown inside a square). A penalty term is introduced for misclassified data, with the property that it is zero for objects correctly, and it has a positive value for objects that are not situated on the correct side of the classifier. The value of the penalty increases when the distance of the corresponding hyperplane increases. The penalty is called a slack variable and it is denoted with  $\xi$ . For the situation depicted in Fig. 4b the penalty for the misclassified object is measured starting from hyperplane  $H_1$  because this hyperplane border the +1 region. For +1 molecules situated in the buffer zone between  $H$  and  $H_1$ , and for  $-1$  molecules situated in the buffer zone between  $H$  and  $H_2$ , the slack variable takes values between 0 and 1.





Drug Design with Machine Learning, Figure 7

SVM classification models for linearly separable data with the hyperbolic tangent (tanh) kernel: **a**  $a = 0.1, b = 0$ ; **b**  $a = 0.1, b = 0.5$ ; **c**  $a = 0.5, b = 0$ ; **d**  $a = 0.5, b = 0.5$

1242 Such patterns are not considered to be misclassified, but  
 1243 have a penalty added to the objective function. The con-  
 1244 straints for the objective function for soft margin linear  
 1245 SVM are:

$$1246 \begin{cases} \mathbf{w} \cdot \mathbf{x}_i + b \geq +1 - \xi_i & \text{if } y_i = +1 \\ \mathbf{w} \cdot \mathbf{x}_i + b \leq -1 + \xi_i & \text{if } y_i = -1 \\ \xi_i > 0, \quad \forall i. \end{cases} \quad (8)$$

1247 The soft margin SVM should balance two opposite con-  
 1248 ditions, namely a maximum margin and a minimum of  
 1249 errors that is equivalent with minimizing a sum of slack  
 1250 variables. The optimization problem for a linear SVM with

classification errors is;

$$\text{minimize } \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^m \xi_i$$

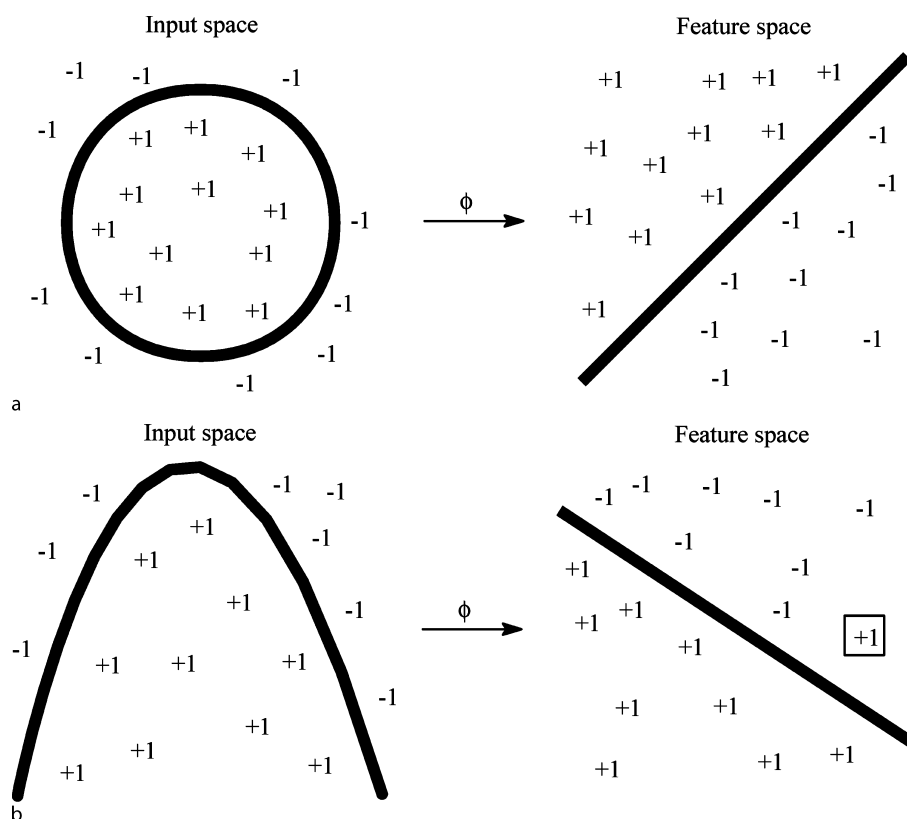
with the constraints

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq +1 - \xi_i, \quad i = 1, \dots, m$$

$$\xi_i \geq 0, \quad i = 1, \dots, m$$

where  $C$  is called capacity and is a parameter that weights  
 the penalty for classification errors. A large value for the  
 capacity parameter  $C$  represents a high penalty for classifi-  
 cation errors, which forces the SVM solution to minimize  
 the number of misclassified objects by reducing the margin  
 of the classifier. A small capacity  $C$  lowers the value of  
 the penalty term and maximizes the margin, which makes

Unauthenticated Proof



Drug Design with Machine Learning, Figure 8

Patterns that are nonlinear separable in the input space are linear separable in the feature space: **a** the class separation is complete; **b** a pattern from class +1 (depicted in a square) is situated in the -1 region of the hyperplane

1260 the solution less sensitive to classification errors. The optimization problem is solved with Lagrange multipliers, similar with the procedure outlined for the hard margin linear SVM, when the primal Lagrangian expression is:

$$1265 \quad L_P(\mathbf{w}, b, \mathbf{A}, \mathbf{M}) =$$

$$1266 \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i]$$

$$1267 \quad - \sum_{i=1}^m \mu_i \xi_i \quad (10)$$

1269 where the Lagrange multipliers  $\mathbf{A} = (\alpha_1, \alpha_2, \dots, \alpha_m)$  are assigned to each constraint  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq +1 - \xi_i$ , and the Lagrange multipliers  $\mathbf{M} = (\mu_1, \mu_2, \dots, \mu_m)$  are assigned to each constraint  $\xi_i \geq 0, \forall i = 1, \dots, m$ . The dual optimization problem is:

$$1274 \quad \text{maximize}$$

$$L_D(\mathbf{w}, b, \mathbf{A}, \mathbf{M}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

subject to  $0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$  1275

$$1266 \quad \text{and } \sum_{i=1}^m \alpha_i y_i = 0. \quad (11)$$

The expression for the vector  $\mathbf{w}$  is identical with that obtained for hard margin SVM: 1278

$$1279 \quad \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i. \quad (12)$$

The SVM classifier, which is determined by the pair  $(\mathbf{w}, b)$ , may be used for the classification of new molecules. The class of a molecule  $k$  with the structural descriptors  $\mathbf{x}_k$  is: 1281

$$1282 \quad \text{class}(k) = \text{sign} \left( \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_k + b \right) \quad (13)$$

where the summation goes only over the support vectors. 1283

## 1286 Nonlinear Classification with Kernels

1287 Nonlinear relationships between structural descriptors  
1288 and drug properties may be modeled with SVM equations  
1289 that replace the linear classifier with nonlinear kernel func-  
1290 tions. The SVM extension to nonlinear relationships is  
1291 based on the property of feature functions  $\phi(\mathbf{x})$  that trans-  
1292 form the input spaces into a feature space of higher dimen-  
1293 sionality, in which a linear separation of the classes  
1294 is possible. This transformation is depicted in Fig. 8a. In  
1295 the input space, which is represented by the structural de-  
1296 scriptors, a circle can discriminate the two classes of ob-  
1297 jects. Using a proper feature function  $\phi$  the input space is  
1298 transformed into a feature space in which the two classes  
1299 can be separated with a linear classifier. Figure. 8b illus-  
1300 trates another situation in which another nonlinear func-  
1301 tion can discriminate the classes in input space, which is  
1302 then transformed by the function  $\phi$  into a linear classifier  
1303 in feature space. The challenge is to find that feature func-  
1304 tion  $\phi$  that maps the objects into a feature space in which  
1305 the classes are linearly separable. In real-life applications  
1306 it is not even desirable to achieve a perfect separation, be-  
1307 cause the function may map noise or errors in data (see  
1308 Fig. 8b, where a +1 object, depicted in a square, is situated  
1309 in the -1 region of the hyperplane). Due to the fact that  
1310 in feature space the objects may be separated with a linear  
1311 classifier, it is possible to use the results for linear SVM  
1312 presented in the previous sections, which is a fundamental  
1313 property of SVM.

1314 The mapping from the input space to the feature space  
1315 may be expressed as:

$$1316 \quad \mathbf{x} = (x_1, x_2, \dots, x_n)$$

$$1317 \quad \rightarrow \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_h(\mathbf{x})). \quad (14)$$

1320 A soft margin SVM in the feature space is obtained by sim-  
1321 ply substituting the input vectors  $\mathbf{x}$  with the feature vectors  
1322  $\phi(\mathbf{x})$ . We have to point here that a good linear classifier in  
1323 feature space depends on a proper form of the function  $\phi$ .  
1324 The class of a molecule  $k$  with the structural descriptors  $\mathbf{x}_k$   
1325 is:

$$1326 \quad \text{class}(k) = \text{sign}[\mathbf{w} \cdot \phi(\mathbf{x}_k) + b]$$

$$= \text{sign} \left( \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_k) + b \right) \quad (15)$$

1327 The prediction for a molecule  $k$  requires the computation  
1328 of the dot product  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_k)$  for all support vectors  $\mathbf{x}_i$ .  
1329 This is an important observation showing that we do not  
1330 have to know the actual expression of the feature func-  
1331 tion  $\phi$ . Furthermore, nonlinear SVM use kernels, which

are a special class of functions that compute the dot prod- 1332  
uct  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_k)$  in the input space: 1333

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j). \quad (16) \quad 1334$$

Several popular kernel functions are presented below. 1335  
The inner product of two vectors defines the linear (dot) 1336  
kernel: 1337

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j. \quad (17) \quad 1338$$

The dot kernel is a linear classifier, and should be used 1339  
as a reference to demonstrate an eventual improvement 1340  
of the classification with nonlinear kernels. The polyno- 1341  
mial kernel is useful mainly for lower degrees, because for 1342  
higher degrees it has the tendency to overfit the data: 1343

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i \cdot \mathbf{x}_j)^d. \quad (18) \quad 1344$$

The Gaussian radial basis functions (RBF) kernel is very 1345  
flexible, and depending on the values of the parameter  $\sigma$  it 1346  
can display a wide range of nonlinearity: 1347

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\sigma \|\mathbf{x}_i - \mathbf{x}_j\|^2). \quad (19) \quad 1348$$

The hyperbolic tangent (tanh) function has a sigmoid 1349  
shape and it is the most used transfer function for arti- 1350  
ficial neural networks. The corresponding kernel has the 1351  
formula: 1352

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(a\mathbf{x}_i \cdot \mathbf{x}_j + b). \quad (20) \quad 1353$$

The nonlinearity of the anova kernel is controlled by the 1354  
parameters  $\gamma$  and  $d$ : 1355

$$K(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_i \exp(-\gamma(\mathbf{x}_i - \mathbf{x}_j)) \right)^d. \quad (21) \quad 1356$$

## 1357 Hard Margin Nonlinear SVM Classification

The formulation of the hard margin nonlinear SVM clas- 1358  
sification is obtained from the linear SVM by substituting 1359  
input vectors  $\mathbf{x}$  with feature functions  $\phi(\mathbf{x})$ . The dual prob- 1360  
lem is: 1361

maximize

$$L_D(\mathbf{w}, b, \mathbf{A}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

subject to  $\alpha_i \geq 0, \quad i = 1, \dots, m$

$$\text{and } \sum_{i=1}^m \alpha_i y_i = 0. \quad (22) \quad 1362$$

1363 The optimum separation hyperplane is determined by the  
1364 vector  $\mathbf{w}$ :

$$1365 \quad \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i). \quad (23)$$

1366 In the final expression of the classifier, the dot product for  
1367 two feature functions  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_k)$  is replaced with a kernel  
1368 function  $K(\mathbf{x}_i, \mathbf{x}_k)$ :

$$1369 \quad \text{class}(k) = \text{sign} \left( \sum_{i=1}^m \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_k) + b \right). \quad (24)$$

1370 The structure of a nonlinear SVM may be represented in  
1371 a network format (Fig. 9). The network input is repre-  
1372 sented by a set of support vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and a predic-  
1373 tion object  $\mathbf{x}_p$ . The input space is transformed with a fea-  
1374 ture function  $\phi$  into the feature space, and the dot product  
1375 of feature functions is computed and multiplied with the  
1376 Lagrangian multipliers  $\alpha$ . The output is Eq. (24) in which  
1377 dot product of feature functions is substituted with a ker-  
1378 nel function  $K$ .

### 1379 Soft Margin Nonlinear SVM Classification

1380 The expression of the soft margin nonlinear SVM classifi-  
1381 cation is obtained from the corresponding linear SVM for-  
1382 mula by substituting input vectors  $\mathbf{x}$  with feature functions  
1383  $\phi(\mathbf{x})$ . The dual problem is:

maximize

$$L_D(\mathbf{w}, b, \mathbf{A}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

subject to  $0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$

$$1384 \quad \text{and} \quad \sum_{i=1}^m \alpha_i y_i = 0 \quad (25)$$

1385 with a solution identical with that for hard margin nonlin-  
1386 ear SVM classification from Eq. (24).

1387 In Fig. 10 we demonstrate the nonlinear SVM classifi-  
1388 cation for the Gaussian RBF kernel. The best separation  
1389 of the two classes is obtained for a low  $\sigma$  value (Fig. 10a,  
1390  $\sigma = 0.1$ ). As the value of  $\sigma$  increases the kernel nonlin-  
1391 earity increases and the separation surface becomes more  
1392 complex. It is clear from this example that the kernel non-  
1393 linearity must match the structure of the data. Also, highly  
1394 nonlinear functions do not translate in better SVM mod-  
1395 els. In practical applications the parameters that determine  
1396 the kernel shape must be optimized to provide the best  
1397 predictions in a cross-validation test.

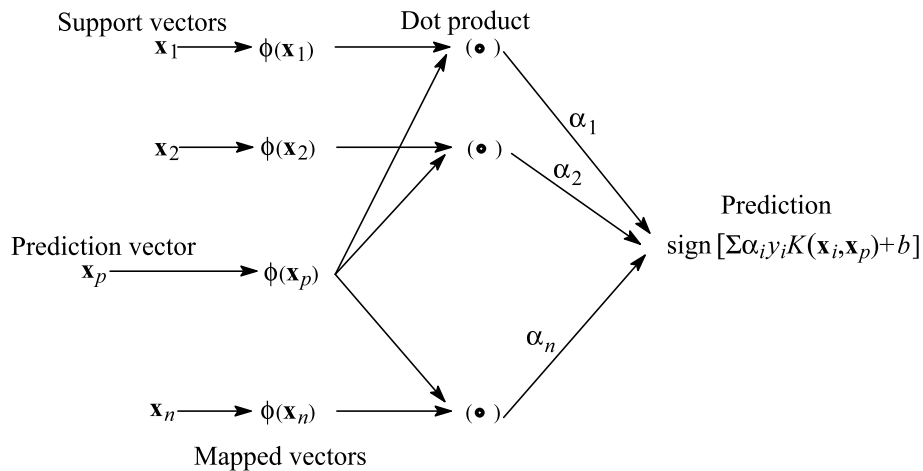
### SVM Regression

Vapnik extended the SVM algorithm for regression [139] 1399  
by using an  $\varepsilon$ -insensitive loss function (Fig. 11). The SVM 1400  
regression (SVMR) determines a function  $f(\mathbf{x})$  with the 1401  
property that for all learning objects  $\mathbf{x}$  it has a maxi- 1402  
mum deviation  $\varepsilon$  from the target (experimental) values  $y$ , 1403  
and it has a maximum margin. Starting from the training 1404  
objects, SVMR computes a model representing a tube 1405  
with radius  $\varepsilon$  fitted to the data. All object situated in- 1406  
side the regression tube have an error equal to zero. For 1407  
the hard margin SVMR no object is allowed outside the 1408  
tube, whereas for the soft margin SVMR uses positive 1409  
slack variables to penalize the objected situated outside the 1410  
tube [44,45,93,94,122]. If Fig. 11 the objects situated inside 1411  
the regression tube are shown as white circles, whereas the 1412  
objects outside the regression tube are depicted as black 1413  
circles. The slack variable  $\xi$  introduces a penalty propor- 1414  
tional with the distance between the object and the margin 1415  
of the regression tube. Support vectors are selected only 1416  
from objects situated outside the tube. 1417

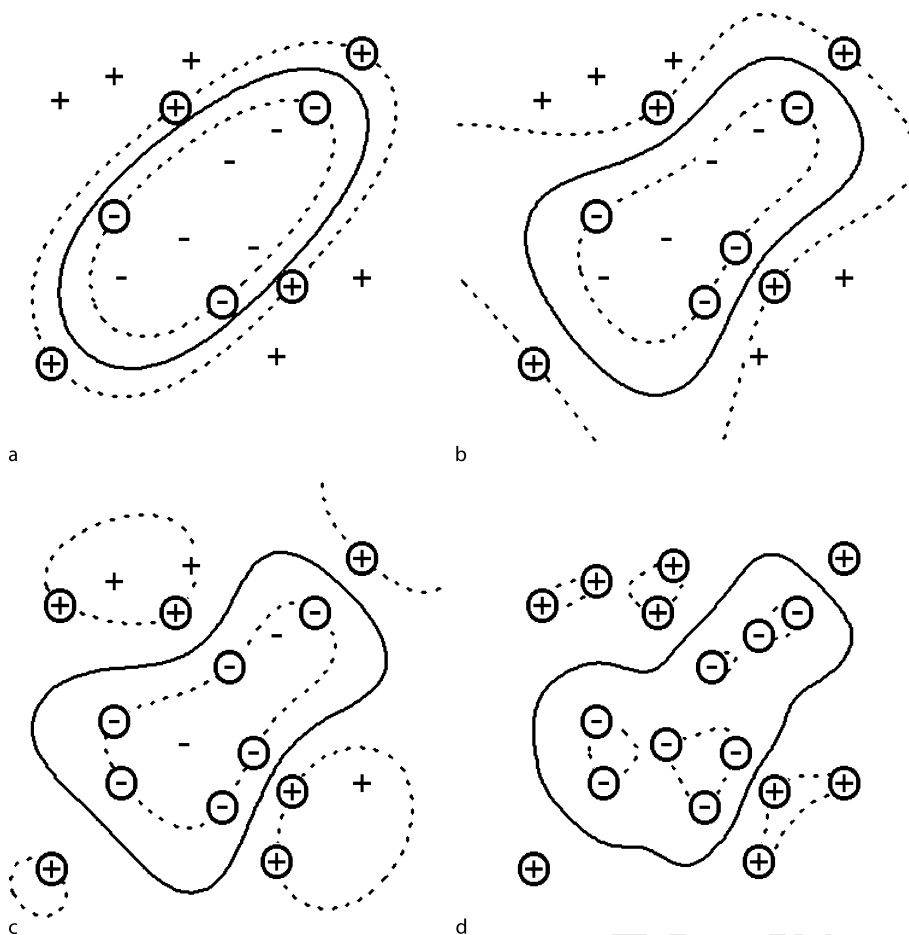
To demonstrate the kernel influence on the SVM re- 1418  
gression we model the inhibition constant  $K_i$  for bovine 1419  
milk xanthine oxidase [51]. The structural descriptor used 1420  
is ClogP, the computed hydrophobicity. To model the 1421  
nonlinear dependence between ClogP and  $K_i$  we apply 1422  
a degree 2 polynomial kernel (Fig. 12) with different values 1423  
for the  $\varepsilon$ -insensitive parameter. The molecules situated in- 1424  
side the regression tube are indicated with black triangles, 1425  
whereas those located outside the regression tube are rep- 1426  
resented as black dots (support vectors are shown as black 1427  
dots inside a circle). A low value for  $\varepsilon$  produces a SVMR 1428  
model with a larger number of support vectors, but an 1429  
overall good shape for the relationship between ClogP and 1430  
 $K_i$  (Fig. 12a,  $\varepsilon = 0.1$ ). As  $\varepsilon$  increases the diameter of the 1431  
regression tube includes almost all molecules, and thus the 1432  
model depends on a small number of support vectors be- 1433  
cause objects situated inside the tube do not influence the 1434  
regression model. 1435

We explore further two higher order polynomials, 1436  
namely a degree 5 polynomial kernel (Fig. 13a) and a de- 1437  
gree 8 polynomial kernel (Fig. 13b). Both kernels overfit 1438  
the data, and have shapes that are too complex. The Gaus- 1439  
sian RBF kernel with  $\sigma = 0.2$  is a good fit for the  $K_i$  data 1440  
(Fig. 13c), but an increase in  $\sigma$  results in a regression that 1441  
start to fit the errors in the experimental data (Fig. 13d). 1442

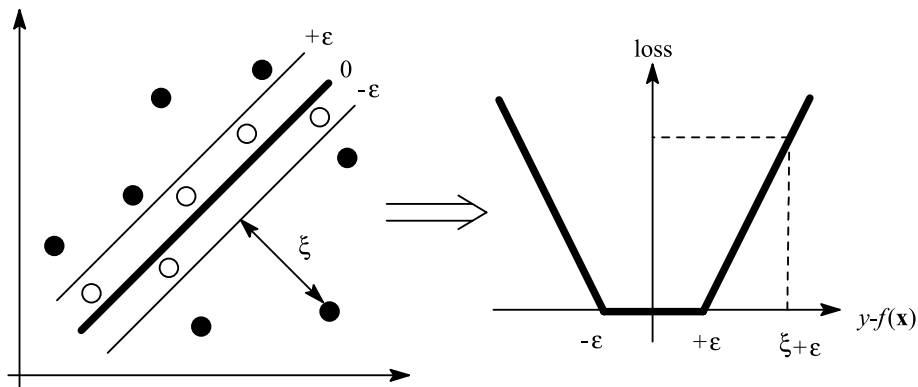
The next set of SVM regression QSAR evaluates the 1443  
tanh kernel (Fig. 14). The first combination of paramet- 1444  
ers is insensitive to the relationship between ClogP and 1445  
 $K_i$  (Fig. 14a). Although the  $K_i$  data are best represented 1446  
by a downward parabola, the second tanh QSAR has the 1447



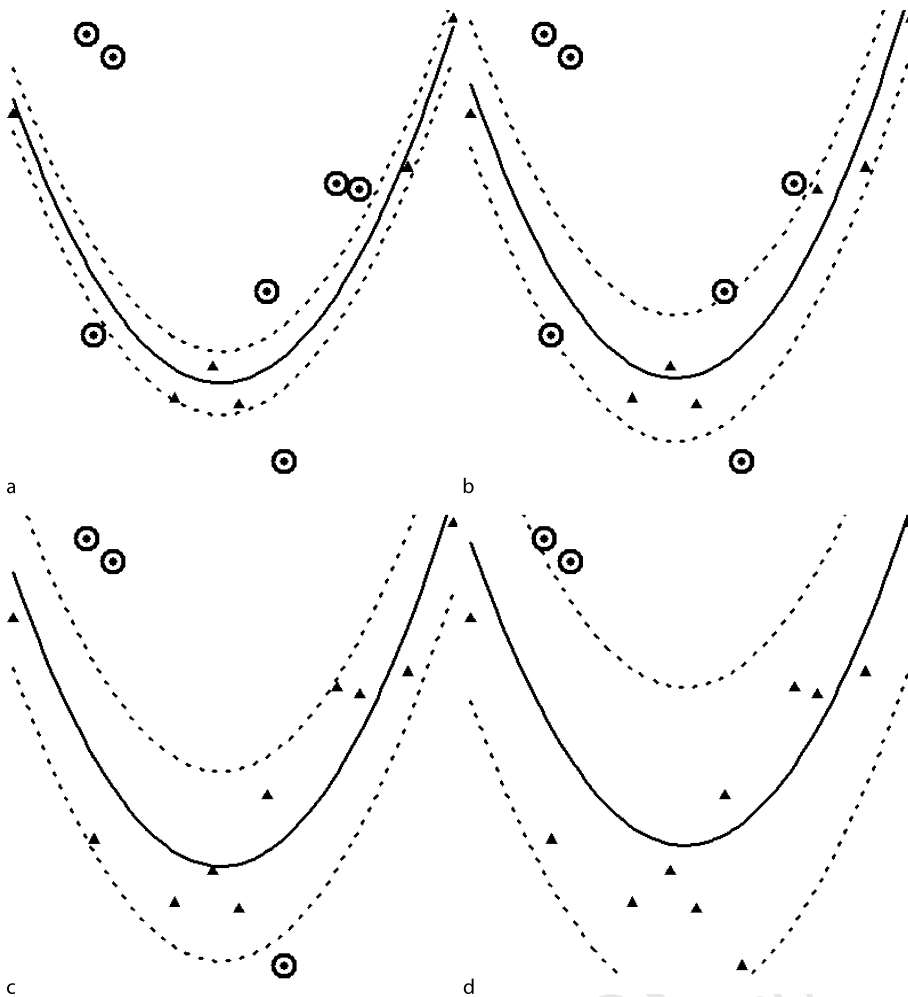
Drug Design with Machine Learning, Figure 9  
Network representation of support vector machines



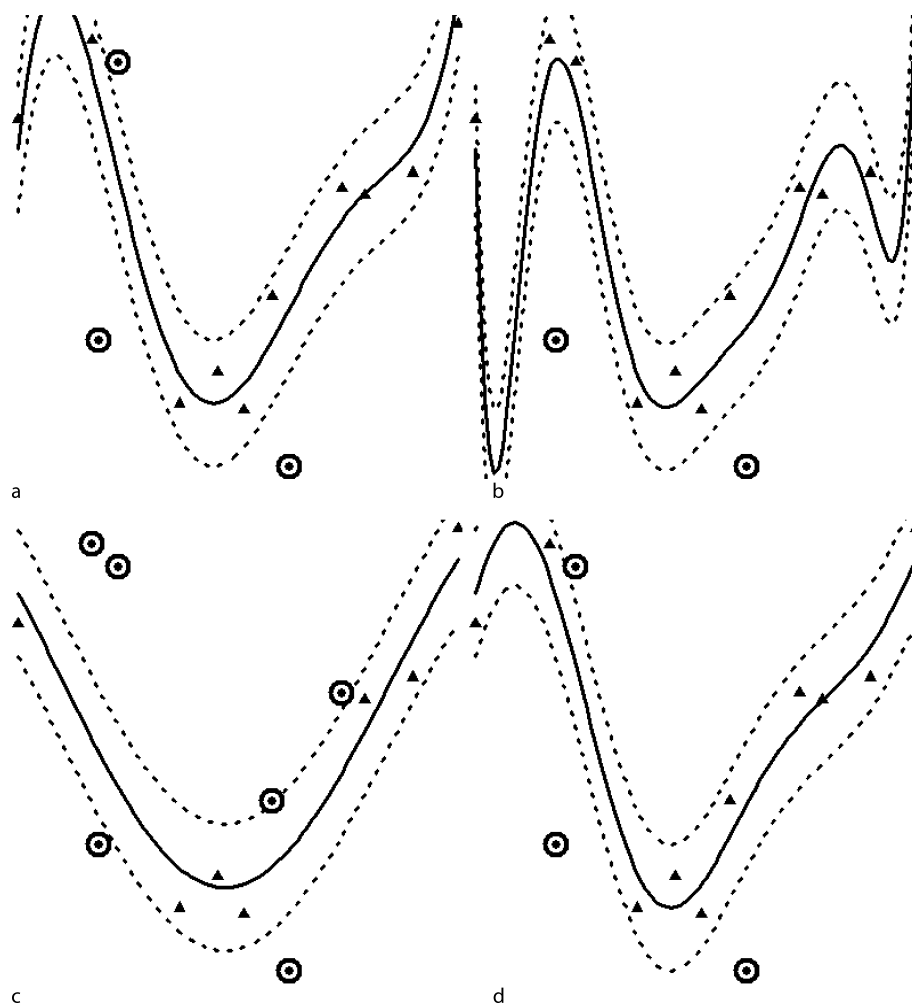
Drug Design with Machine Learning, Figure 10  
SVM classification models for linearly non-separable data with the Gaussian RBF kernel: a  $\sigma = 0.1$ ; b  $\sigma = 0.5$ ; c  $\sigma = 1$ ; d  $\sigma = 5$



Drug Design with Machine Learning, Figure 11  
SVM regression with  $\epsilon$ -insensitive loss function



Drug Design with Machine Learning, Figure 12  
SVM regression models with a degree 2 polynomial kernel ( $C = 1000$ ): a  $\epsilon = 0.1$ ; b  $\epsilon = 0.2$ ; c  $\epsilon = 0.3$ ; d  $\epsilon = 0.5$



Drug Design with Machine Learning, Figure 13

SVM regression models ( $C = 1000$ ,  $\epsilon = 0.2$ ): a degree 5 polynomial kernel; b degree 8 polynomial kernel; c Gaussian RBF kernel,  $\sigma = 0.1$ ; d Gaussian RBF kernel,  $\sigma = 0.5$

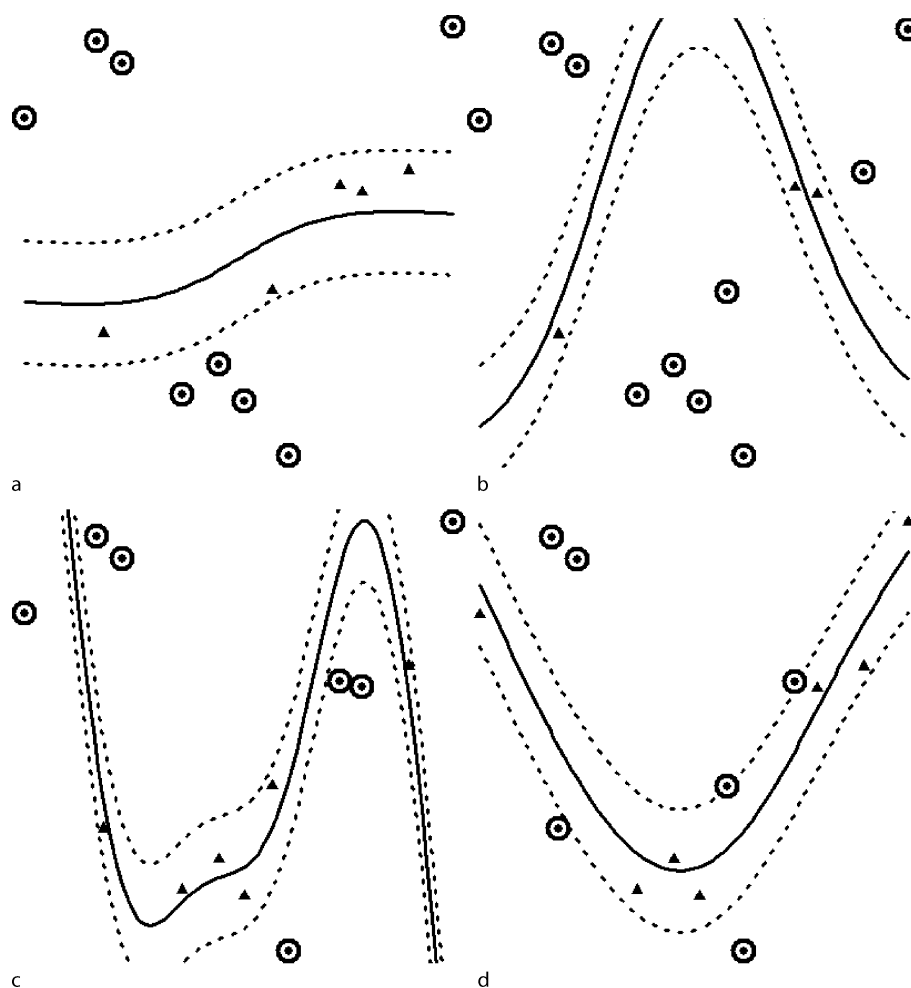
1448 shape of an upward parabola, and misses the most appar- 1463  
 1449 ent correlation between ClogP and  $K_i$  (Fig. 14b). The third 1464  
 1450 tanh QSAR is overfitted, mainly in its right side (Fig. 14c). 1465  
 1451 The last combination of parameters for the tanh QSAR 1466  
 1452 shows a very good fit of the data and obviously it is a good 1467  
 1453 candidate for the prediction of new molecules (Fig. 14d). 1468  
 1454 Considering the large variation in shape of the tanh QSAR 1469  
 1455 it is apparent that its parameters should be varied over 1470  
 1456 a large range to find the best regression model. 1471

1457 Since the dataset considered is a real QSAR problem, 1472  
 1458 it is interesting to examine the behavior of other, less 1473  
 1459 common, kernels. The anova RBF kernel shows an over- 1474  
 1460 fit of the data, which is more prominent in its left side 1475  
 1461 (Fig. 15a). The spline kernel has a too low nonlinearity, 1476  
 1462 and the model has no value in predicting new chemicals 1477

(Fig. 15b). The degree 1 Bessel kernel is a very good fit for 1463  
 the correlation between ClogP and  $K_i$  (Fig. 15c), whereas 1464  
 a degree 5 Bessel kernel overfits the data in the left side 1465  
 (Fig. 15d). The results presented here are only to illustrate 1466  
 the shape of various kernels, and to demonstrate their ability 1467  
 to model nonlinear relationships. Other combinations 1468  
 of parameters and capacity  $C$  may give better or worse cor- 1469  
 relations, and the most predictive QSAR may be identified 1470  
 only by optimizing the kernel parameters. 1471

#### Kernel Comparison in QSAR 1472

Several kernels have the property to give SVM models that 1473  
 are universal approximators to arbitrary functions, namely 1474  
 the SVM model can approximate any function to any level 1475



Drug Design with Machine Learning, Figure 14

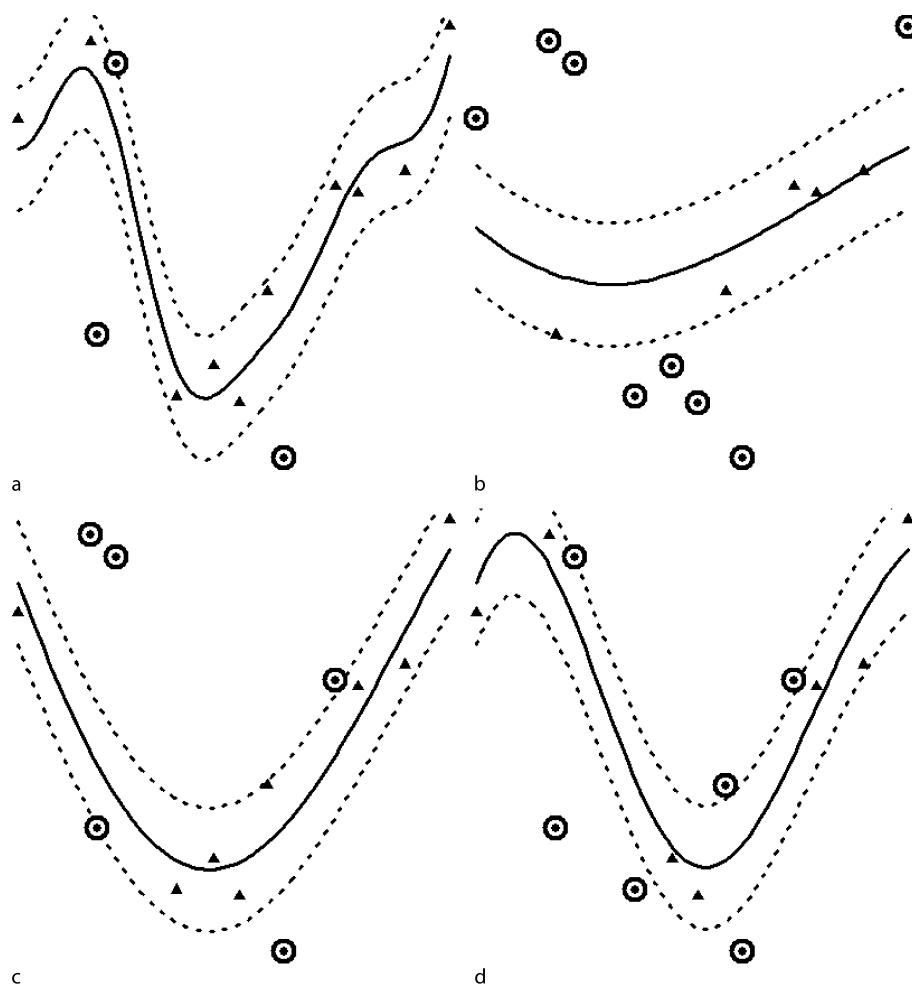
SVM regression models with a hyperbolic tangent (tanh) kernel ( $\epsilon = 0.2$ ): a  $a = 1, b = 0, C = 1$ ; b  $a = 1, b = 1, C = 1$ ; c  $a = 0.9, b = -1, C = 10$ ; d  $a = 0.5, b = -0.7, C = 10$

1476 of precision. However, the theory cannot guarantee that  
 1477 such a model will be also optimally predictive. The essen-  
 1478 tial requirement for a SAR or QSAR model is its ability to  
 1479 provide reliable predictions for new molecules. The sim-  
 1480 ple data fitting ability with a low error is not enough for  
 1481 an SVM model, because such models are usually overfit-  
 1482 ted and give bad predictions. The experimental data used  
 1483 in SAR and QSAR are affected by experimental errors  
 1484 and sometimes by the fact that good structural descrip-  
 1485 tors are missing from the model. Therefore it is impor-  
 1486 tant to evaluate a number of kernels in order to find the  
 1487 most predictive SVM model and to avoid fitting the noise  
 1488 or the errors from the data. To offer a comparative eval-  
 1489 uation of several kernels we review here a number of SVM  
 1490 classification and regression models in which five kernels  
 1491 are compared, namely linear, polynomial, Gaussian RBF,

1492 tanh, and anova. Each kernel was evaluated for a large  
 1493 range of parameters, and all results are compared for  
 1494 cross-validation. The SVM models were computed with  
 1495 mySVM, by Rüping ([http://www-ai.cs.uni-dortmund.de/](http://www-ai.cs.uni-dortmund.de/SOFTWARE/MYSVM/index.html)  
 1496 SOFTWARE/MYSVM/index.html).

1497 The carcinogenic activity of a group of 46 methylated  
 1498 and 32 non-methylated polycyclic aromatic hydrocarbons  
 1499 was modeled with semiempirical quantum indices [61].  
 1500 Form the set of 78 chemicals, 34 are carcinogenic and  
 1501 44 are non-carcinogenic. The accuracy for leave-10%-out  
 1502 (L10%O) cross-validation tests shows that the best pre-  
 1503 dictions are obtained with the Gaussian RBF kernel: RBF,  
 1504  $\sigma = 0.5, AC = 0.86$ ; anova,  $\gamma = 0.5, d = 1, AC = 0.84$ ;  
 1505 degree 2 polynomial,  $AC = 0.82$ ; linear,  $AC = 0.76$ ; tanh,  
 1506  $a = 2, b = 0, AC = 0.66$ . All SVM classification models  
 1507 were obtained with  $C = 10$ . The increase in prediction ac-





Drug Design with Machine Learning, Figure 15

SVM regression models ( $\epsilon = 0.2$ ): a anova RBF kernel,  $\sigma = 0.2$ , degree = 5,  $C = 1000$ ; b spline kernel,  $C = 0.05$ ; c Bessel kernel,  $\sigma = 1$ , order = 1, degree = 1,  $C = 1000$ ; d Bessel kernel,  $\sigma = 0.6$ , order = 1, degree = 5,  $C = 1000$

1508 curacy for RBF and anova kernel, compared with the linear SVM, indicates that there is a nonlinear relationship  
 1509 between the quantum indices and hydrocarbon carcinogenicity.  
 1510  
 1511

1512 A structure-odor classification was performed for  
 1513 98 tetra-substituted pyrazines representing three odor  
 1514 classes, namely 32 green, 23 nutty, and 43 bell-pepper [62].  
 1515 L10%O predictions were obtained for three classification  
 1516 tests, in which one odor type is selected as class +1 and  
 1517 the remaining two odor types form class -1. For green  
 1518 aroma compounds the best predictions are obtained with  
 1519 a polynomial kernel: degree 2 polynomial,  $C = 1000$ ,  
 1520  $AC = 0.86$ ; anova,  $C = 100$ ,  $\gamma = 0.5$ ,  $d = 1$ ,  $AC = 0.84$ ;  
 1521 linear,  $C = 10$ ,  $AC = 0.80$ ; RBF,  $C = 10$ ,  $\sigma = 0.5$ ,  
 1522  $AC = 0.79$ ; tanh,  $C = 10$ ,  $a = 0.5$ ,  $b = 0$ ,  $AC = 0.73$ .

1523 For nutty aroma compounds the anova kernel has slightly  
 1524 better predictions: anova,  $C = 100$ ,  $\gamma = 0.5$ ,  $d = 1$ ,  
 1525  $AC = 0.92$ ; linear,  $C = 10$ ,  $AC = 0.89$ ; degree 2 poly-  
 1526 nomial,  $C = 10$ ,  $AC = 0.89$ ; RBF,  $C = 10$ ,  $\sigma = 0.5$ ,  
 1527  $AC = 0.89$ ; tanh,  $C = 100$ ,  $a = 0.5$ ,  $b = 0$ ,  $AC = 0.79$ .  
 1528 The SVM models for bell-pepper aroma compounds show  
 1529 that three kernels (RBF, polynomial, and anova) give good  
 1530 predictions: RBF,  $C = 10$ ,  $\sigma = 0.5$ ,  $AC = 0.89$ ; degree 2  
 1531 polynomial,  $C = 10$ ,  $AC = 0.88$ ; anova,  $C = 10$ ,  $\gamma = 0.5$ ,  
 1532  $d = 1$ ,  $AC = 0.87$ . linear,  $C = 10$ ,  $AC = 0.74$ ; tanh,  
 1533  $C = 100$ ,  $a = 2$ ,  $b = 1$ ,  $AC = 0.68$ . In all three SVM ex-  
 1534 periments the best model outperformed the linear kernel,  
 1535 showing that nonlinear relationships are useful in struc-  
 1536 ture-aroma studies.

1537 The relationships between the structure of a chemical

1538 compounds and their aqueous toxicity is highly relevant  
1539 because many chemicals are environmental pollutants.  
1540 The mechanism of toxic action (MOA) of organic chemi-  
1541 cals may be predicted from hydrophobicity and experi-  
1542 mental toxicity against *Pimephales promelas* and *Tetrahy-*  
1543 *mena pyriformis* [67]. SVM classification was used to dis-  
1544 criminate between 126 nonpolar narcotics and 211 chemi-  
1545 cals that have other MOA. The RBF, anova, linear and  
1546 polynomial kernels have similar statistics: RBF,  $\sigma = 1$ ,  
1547  $AC = 0.80$ ; anova,  $\gamma = 0.5$ ,  $d = 2$ ,  $AC = 0.80$ ;  
1548 linear,  $AC = 0.79$ ; degree 2 polynomial,  $AC = 0.79$ ;  
1549 tanh,  $a = 0.5$ ,  $b = 0$ ,  $AC = 0.53$ . These results indi-  
1550 cated that there is a linear relationship between descrip-  
1551 tors and MOA. A similar conclusion was obtained for the  
1552 MOA classification of polar and nonpolar narcotic com-  
1553 pounds [66]. In a related study for the classification of  
1554 MOA for narcotic and reactive compounds [63], the de-  
1555 gree 2 polynomial offered the best predictions: degree 2  
1556 polynomial,  $C = 10$ ,  $AC = 0.92$ ; anova,  $C = 10$ ,  $\gamma = 0.5$ ,  
1557  $d = 1$ ,  $AC = 0.87$ ; linear,  $C = 1000$ ,  $AC = 0.86$ ; RBF,  
1558  $C = 100$ ,  $\gamma = 0.5$ ,  $AC = 0.83$ ; tanh,  $C = 10$ ,  $a = 0.5$ ,  
1559  $b = 0$ ,  $AC = 0.78$ .

1560 SVM regression is demonstrated for QSAR models ob-  
1561 tained for 52 benzodiazepine receptor ligands [68]. Five  
1562 structural descriptors were used as SVM input, and the  
1563 prediction was evaluated with leave-5%-out (L5%O) and  
1564 L10%O cross-validation. A number of 34 experiments  
1565 were conducted with five kernels (Tab. 1). The predic-  
1566 tions obtained with the linear kernel L are used as base-  
1567 line to compare the results provided by the other kernels.  
1568 All four polynomial kernels P give very bad predictions,  
1569 as indicated by the statistical indices. The RBF kernel R  
1570 has better predictions than L only for  $\sigma = 0.25$  (experi-  
1571 ment 6), whereas all the other RBF models are worse than  
1572 L. Overall, the best results are obtained with the tanh ker-  
1573 nel T for  $a = 0.5$  and  $b = 0$  (experiment 11) which is  
1574 a clear improvement compared to the linear kernel. Fi-  
1575 nally, the anova kernel A has uneven predictions. Some are  
1576 good (experiments 22 and 20) whereas experiments 25–  
1577 34 have very bad predictions. The SVM regression experi-  
1578 ments presented in Tab. 1 show that kernel parameters are  
1579 of paramount importance in obtaining a predictive QSAR.

1580 In this section we presented comparative SAR and  
1581 QSAR models obtained with five kernels. In developing  
1582 SVM models it is important to include the linear kernel  
1583 as a baseline comparison. If nonlinear kernels do not im-  
1584 prove the predictions of the linear kernel, then the struc-  
1585 ture-activity relationship is linear and a linear SVM should  
1586 be used. The parameters of nonlinear kernels should be  
1587 optimized in order to obtain the best prediction. The most  
1588 predictive SAR and QSAR model may be obtained with

1589 anyone of the five kernels investigated, and there is no ap-  
1590 parent rule that can predict which kernel is the best. There-  
1591 fore, it is important to experiment with a large diversity of  
1592 kernels in order to find a predictive SVM model.

### 1593 Applications in Drug Design

1594 Torsade de pointes (TdP) is an important adverse drug re-  
1595 action that is responsible for almost one-third of all drug  
1596 failures during drug development, and resulted in several  
1597 drugs being withdrawn from the market. Yap et al com-  
1598 pared several ML for their ability to separate TdP com-  
1599 pounds from those that do not induce TdP [156]. The  
1600 prediction accuracy shows that SVM is the most reliable  
1601 ML: SVM with Gaussian RBF kernel 91.0%,  $k$ -NN 88.5%,  
1602 probabilistic neural network 78.2%, and C4.5 decision tree  
1603 65.4%. Inhibitors of the hERG (human ether-a-go-go-  
1604 related gene) potassium channel can lead to a prolonga-  
1605 tion of the QT interval that can trigger TdP, an atypical  
1606 ventricular tachycardia. A dataset of 39 hERG inhibitors  
1607 (26 for learning and 13 for test) was evaluated by a phar-  
1608 macophore ensemble coupled with SVM regression [87].  
1609 The pharmacophore ensemble encodes the protein con-  
1610 formational flexibility, and the SVM brings the capability  
1611 of modeling nonlinear relationships. The QSAR model has  
1612 very good statistics, with  $q^2 = 0.89$  and  $r^2 = 0.94$  for  
1613 the test set. A structure-activity SVM study for hERG in-  
1614 hibitors was performed by Tobita et al for 73 drugs [132].  
1615 Using a combination of structural descriptors and molec-  
1616 ular fragments, the SVM model had a prediction accuracy  
1617 higher than 90%. Selecting an optimum group of descrip-  
1618 tors is an important phase in computing a predictive struc-  
1619 ture-activity model. The SVM recursive feature elimina-  
1620 tion algorithm was evaluated by Xue et al for three SAR  
1621 datasets, namely  $P$ -glycoprotein substrates, human intes-  
1622 tinal absorption, and compounds that cause TdP [153]. The  
1623 prediction statistics show that a proper descriptor selec-  
1624 tion may significantly increase the SVM performance.

1625 The drug metabolism by human cytochrome P450 iso-  
1626 forms 3A4, 2D6, and 2C9 was studied by Terfloeth et al for  
1627 379 drugs and drug analogs [131]. The best SVM model  
1628 has a cross-validation accuracy of 89% and 83% accu-  
1629 racy for 233 validation compounds. Unbalanced data sets,  
1630 with a majority of inactive compounds and only a few  
1631 active chemicals, represent the usual situation in library  
1632 screening. To overcome the ML bias toward the larger  
1633 class, Eitrich et al used SVM with oversampling to iden-  
1634 tify cytochrome P450 2D6 inhibitors [37]. An ensemble of  
1635 SVM models was presented by Yap and Chen as suitable  
1636 classifiers for inhibitors and substrates of cytochromes  
1637 P450 3A4, 2D6, and 2C9 [155]. Descriptors computed

Drug Design with Machine Learning, Table 1  
Support vector regression statistics for benzodiazepine receptor ligands

Exp	Kernel	$\rho_1$	$\rho_2$	$q_{L5\%O}^2$	RMSE <sub>L5%O</sub>	$q_{L10\%O}^2$	RMSE <sub>L10%O</sub>
1	L			0.261	0.98	0.273	0.97
2	P	2		< -100	> 10	< -100	> 10
3	P	3		< -100	> 10	< -100	> 10
4	P	4		< -100	> 10	< -100	> 10
5	P	5		< -100	> 10	< -100	> 10
6	R	0.25		0.368	0.91	0.370	0.91
7	R	0.5		0.226	1.00	0.324	0.94
8	R	1.0		0.221	1.01	0.310	0.95
9	R	1.5		0.236	1.00	0.258	0.98
10	R	2.0		0.208	1.01	0.205	1.02
11	T	0.5	0.0	0.453	0.84	0.498	0.81
12	T	1.0	0.0	0.416	0.87	0.411	0.87
13	T	2.0	0.0	0.396	0.89	0.394	0.89
14	T	0.5	1.0	0.070	1.10	0.120	1.07
15	T	1.0	1.0	0.389	0.89	< -10	5.36
16	T	2.0	1.0	0.297	0.96	0.348	0.92
17	T	0.5	2.0	-0.365	1.33	-0.405	1.35
18	T	1.0	2.0	0.106	1.08	0.243	0.99
19	T	2.0	2.0	0.345	0.92	0.376	0.90
20	A	0.25	1	0.397	0.89	0.377	0.90
21	A	0.5	1	0.324	0.94	0.331	0.93
22	A	1.0	1	0.398	0.88	0.412	0.87
23	A	1.5	1	0.299	0.95	0.374	0.90
24	A	2.0	1	0.293	0.96	0.339	0.93
25	A	0.25	2	-1.289	1.73	-0.921	1.58
26	A	0.5	2	-1.375	1.76	-0.873	1.56
27	A	1.0	2	-0.606	1.44	-0.465	1.38
28	A	1.5	2	-0.241	1.27	-0.195	1.25
29	A	2.0	2	-0.071	1.18	-0.060	1.17
30	A	0.25	3	-2.998	2.28	-1.934	1.95
31	A	0.5	3	-1.282	1.72	-0.983	1.61
32	A	1.0	3	-0.060	1.17	-0.094	1.19
33	A	1.5	3	0.060	1.11	-0.097	1.19
34	A	2.0	3	-0.062	1.18	-0.073	1.18

1638 with Dragon were used to develop predictive models,  
 1639 with MCC (Matthews correlation coefficient) statistics,  
 1640 namely 899 for 3A4, 0.884 for 2D6, and 0.872 for 2C9.  
 1641 The classification of of cytochrome P450 3A4 inhibitors  
 1642 and non-inhibitors was investigated by Arimoto et al with  
 1643 SVM,  $k$ -NN, recursive partitioning, logistic regression,  
 1644 and Bayesian classifier [7]. The chemical structure of the  
 1645 4470 compounds was encoded with fingerprints and Mol-  
 1646 connZ topological indices, and the best predictions were  
 1647 obtained with SVM trained with fingerprints. A compar-  
 1648 ison of single and ensemble classifiers was made by Merk-  
 1649 wirth et al for SVM,  $k$ -NN, and ridge regression [95]. The

1650 tests conducted for two libraries of chemicals, one for un-  
 1651 specific protein inhibition and another one for inhibitors  
 1652 of cytochrome P450 3A4, showed that single and ensemble  
 1653 SVM models give similar results, with better predictions  
 1654 than  $k$ -NN and ridge regression.

1655 Du et al predicted the genotoxicity of 140 thiophene  
 1656 derivatives with SVM and linear discriminant analysis  
 1657 (LDA) [34]. Seven structural descriptors were used as in-  
 1658 put data for the classification models, and the SVM param-  
 1659 eters were identified by grid search. The accuracy of the  
 1660 SVM model (92.9% in training and 92.6% in prediction)  
 1661 is significantly higher than the LDA accuracy (81.4% in

1662 training and 85.2% in prediction). Recursive partitioning  
1663 and SVM were compared in a test for mutagenicity predic-  
1664 tion from substructure descriptors [92]. The prediction ac-  
1665 curacy for 2199 mutagens is similar for SVM (81.4%) and  
1666 recursive partitioning (80.2%).

1667 SVM classification was compared with linear discrim-  
1668 inant analysis in predicting drug bioavailability [149]. The  
1669 learning set comprises 167 compounds characterized by  
1670 five structural descriptors. The SVM parameters were op-  
1671 timized with grid search. The classification results indi-  
1672 cate that SVM (accuracy 85.6%) provides better models  
1673 for bioavailability compared to LDA (accuracy 72.4%).  
1674 Sakiyama et al studied SAR models for the metabolic  
1675 stability of drug candidates [113]. The structure-activity  
1676 models were obtained with random forest, support vec-  
1677 tor machine, logistic regression, and recursive partition-  
1678 ing. Classification models were obtained for a collection of  
1679 1952 compounds characterized by 193 descriptors. All ML  
1680 models had accuracy > 0.8, with slightly higher results for  
1681 random forest and SVM.

1682 The identification of drug-like compounds based on  
1683 atom type counts was investigated with SVM, linear pro-  
1684 gramming machines, linear discriminant analysis, bagged  
1685  $k$ -nearest neighbors, and bagged decision trees C4.5 [97].  
1686 The classification dataset consisted of drug-like com-  
1687 pounds from World Drug Index and non-drug com-  
1688 pounds from Available Chemicals Directory. The small-  
1689 est errors were obtained with SVM trained with poly-  
1690 nomial and Gaussian RBF kernels. A computational filter for  
1691 chemical compounds with antidepressant activity was de-  
1692 veloped by Lepp et al based on 21 biological targets related  
1693 to depression [88]. An SVM model obtained with atom  
1694 type descriptors gave high enrichment levels, and showed  
1695 that a compound selected as active interacts on average  
1696 with 2.3 targets. SVM, artificial neural networks, and mul-  
1697 tiple linear regression were compared in a model for apop-  
1698 tosis induction by 4-aryl-4-H-chromenes [40]. A leave-  
1699 one-out cross-validation test showed that SVM gives the  
1700 best predictions. Briem and Günther used a dataset of 565  
1701 kinase inhibitors and 7194 inactive compounds to com-  
1702 pare four machine learning algorithms, namely SVM with  
1703 a Gaussian RBF kernel, artificial neural networks,  $k$ -near-  
1704 est neighbors, and recursive partitioning [23]. The best  
1705 predictions were obtained with SVM, followed closely by  
1706  $k$ -NN. The SVM ability to identify active compounds from  
1707 a chemical library was investigated by Jorissen and Gilson  
1708 for sets of molecules that interact with the  $\alpha_{1A}$  adreno-  
1709 ceptor, and cyclin-dependent kinase, cyclooxygenase-2,  
1710 factor Xa, and phosphodiesterase-5 [72]. A comparison  
1711 over ten QSAR datasets was performed to compare the  
1712 stochastic gradient boosting method (SGB, an ensem-

1713 ble of classification and regression trees) with decision  
1714 tree, random forest, partial least squares,  $k$ -NN, Bayes  
1715 classifier, and SVM [129]. The best results predictions  
1716 were obtained with random forest, followed by SVM and  
1717 SGB.

1718 The enzyme-substrate interaction is very specific,  
1719 which points to the existence of recognition elements lo-  
1720 cated on the enzyme interface. A surface patch rank-  
1721 ing (SPR) method was proposed for the identification of  
1722 those protein residues that determine the substrate speci-  
1723 ficity [157]. SPR incorporates an SVM module that can  
1724 highlight the residues important in ligand binding. The  
1725 method was applied with good results for several homol-  
1726 ogous enzymes, namely guanylyl/adenylyl cyclases, lac-  
1727 tate/malate dehydrogenases, and trypsin/chymotrypsin.  
1728 Potential applications are residue selection for mutagenesis  
1729 experiments and functional annotation of proteins. Sing-  
1730 le base DNA mutations that result in an amino acid sub-  
1731 stitution are a common cause of monogenic disease. These  
1732 substitutions reduce the protein stability, as measured by  
1733 reduction in hydrophobic area, overpacking, backbone  
1734 strain, and loss of electrostatic interactions [158]. An SVM  
1735 classifier was trained with a set of mutations causative of  
1736 disease and another set of non-disease causing mutations.  
1737 The cross-validation tests show that the SVM identifies  
1738 74% of disease mutations, and confirms that loss of pro-  
1739 tein stability is a determinant factor in monogenic disease.

### 1740 Comparative Studies

1741 SAR and QSAR models have two major components,  
1742 namely a set of structural descriptors and a machine learn-  
1743 ing or ensemble of ML algorithms. The studies reviewed  
1744 in previous sections show that there is no universal fam-  
1745 ily of descriptors that gives the best predictions for any  
1746 molecular property, which justifies the continuous devel-  
1747 opment of novel structural descriptors. Similarly, no ML  
1748 algorithm performs optimally for any property and any set  
1749 of descriptors. Extensive computational experiments show  
1750 that for a given property, the best model can be identi-  
1751 fied only by an empirical comparison of a large number of  
1752 ML methods. In this section we review several studies that  
1753 sought to compare the predictive abilities of ML models.

1754 The aqueous solubility of drugs and drug-related com-  
1755 pounds was modeled by Schroeter et al with Gaussian pro-  
1756 cess, random forest, SVM, and ridge regression [118]. The  
1757 domain of applicability for each model was estimate by  
1758 computing error levels. The aqueous solubility of 988 or-  
1759 ganic chemicals was modeled with four regression algo-  
1760 rithms, namely partial least squares, random forest, SVM,  
1761 and MLF artificial neural networks [104]. The best predic-

tions were obtained with random forest, which also gave good results for an external validation set of 330 chemicals. The melting temperatures of ionic liquids were modeled with SVM, associative neural networks,  $k$ -NN, partial least squares, multilayer feedforward (MLF) artificial neural networks, and multiple linear regression [145]. Slightly better results were obtained with SVM, associative neural networks, and MLF artificial neural networks.

The human serum protein binding of 808 compounds was simulated with multiple linear regression, MLF artificial neural networks,  $k$ -NN, and SVM [148]. An ensemble of MLF artificial neural networks provided the best predictions for an external validation set of 200 compounds. The blood-brain barrier permeability of 415 chemicals was modeled by Li et al with logistic regression, linear discriminant analysis,  $k$ -NN, C4.5 decision tree, probabilistic neural network, and SVM [89]. It was found that all ML models are improved by using a descriptor set selected with recursive feature elimination. Plewczynski et al compared SAR models for five biological targets obtained with SVM, random forest, artificial neural networks,  $k$ -NN, trend vectors, Bayesian classifier, and decision tree [107]. There are significant differences in performance depending on the target and objective, i. e., high enrichment or maximum number of actives. The toxicity risk assessment of organic chemicals was evaluated with SVM,  $k$ -NN, Bayesian classifier, and self-organizing maps (SOM) [147]. Both SOM and SVM can separate toxic from nontoxic compounds based on structural fragments. A new ML algorithm that shows promising results is kScore, which in several tests performed better than SVM,  $k$ -NN, recursive partitioning, artificial neural networks, Gaussian process, and Bayesian classifier [102].

A comparison of 21 machine learning algorithms is presented for data selected from the National Cancer Institute 60-cell line screening panel (NCI-60) [69]. The SMILES codes of the chemical compounds [130] were used to compute fingerprints with Open Babel (<http://openbabel.org/>). The descriptor selection (based on Cfs-SubsetEval and BestFirst) and all machine learning models were computed with Weka [41,151], and all prediction results are for 10-fold (leave-10%-out) cross-validation. The machine learning algorithms used are briefly listed here, using their notation in Weka: BayesNet, Bayesian network; NaiveBayes, naïve Bayesian classifier [71]; NaiveBayesUpdateable, naïve Bayesian classifier with estimator classes [71]; Logistic, logistic regression with a ridge estimator [86]; RBFNetwork, Gaussian radial basis function network; IBk,  $k$ -NN classifier with distance weight (NoW – no distance weight,  $W(1/d)$  – weighted with  $1/d$ ,  $W(1 - d)$  – weighted with  $(1 - d)$ ) and  $k$  an odd

number between 1 and 9 [1]; KStar,  $K^*$  lazy learner with entropy-based distance function [26]; ADTree, alternating decision tree (search type: Eap, expand all paths; Ehp, expand the heaviest path; Ezp, expand the best  $z$ -pure path; Erp, expand a random path) [42]; DecisionStump, one-level binary decision tree with categorical or numerical class label; J48, C4.5 decision tree [108]; NBTree, decision tree with naïve Bayes classifiers at the leaves [82]; RandomForest, random forest (the number of random trees is between 10 and 50) [22]; RandomTree, a tree that considers  $k$  randomly chosen attributes at each node; REP-Tree, fast decision tree learner; ConjunctiveRule, conjunctive rule learner; DecisionTable, decision table majority classifier [81]; JRip, a propositional rule learner based on RIPPER [27]; OneR, rule classifier that uses the minimum-error attribute for prediction [56]; PART, a PART decision list that builds a partial C4.5 decision tree in each iteration and transforms the best leaf into a rule; Ridor, a Ripple-Down Rule learner [43]; SVM, LibSVM. More details for each ML may be found in Weka.

The first series of SAR models for anticancer compounds considers the cell line NCI-H460 (lung large cell carcinoma) with 3599 chemicals (2049 in class +1, 1550 in class -1) and 30 descriptors (Tab. 2). For an easier comparison of the ranking of all ML methods, the results are ordered after the Matthews correlation coefficient MCC. The prediction statistics show that  $k$ -NN, in its several variants evaluated here, constantly gives the best predictions. Three other ML with good predictions are KStar, SVM, and RandomForest.

The second group of structure-activity models is obtained for the cell line SF-268 (glioma) with 3721 chemicals (2020 in class +1, 1701 in class -1) and 36 descriptors (Tab. 3). A  $k$ -NN classifier is in the first position, followed by SVM and other  $k$ -NN variants. In decreasing order of the predictions, the list continues with another lazy learning method, KStar, followed by RandomForest and J48.

The third set of experiments considers the cell line SK-MEL-5 (melanoma) with 3685 chemicals (2034 in class +1, 1651 in class -1) and 30 descriptors (Tab. 4). All  $k$ -NN give good predictions, including the top five results. RandomForest ranks better compared with the previous two cell lines, and it is followed by SVM and KStar. As a conclusion of these three sets of experiments, we find that a small group of ML algorithms gives constantly the best predictions. The top results are obtained with  $k$ -NN, which has good statistics for the whole range of parameters tested here. This is an important result, because  $k$ -NN is a robust method that is much easier to compute than SVM. Three other ML algorithms have good predictions, namely SVM, KStar, and RandomForest. Although KStar

**Drug Design with Machine Learning, Table 2**  
**Machine learning anticancer SAR for lung large cell carcinoma NCI-H460**

Exp	Model	TP <sub>p</sub>	FN <sub>p</sub>	TN <sub>p</sub>	FP <sub>p</sub>	Ac <sub>p</sub>	MCC <sub>p</sub>
1	IBk $k = 9W(1 - d)$	1386	663	1035	515	0.6727	0.3414
2	IBk $k = 5W(1 - d)$	1410	639	1011	539	0.6727	0.3383
3	IBk $k = 9$ NoW	1404	645	1015	535	0.6721	0.3378
4	IBk $k = 7W(1 - d)$	1386	663	1027	523	0.6705	0.3364
5	IBk $k = 9W(1/d)$	1399	650	1012	538	0.6699	0.3334
6	IBk $k = 5$ NoW	1433	616	983	567	0.6713	0.3324
7	KStar	1364	685	1034	516	0.6663	0.3299
8	IBk $k = 7$ NoW	1404	645	1001	549	0.6682	0.3290
9	SVM RBF $\sigma = 0.01$	1494	555	926	624	0.6724	0.3286
10	IBk $k = 7W(1/d)$	1396	653	1003	547	0.6666	0.3263
11	IBk $k = 3W(1 - d)$	1420	629	982	568	0.6674	0.3252
12	RandomForest $T = 40$	1509	540	902	648	0.6699	0.3217
13	IBk $k = 3$ NoW	1446	603	954	596	0.6669	0.3210
14	IBk $k = 5W(1/d)$	1408	641	985	565	0.6649	0.3210
15	RandomForest $T = 50$	1499	550	903	647	0.6674	0.3171
16	RandomForest $T = 30$	1499	550	902	648	0.6671	0.3164
17	IBk $k = 3W(1/d)$	1428	621	955	595	0.6621	0.3125
18	RandomForest $T = 20$	1487	562	889	661	0.6602	0.3021
19	IBk $k = 1$ NoW	1457	592	911	639	0.6580	0.3000
20	IBk $k = 1W(1/d)$	1457	592	911	639	0.6580	0.3000
21	IBk $k = 1W(1 - d)$	1457	592	911	639	0.6580	0.3000
22	RandomForest $T = 10$	1494	555	878	672	0.6591	0.2990
23	JRip	1568	481	810	740	0.6607	0.2972
24	J48	1481	568	884	666	0.6571	0.2959
25	ADTree Ehp	1536	513	836	714	0.6591	0.2956
26	ADTree Ezp	1536	513	836	714	0.6591	0.2956
27	ADTree Erp	1580	469	791	759	0.6588	0.2922
28	NaiveBayes	1130	919	1147	403	0.6327	0.2919
29	NaiveBayesUpdateable	1130	919	1147	403	0.6327	0.2919
30	ADTree Eap	1563	486	806	744	0.6582	0.2919
31	Logistic	1514	535	848	702	0.6563	0.2911
32	BayesNet	1256	793	1054	496	0.6418	0.2903
33	NBTree	1406	643	930	620	0.6491	0.2857
34	RBFNetwork	1578	471	769	781	0.6521	0.2774
35	PART	1404	645	917	633	0.6449	0.2766
36	REPTree	1551	498	785	765	0.6491	0.2723
37	DecisionTable	1549	500	752	798	0.6393	0.2507
38	RandomTree	1384	665	885	665	0.6305	0.2464
39	Ridor	1617	432	642	908	0.6277	0.2201
40	OneR	1745	304	427	1123	0.6035	0.1565
41	ConjunctiveRule	1899	150	163	1387	0.5729	0.0562
42	DecisionStump	2049	0	0	1550	0.5693	0.0000

1864 is not used in drug design, the results reported here indi-  
 1865 cate that this ML algorithm should be added to the toolbox  
 1866 of structure-activity methods. We have to note the absence  
 1867 of the Bayesian classifiers, although they are the preferred

method in some library screening experiments.

The studies reviewed in this section show that it is dif-  
 ficult to perform comprehensive comparisons for a large  
 number of machine learning algorithms. To address this

1868

1869

1870

1871

**Drug Design with Machine Learning, Table 3**  
**Machine learning anticancer SAR for glioma SF-268**

Exp	Model	TP <sub>p</sub>	FN <sub>p</sub>	TN <sub>p</sub>	FP <sub>p</sub>	Ac <sub>p</sub>	MCC <sub>p</sub>
1	IBk $k = 9W(1/d)$	1338	682	1177	524	0.6759	0.3530
2	SVM RBF $\sigma = 0.01$	1403	617	1118	583	0.6775	0.3513
3	IBk $k = 7W(1/d)$	1344	676	1165	536	0.6743	0.3490
4	IBk $k = 5W(1/d)$	1351	669	1152	549	0.6727	0.3449
5	IBk $k = 7W(1 - d)$	1325	695	1173	528	0.6713	0.3443
6	IBk $k = 9W(1 - d)$	1301	719	1185	516	0.6681	0.3395
7	IBk $k = 5W(1 - d)$	1330	690	1160	541	0.6692	0.3391
8	IBk $k = 5$ NoW	1372	648	1120	581	0.6697	0.3368
9	IBk $k = 7$ NoW	1353	667	1135	566	0.6686	0.3360
10	KStar	1307	713	1168	533	0.6651	0.3325
11	IBk $k = 3$ NoW	1396	624	1089	612	0.6678	0.3311
12	IBk $k = 3W(1 - d)$	1358	662	1121	580	0.6662	0.3304
13	IBk $k = 9$ NoW	1323	697	1150	551	0.6646	0.3298
14	RandomForest $T = 40$	1408	612	1074	627	0.6670	0.3287
15	RandomForest $T = 50$	1408	612	1072	629	0.6665	0.3275
16	RandomForest $T = 30$	1418	602	1062	639	0.6665	0.3269
17	IBk $k = 3W(1/d)$	1363	657	1108	593	0.6641	0.3254
18	J48	1376	644	1082	619	0.6606	0.3169
19	RandomForest $T = 20$	1403	617	1045	656	0.6579	0.3095
20	Logistic	1404	616	1028	673	0.6536	0.3003
21	NBTree	1240	780	1164	537	0.6461	0.2974
22	RandomForest $T = 10$	1411	609	1015	686	0.6520	0.2965
23	BayesNet	1220	800	1176	525	0.6439	0.2948
24	IBk $k = 1$ NoW	1396	624	1020	681	0.6493	0.2916
25	IBk $k = 1W(1/d)$	1396	624	1020	681	0.6493	0.2916
26	IBk $k = 1W(1 - d)$	1396	624	1020	681	0.6493	0.2916
27	ADTree Erp	1407	613	993	708	0.6450	0.2819
28	RBFNetwork	1519	501	884	817	0.6458	0.2800
29	NaiveBayes	1064	956	1272	429	0.6278	0.2790
30	NaiveBayesUpdateable	1064	956	1272	429	0.6278	0.2790
31	ADTree Ehp	1406	614	986	715	0.6428	0.2774
32	ADTree Ezp	1406	614	986	715	0.6428	0.2774
33	REPTree	1380	640	1009	692	0.6420	0.2771
34	JRip	1418	602	960	741	0.6391	0.2689
35	ADTree Eap	1443	577	936	765	0.6393	0.2684
36	PART	1328	692	1024	677	0.6321	0.2593
37	DecisionTable	1438	582	915	786	0.6324	0.2538
38	RandomTree	1357	663	984	717	0.6291	0.2510
39	Ridor	1705	315	630	1071	0.6275	0.2454
40	DecisionStump	674	1346	1414	287	0.5611	0.1877
41	ConjunctiveRule	665	1355	1419	282	0.5601	0.1869
42	OneR	1257	763	932	769	0.5883	0.1702

1872 issue, we organized a competition for the ML evaluation in  
 1873 blind predictions, CoEPrA 2006 (Comparative Evaluation  
 1874 of Prediction Algorithms, <http://www.coepra.org/>). For  
 1875 each prediction task the participants received a training

dataset (structure, descriptors, and class attribution or ac-  
 tivity) and a prediction dataset (structure and descriptors).  
 All predictions received before deadline were compared  
 with the experimental results, and then the ML models

1876  
 1877  
 1878  
 1879

**Drug Design with Machine Learning, Table 4**  
**Machine learning anticancer SAR for melanoma SK-MEL-5**

Exp	Model	$TP_p$	$FN_p$	$TN_p$	$FP_p$	$Ac_p$	$MCC_p$
1	IBk $k = 9W(1 - d)$	1379	655	1119	532	0.6779	0.3541
2	IBk $k = 9$ NoW	1401	633	1099	552	0.6784	0.3532
3	IBk $k = 7$ NoW	1422	612	1081	570	0.6792	0.3531
4	IBk $k = 5W(1 - d)$	1401	633	1095	556	0.6773	0.3508
5	IBk $k = 7W(1 - d)$	1390	644	1104	547	0.6768	0.3506
6	RandomForest $T = 50$	1487	547	1019	632	0.6801	0.3504
7	IBk $k = 7W(1/d)$	1414	620	1079	572	0.6765	0.3479
8	IBk $k = 5W(1/d)$	1417	617	1074	577	0.6760	0.3465
9	IBk $k = 9W(1/d)$	1404	630	1084	567	0.6752	0.3458
10	RandomForest $T = 30$	1478	556	1019	632	0.6776	0.3456
11	IBk $k = 5$ NoW	1427	607	1063	588	0.6757	0.3451
12	RandomForest $T = 20$	1482	552	1013	638	0.6771	0.3443
13	RandomForest $T = 40$	1481	553	1013	638	0.6768	0.3437
14	SVM RBF $\sigma = 0.03$	1414	620	1061	590	0.6716	0.3373
15	KStar	1326	708	1119	532	0.6635	0.3279
16	RandomForest $T = 10$	1481	553	976	675	0.6668	0.3222
17	IBk $k = 3W(1 - d)$	1390	644	1053	598	0.6630	0.3204
18	IBk $k = 3$ NoW	1416	618	1025	626	0.6624	0.3171
19	IBk $k = 3W(1/d)$	1393	641	1044	607	0.6613	0.3166
20	IBk $k = 1$ NoW	1427	607	987	664	0.6551	0.3005
21	IBk $k = 1W(1/d)$	1427	607	987	664	0.6551	0.3005
22	IBk $k = 1W(1 - d)$	1427	607	987	664	0.6551	0.3005
23	BayesNet	1232	802	1132	519	0.6415	0.2901
24	PART	1348	686	1033	618	0.6461	0.2875
25	RBFNetwork	1541	493	857	794	0.6507	0.2856
26	NBTree	1313	721	1055	596	0.6426	0.2832
27	J48	1378	656	1000	651	0.6453	0.2831
28	Logistic	1426	608	958	693	0.6469	0.2830
29	JRip	1483	551	899	752	0.6464	0.2785
30	RandomTree	1401	633	963	688	0.6415	0.2731
31	NaiveBayes	1068	966	1224	427	0.6220	0.2698
32	NaiveBayesUpdateable	1068	966	1224	427	0.6220	0.2698
33	DecisionTable	1494	540	863	788	0.6396	0.2634
34	REPTree	1431	603	918	733	0.6374	0.2622
35	ADTree Erp	1436	598	889	762	0.6309	0.2478
36	ADTree Ehp	1429	605	892	759	0.6299	0.2459
37	ADTree Ezp	1429	605	892	759	0.6299	0.2459
38	ADTree Eap	1323	711	981	670	0.6252	0.2441
39	Ridor	1469	565	817	834	0.6204	0.2230
40	DecisionStump	1279	755	932	719	0.6000	0.1930
41	OneR	1279	755	932	719	0.6000	0.1930
42	ConjunctiveRule	1404	630	770	881	0.5900	0.1605

1880 were ranked after different statistical indices (Tab. 5). The  
 1881 competition had four classification (C in Tab. 5) tasks and  
 1882 four regression (R in Tab. 5) tasks. All datasets consisted  
 1883 of peptides evaluated for their binding affinity to the ma-

1884 jor histocompatibility complex (MHC) proteins. To evalu-  
 1885 ate the general performance of the algorithms tested in  
 1886 competition, for each task we present the ML from the top  
 1887 five positions. In general, each task has a different win-



**Drug Design with Machine Learning, Table 5**

Machine learning algorithms with best results in the CoEPrA 2006 (comparative evaluation of prediction algorithms, <http://www.coepra.org/>) competition

Task	Rank	Machine Learning	Task	Rank	Machine Learning
C1	1	naïve Bayes	R1	1	LS-SVM linear kernel
	2	LS-SVM RBF kernel		2	random forest
	3	CART		3	Gaussian process
	4	C4.5		4	kernel PLS
	5	SVM string kernel		5	SVM RBF kernel
C2	1	SVM linear kernel	R2	1	kernel PLS
	2	kScore		2	kNN
	3	SVM RBF kernel		3	PLS
	4	ClassificationViaRegression		4	SVM string kernel
	5	random forest		5	kScore
C3	1	SVM linear kernel	R3	1	random forest
	2	LS-SVM RBF kernel		2	kernel PLS
	3	SVM		3	SVM string kernel
	4	C4.5		4	PLS
	5	random forest		5	Gaussian process
C4	1	LS-SVM quadratic kernel	R3D2	1	kernel PLS
	2	SVM with AA sequence		2	SVM string kernel
	3	SVM binary encoding of AA		3	random forest
	4	kScore		4	Gaussian process
	5	random forest		5	SVM binary encoding of AA

ner and order for the top ML. The most frequent ML is SVM, in several variants. The least squares SVM (LS-SVM) is particularly robust, insensitive to noise, and with good performance. Currently, LS-SVM is not used in virtual screening, but the results from the competition show that it should be included in the set of preferred methods. Other predictive ML methods are random forest, Gaussian process, and kScore. The CoEPrA offered interesting conclusions with direct implications in drug design, SAR, and QSAR. Similar experiments are necessary to evaluate datasets consisting of drugs and drug-like chemicals.

**Future Directions**

For a long period of time the machine learning selection available to the computational chemist was restricted to a small number of algorithms. However, recent publications explore a much larger diversity of statistical models, with the intention to help the drug discovery process with reliable predictions. As an example of the current trend, one should consider the fast transition of the support vector machines from small-scale experiments to the preferred method in many industrial applications. Following the SVM success, other kernel methods were adopted in chemoinformatics and drug design. The results reviewed

here show that it is not possible to predict which combination of structural descriptors and machine learning will give the most reliable predictions for novel chemicals. Therefore, in order to identify good models, it is always a good idea to compare several types of descriptors and as many ML algorithms as computationally possible. Equally important is a good plan for the experiments, that includes a proper validation and whenever possible tests with external data that were not used in training or cross-validation. Of great importance is an objective and comprehensive evaluation of ML algorithms, which can be performed in settings similar with the CoEPrA competition. Such blind predictions may offer an unbiased ranking of the machine learning algorithms.

**Bibliography**

1. Aha DW, Kibler D, Albert MK (1991) Instance-based learning algorithms. *Mach Learn* 6:37–66
2. Ajmani S, Jadhav K, Kulkarni SA (2006) Three-dimensional QSAR using the *k*-nearest neighbor method and its interpretation. *J Chem Inf Model* 46:24–31
3. Andres C, Hutter MC (2006) CNS permeability of drugs predicted by a decision tree. *QSAR Comb Sci* 25:305–309
4. Alpaydin E (2004) Introduction to machine learning. MIT Press, Cambridge, p 445

- 1935 5. Atkeson CG, Moore AW, Schaal S (1997) Locally weighted  
1936 learning. *Artif Intell Rev* 11:11–73
- 1937 6. Atkeson CG, Moore AW, Schaal S (1997) Locally weighted  
1938 learning for control. *Artif Intell Rev* 11:75–113
- 1939 7. Arimoto R, Prasad MA, Gifford EM (2005) Development of  
1940 CYP3A4 inhibition models: comparisons of machine-learning  
1941 techniques and molecular descriptors. *J Biomol Screen*  
1942 10:197–205
- 1943 8. Balaban AT, Ivanciuc O (1999) Historical development of  
1944 topological indices. In: Devillers J, Balaban AT (eds) *Topologi-*  
1945 *cal indices and related descriptors in QSAR and QSPR*. Gordon  
1946 & Breach Science Publishers, Amsterdam, pp 21–57
- 1947 9. Basak SC, Grunwald GD (1995) Molecular similarity and es-  
1948 timation of molecular properties. *J Chem Inf Comput Sci*  
1949 35:366–372
- 1950 10. Basak SC, Bertelsen S, Grunwald GD (1994) Application of  
1951 graph theoretical parameters in quantifying molecular simi-  
1952 larity and structure-activity relationships. *J Chem Inf Comput*  
1953 *Sci* 34:270–276
- 1954 11. Basak SC, Bertelsen S, Grunwald GD (1995) Use of graph theo-  
1955 retic parameters in risk assessment of chemicals. *Toxicol Lett*  
1956 79:239–250
- 1957 12. Bayes T (1763) An essay towards solving a problem in the doc-  
1958 trine of chances. *Philos Trans Roy Soc London* 53:370–418
- 1959 13. Bender A, Jenkins JL, Glick M, Deng Z, Nettles JH, Davies  
1960 JW (2006) “Bayes affinity fingerprints” improve retrieval rates  
1961 in virtual screening and define orthogonal bioactivity space:  
1962 when are multitarget drugs a feasible concept? *J Chem Inf*  
1963 *Model* 46:2445–2456
- 1964 14. Bender A, Scheiber J, Glick M, Davies JW, Azzaoui K, Hamon J,  
1965 Urban L, Whitebread S, Jenkins JL (2007) Analysis of pharma-  
1966 cology data and the prediction of adverse drug reactions and  
1967 off-target effects from chemical structure. *Chem Med Chem*  
1968 2:861–873
- 1969 15. Bishop CM (2006) *Pattern recognition and machine learning*.  
1970 Springer, Berlin, p 740
- 1971 16. Bishop CM (1996) *Neural networks for pattern recognition*.  
1972 Oxford University Press, Oxford, p 504
- 1973 17. Boid DB (2007) How computational chemistry became impor-  
1974 tant in the pharmaceutical industry. In: Lipkowitz KB, Cundari  
1975 TR (eds) *Reviews in computational chemistry*, vol 23. Wiley,  
1976 Weinheim, pp 401–451
- 1977 18. Bonchev D (1983) *Information theoretic indices for character-*  
1978 *ization of chemical structure*. Research Studies Press, Chich-  
1979 ester
- 1980 19. Bonchev D, Rouvray DH (eds) (1991) *Chemical graph theo-*  
1981 *ry. Introduction and fundamentals*. Abacus Press/Gordon &  
1982 Breach Science Publishers, New York
- 1983 20. Boser BE, Guyon IM, Vapnik VN (1992) A training algorithm  
1984 for optimal margin classifiers. In: Haussler D (ed) *Proc of the*  
1985 *5th annual ACM workshop on computational learning theory*.  
1986 ACM Press, Pittsburgh, pp 144–152
- 1987 21. Bottou L, Chapelle O, DeCoste D, Weston J (2007) *Large-scale*  
1988 *kernel machines*. MIT Press, Cambridge, p 416
- 1989 22. Breiman L (2001) Random forests. *Mach Learn* 45:5–32
- 1990 23. Briem H, Günther J (2005) Classifying “kinase inhibitor-like-  
1991 ness” by using machine-learning methods. *Chem Bio Chem*  
1992 6:558–566
- 1993 24. Cash GG (1999) Prediction of physicochemical properties  
1994 from Euclidean distance methods based on electrotopologi-  
1995 cal state indices. *Chemosphere* 39:2583–2591
25. Chapelle O, Haffner P, Vapnik VN (1999) Support vector ma-  
1996 chines for histogram-based image classification. *IEEE Trans*  
1997 *Neural Netw* 10:1055–1064
- 1998 26. Cleary JG, Trigg LE (1995)  $K^*$ : an instance-based learner us-  
1999 ing and entropic distance measure. In: Prieditis A, Russell SJ  
2000 (eds) *Proc of the 12th international conference on machine*  
2001 *learning*. Morgan Kaufmann, Tahoe City, pp 108–114
- 2002 27. Cohen WW (1995) Fast effective rule induction. In: Prieditis A,  
2003 Russell SJ (eds) *Proc of the 12th international conference on*  
2004 *machine learning*. Morgan Kaufmann, Tahoe City, pp 115–  
2005 123
- 2006 28. Cortes C, Vapnik V (1995) Support vector networks. *Mach*  
2007 *Learn* 20:273–297
- 2008 29. Cristianini N, Shawe-Taylor J (2000) *An introduction to sup-*  
2009 *port vector machines*. Cambridge University Press, Cam-  
2010 bridge
- 2011 30. Deconinck E, Zhang MH, Coomans D, Vander Heyden Y  
2012 (2006) Classification tree models for the prediction of blood-  
2013 brain barrier passage of drugs. *J Chem Inf Model* 46:1410–  
2014 1419
- 2015 31. Deng Z, Chuaqui C, Singh J (2006) Knowledge-based design  
2016 of target-focused libraries using protein-ligand interaction  
2017 constraints. *J Med Chem* 49:490–500
- 2018 32. Doddareddy MR, Cho YS, Koh HY, Kim DH, Pae AN (2006) In  
2019 silico renal clearance model using classical Volsurf approach.  
2020 *J Chem Inf Model* 46:1312–1320
- 2021 33. Drucker H, Wu DH, Vapnik VN (1999) Support vector ma-  
2022 chines for spam categorization. *IEEE Trans Neural Netw*  
2023 10:1048–1054
- 2024 34. Du H, Wang J, Watzl J, Zhang X, Hu Z (2008) Classification  
2025 structure-activity relationship (CSAR) studies for prediction  
2026 of genotoxicity of thiophene derivatives. *Toxicol Lett* 177:  
2027 10–19
- 2028 35. Duda RO, Hart PE, Stork DG (2000) *Pattern classification*. 2nd  
2029 edn. Wiley, New York
- 2030 36. Ehrman TM, Barlow DJ, Hylands PJ (2007) Virtual screening of  
2031 chinese herbs with random forest. *J Chem Inf Model* 47:264–  
2032 278
- 2033 37. Eitrich T, Kless A, Druska C, Meyer W, Grotendorst J (2007)  
2034 Classification of highly unbalanced CYP450 data of drugs us-  
2035 ing cost sensitive machine learning techniques. *J Chem Inf*  
2036 *Model* 47:92–103
- 2037 38. Ekins S, Balakin KV, Savchuk N, Ivanenkov Y (2006) Insights  
2038 for human ether-a-go-go-related gene potassium channel in-  
2039 hibition using recursive partitioning and Kohonen and Sam-  
2040 mon mapping techniques. *J Med Chem* 49:5059–5071
- 2041 39. Ertl P, Roggo S, Schuffenhauer A (2008) Natural product-like-  
2042 ness score and its application for prioritization of compound  
2043 libraries. *J Chem Inf Model* 48:68–74
- 2044 40. Fatemi MH, Gharaghani S (2007) A novel QSAR model  
2045 for prediction of apoptosis-inducing activity of 4-aryl-4-  
2046 H-chromenes based on support vector machine. *Bioorg Med*  
2047 *Chem* 15:7746–7754
- 2048 41. Frank E, Hall M, Trigg L, Holmes G, Witten IH (2004) Data min-  
2049 ing in bioinformatics using Weka. *Bioinformatics* 20:2479–  
2050 2481
- 2051 42. Freund Y, Mason L (1999) The alternating decision tree learn-  
2052 ing algorithm. In: Bratko I, Dzeroski S (eds) *Proc of the 16th*  
2053 *international conference on machine learning (ICML (1999))*.  
2054 Morgan Kaufmann, Bled, pp 124–133
- 2055 43. Gaines BR, Compton P (1995) Induction of ripple-down rules  
2056

- 2057 applied to modeling large databases. *Intell J Inf Syst* 5:211–  
2058 228
- 2059 44. Gao JB, Gunn SR, Harris CJ (2003) SVM regression through  
2060 variational methods and its sequential implementation. *Neurocomputing* 55:151–167  
2061
- 2062 45. Gao JB, Gunn SR, Harris CJ (2003) Mean field method for  
2063 the support vector machine regression. *Neurocomputing*  
2064 50:391–405
- 2065 46. Gepp MM, Hutter MC (2006) Determination of hERG channel  
2066 blockers using a decision tree. *Bioorg Med Chem* 14:5325–  
2067 5332
- 2068 47. Guha R, Dutta D, Jurs PC, Chen T (2006) Local lazy regression:  
2069 making use of the neighborhood to improve QSAR predic-  
2070 tions. *J Chem Inf Model* 46:1836–1847
- 2071 48. Gute BD, Basak SC (2001) Molecular similarity-based estima-  
2072 tion of properties: a comparison of three structure spaces.  
2073 *J Mol Graph Modell* 20:95–109
- 2074 49. Gute BD, Basak SC, Mills D, Hawkins DM (2002) Tailored simi-  
2075 larity spaces for the prediction of physicochemical properties.  
2076 *Internet Electron J Mol Des* 1:374–387
- 2077 50. Guyon I, Weston J, Barnhill S, Vapnik V (2002) Gene selection  
2078 for cancer classification using support vector machines. *Mach*  
2079 *Learn* 46:389–422
- 2080 51. Hansch C, Garg R, Kurup A, Mekapati SB (2003) Allosteric in-  
2081 teractions and QSAR: on the role of ligand hydrophobicity.  
2082 *Bioorg Med Chem* 11:2075–2084
- 2083 52. Hastie T, Tibshirani R, Friedman JH (2003) The elements of sta-  
2084 tistical learning. Springer, Berlin, p 552
- 2085 53. Herbrich R (2002) Learning kernel classifiers. MIT Press, Cam-  
2086 bridge
- 2087 54. Hert J, Willett P, Wilton DJ, Acklin P, Azzaoui K, Jacoby E,  
2088 Schuffenhauer A (2006) New methods for ligand-based vir-  
2089 tual screening: use of data fusion and machine learning to  
2090 enhance the effectiveness of similarity searching. *J Chem Inf*  
2091 *Model* 46:462–470
- 2092 55. Hoffman B, Cho SJ, Zheng W, Wyrick S, Nichols DE, Mail-  
2093 man RB, Tropsha A (1999) Quantitative structure-activity re-  
2094 lationship modeling of dopamine D<sub>1</sub> antagonists using com-  
2095 parative molecular field analysis, genetic algorithms-partial  
2096 least-squares, and *K*-nearest neighbor methods. *J Med Chem*  
2097 42:3217–3226
- 2098 56. Holte RC (1993) Very simple classification rules perform  
2099 well on most commonly used datasets. *Mach Learn* 11:63–  
2100 90
- 2101 57. Hou T, Wang J, Zhang W, Xu X (2007) ADME evaluation in  
2102 drug discovery. 7. Prediction of oral absorption by correlation  
2103 and classification. *J Chem Inf Model* 47:208–218
- 2104 58. Huang T-M, Kecman V, Kopriva I (2006) Kernel based algo-  
2105 rithms for mining huge data sets. Springer, Berlin, p 260
- 2106 59. Hudelson MG, Ketkar NS, Holder LB, Carlson TJ, Peng C-C,  
2107 Waldher BJ, Jones JP (2008) High confidence predictions of  
2108 drug-drug interactions: predicting affinities for cytochrome  
2109 P450 2C9 with multiple computational methods. *J Med Chem*  
2110 51:648–654
- 2111 60. Itskowitz P, Tropsha A (2005) *k*-nearest neighbors QSAR  
2112 modeling as a variational problem: theory and applications.  
2113 *J Chem Inf Model* 45:777–785
- 2114 61. Ivanciuc O (2002) Support vector machine classification of the  
2115 carcinogenic activity of polycyclic aromatic hydrocarbons. *Inter-  
2116 net Electron J Mol Des* 1:203–218
- 2117 62. Ivanciuc O (2002) Structure-odor relationships for pyrazines  
with support vector machines. *Internet Electron J Mol Des* 1:269–284 2118
63. Ivanciuc O (2002) Support vector machine identification of the  
2119 aquatic toxicity mechanism of organic compounds. *Inter-  
2120 net Electron J Mol Des* 1:157–172 2121
64. Ivanciuc O (2003) Graph theory in chemistry. In: Gasteiger J  
2122 (ed) *Handbook of chemoinformatics*, vol 1. Wiley, Weinheim,  
2123 pp 103–138 2124
65. Ivanciuc O (2003) Topological indices. In: Gasteiger J (ed)  
2125 *Handbook of chemoinformatics*, vol 3. Wiley, Weinheim,  
2126 pp 981–1003 2127
66. Ivanciuc O (2003) Aquatic toxicity prediction for polar and  
2128 nonpolar narcotic pollutants with support vector machines.  
2129 *Internet Electron J Mol Des* 2:195–208 2130
67. Ivanciuc O (2004) Support vector machines prediction of the  
2131 mechanism of toxic action from hydrophobicity and experi-  
2132 mental toxicity against pimephales promelas and tetrahy-  
2133 mena pyriformis. *Internet Electron J Mol Des* 3:802–821 2134
68. Ivanciuc O (2005) Support vector regression quantitative  
2135 structure-activity relationships (QSAR) for benzodiazepine re-  
2136 ceptor ligands. *Internet Electron J Mol Des* 4:181–193 2137
69. Ivanciuc O (2005) Machine learning applied to anticancer  
2138 structure-activity relationships for NCI human tumor cell  
2139 lines. *Internet Electron J Mol Des* 4:948–958 2140
70. Ivanciuc O (2007) Applications of support vector machines in  
2141 chemistry. In: Lipkowitz KB, Cundari TR (eds) *Reviews in com-  
2142 putational chemistry*, vol 23. Wiley, Weinheim, pp 291–400 2143
71. John GH, Langley P (1995) Estimating continuous distribu-  
2144 tions in Bayesian classifiers. In: Besnard P, Hanks S (eds) *UAI*  
2145 '95: Proc of the 11th annual conference on uncertainty in ar-  
2146 tificial intelligence. Morgan Kaufmann, Montreal, pp 338–345 2147
72. Jorissen RN, Gilson MK (2005) Virtual screening of molecular  
2148 databases using a support vector machine. *J Chem Inf Model*  
2149 45:549–561 2150
73. Jurs P (2003) Quantitative structure-property relationships.  
2151 In: Gasteiger J (ed) *Handbook of chemoinformatics*, vol 3. Wil-  
2152 ley, Weinheim, pp 1314–1335 2153
74. Kier LB, Hall LH (1976) Molecular connectivity in chemistry  
2154 and drug research. Academic Press, New York 2155
75. Kier LB, Hall LH (1986) Molecular connectivity in structure-  
2156 activity analysis. Research Studies Press, Letchworth 2157
76. Kier LB, Hall LH (1999) Molecular structure description. The  
2158 electrotopological state. Academic Press, San Diego 2159
77. Klon AE, Diller DJ (2007) Library fingerprints: a novel ap-  
2160 proach to the screening of virtual libraries. *J Chem Inf Model*  
2161 47:1354–1365 2162
78. Klon AE, Glick M, Davies JW (2004) Combination of a naive  
2163 Bayes classifier with consensus scoring improves enrichment  
2164 of high-throughput docking results. *J Med Chem* 47:4356–  
2165 4359 2166
79. Klon AE, Glick M, Thoma M, Acklin P, Davies JW (2004) Finding  
2167 more needles in the haystack: a simple and efficient method  
2168 for improving high-throughput docking results. *J Med Chem*  
2169 47:2743–2749 2170
80. Klon AE, Lowrie JF, Diller DJ (2006) Improved naive Bayesian  
2171 modeling of numerical data for absorption, distribution,  
2172 metabolism and excretion (ADME) property prediction.  
2173 *J Chem Inf Model* 46:1945–1956 2174
81. Kohavi R (1995) The power of decision tables. In: Lavrac N,  
2175 Wrobel S (eds) *ECML-95 8th european conference on ma-  
2176 chine learning*, vol 912. Springer, Heraclion, pp 174–189 2177

**CE2** Please provide publishing location.

- 2179 82. Kohavi R (1996) Scaling up the accuracy of naive-Bayes classifiers: a decision-tree hybrid. In: Simoudis E, Han J, Fayyad UM (eds) Proc of the 2nd international conference on knowledge discovery and data mining (KDD-96). AAAI Press, pp 202–207 **CE2** 2240
- 2180 2241
- 2181 2242
- 2182 2243
- 2183 2244
- 2184 83. Kononenko I, Kukar M (2007) Machine learning and data mining: introduction to principles and algorithms. Horwood, Westergate, p 454 2245
- 2185 2246
- 2186 2247
- 2187 84. Konovalov DA, Coomans D, Deconinck E, Vander Heyden Y (2007) Benchmarking of QSAR models for blood-brain barrier permeation. *J Chem Inf Model* 47:1648–1656 2248
- 2188 2249
- 2189 2250
- 2190 85. Kumar R, Kulkarni A, Jayaraman VK, Kulkarni BD (2004) Structure-activity relationships using locally linear embedding assisted by support vector and lazy learning regressors. *Internet Electron J Mol Des* 3:118–133 2251
- 2191 2252
- 2192 2253
- 2193 2254
- 2194 86. le Cessie S, van Houwelingen JC (1992) Ridge estimators in logistic regression. *Appl Statist* 41:191–201 2255
- 2195 2256
- 2196 87. Leong MK (2007) A novel approach using pharmacophore ensemble/support vector machine (PhE/SVM) for prediction of hERG liability. *Chem Res Toxicol* 20:217–226 2257
- 2197 2258
- 2198 88. Lepp Z, Kinoshita T, Chuman H (2006) Screening for new antidepressant leads of multiple activities by support vector machines. *J Chem Inf Model* 46:158–167 2259
- 2199 2260
- 2200 2261
- 2201 2262
- 2202 89. Li H, Yap CW, Ung CY, Xue Y, Cao ZW, Chen YZ (2005) Effect of selection of molecular descriptors on the prediction of blood-brain barrier penetrating and nonpenetrating agents by statistical learning methods. *J Chem Inf Model* 45:1376–1384 2263
- 2203 2264
- 2204 2265
- 2205 2266
- 2206 2267
- 2207 90. Li S, Fedorowicz A, Singh H, Soderholm SC (2005) Application of the random forest method in studies of local lymph node assay based skin sensitization data. *J Chem Inf Model* 45:952–964 2268
- 2208 2269
- 2209 2270
- 2210 2271
- 2211 91. Li W-X, Li L, Eksterowicz J, Ling XB, Cardozo M (2007) Significance analysis and multiple pharmacophore models for differentiating *P*-glycoprotein substrates. *J Chem Inf Model* 47:2429–2438 2272
- 2212 2273
- 2213 2274
- 2214 2275
- 2215 92. Liao Q, Yao J, Yuan S (2007) Prediction of mutagenic toxicity by combination of recursive partitioning and support vector machines. *Mol Divers* 11:59–72 2276
- 2216 2277
- 2217 2278
- 2218 93. Mangasarian OL, Musicant DR (2000) Robust linear and support vector regression. *IEEE Trans Pattern Anal Mach Intell* 22:950–955 2279
- 2219 2280
- 2220 2281
- 2221 94. Mangasarian OL, Musicant DR (2002) Large scale kernel regression via linear programming. *Mach Learn* 46:255–269 2282
- 2222 2283
- 2223 95. Merkwirth C, Mauser HA, Schulz-Gasch T, Roche O, Stahl M, Lengauer T (2004) Ensemble methods for classification in cheminformatics. *J Chem Inf Comput Sci* 44:1971–1978 2284
- 2224 2285
- 2225 2286
- 2226 96. Mitchell TM (1997) Machine learning. McGraw-Hill, Maidenhead, p 432 2287
- 2227 2288
- 2228 97. Müller K-R, Rätsch G, Sonnenburg S, Mika S, Grimm M, Heinrich N (2005) Classifying ‘drug-likeness’ with kernel-based learning methods. *J Chem Inf Model* 45:249–253 2289
- 2229 2290
- 2230 2291
- 2231 98. Neugebauer A, Hartmann RW, Klein CD (2007) Prediction of protein-protein interaction inhibitors by cheminformatics and machine learning methods. *J Med Chem* 50:4665–4668 2292
- 2232 2293
- 2233 2294
- 2234 99. Neumann D, Kohlbacher O, Merkwirth C, Lengauer T (2006) A fully computational model for predicting percutaneous drug absorption. *J Chem Inf Model* 46:424–429 2295
- 2235 2296
- 2236 2297
- 2237 100. Nidhi **CE3**, Glick M, Davies JW, Jenkins JL (2006) Prediction of biological targets for compounds using multiple-category Bayesian models trained on chemogenomics databases. *J Chem Inf Model* 46:1124–1133 2298
- 2238 2299
101. Nigsch F, Bender A, van Buuren B, Tissen J, Nigsch E, Mitchell JBO (2006) Melting point prediction employing *k*-nearest neighbor algorithms and genetic parameter optimization. *J Chem Inf Model* 46:2412–2422 2244
102. Oloff S, Muegge I (2007) kScore: a novel machine learning approach that is not dependent on the data structure of the training set. *J Comput-Aided Mol Des* 21:87–95 2247
103. Oloff S, Zhang S, Sukumar N, Breneman C, Tropsha A (2006) Chemometric analysis of ligand receptor complementarity: Identifying complementary ligands based on receptor information (CoLiBRI). *J Chem Inf Model* 46:844–851 2251
104. Palmer DS, O’Boyle NM, Glen RC, Mitchell JBO (2007) Random forest models to predict aqueous solubility. *J Chem Inf Model* 47:150–158 2254
105. Pelletier DJ, Gehlhaar D, Tilloy-Ellul A, Johnson TO, Greene N (2007) Evaluation of a published in silico model and construction of a novel Bayesian model for predicting phospholipidosis inducing potential. *J Chem Inf Model* 47:1196–1205 2258
106. Platt J (1999) Fast training of support vector machines using sequential minimal optimization. In: Schölkopf B, Burges CJC, Smola AJ (eds) *Advances in kernel methods – support vector learning*. MIT Press, Cambridge, pp 185–208 2262
107. Plewczynski D, Spieser SAH, Koch U (2006) Assessing different classification methods for virtual screening. *J Chem Inf Model* 46:1098–1106 2265
108. Quinlan R (1993) *C4.5: programs for machine learning*. Morgan Kaufmann, San Mateo 2267
109. Ren S (2002) Classifying class I and class II compounds by hydrophobicity and hydrogen bonding descriptors. *Environ Toxicol* 17:415–423 2270
110. Ripley BD (2008) *Pattern recognition and neural networks*. Cambridge University Press, Cambridge, p 416 2272
111. Rodgers S, Glen RC, Bender A (2006) Characterizing bitterness: identification of key structural features and development of a classification model. *J Chem Inf Model* 46:569–576 2275
112. Rusinko A, Farmen MW, Lambert CG, Brown PL, Young SS (1999) Analysis of a large structure/biological activity data set using recursive partitioning. *J Chem Inf Comput Sci* 39:1017–1026 2278
113. Sakiyama Y, Yuki H, Moriya T, Hattori K, Suzuki M, Shimada K, Honma T (2008) Predicting human liver microsomal stability with machine learning techniques. *J Mol Graph Modell* 26:907–915 2282
114. Schneider N, Jäckels C, Andres C, Hutter MC (2008) Gradual in silico filtering for druglike substances. *J Chem Inf Model* 48:613–628 2286
115. Schölkopf B, Smola AJ (2002) *Learning with kernels*. MIT Press, Cambridge 2287
116. Schölkopf B, Sung KK, Burges CJC, Girosi F, Niyogi P, Poggio T, Vapnik V (1997) Comparing support vector machines with gaussian kernels to radial basis function classifiers. *IEEE Trans Signal Process* 45:2758–2765 2292
117. Schölkopf B, Burges CJC, Smola AJ (1999) *Advances in kernel methods: support vector learning*. MIT Press, Cambridge 2294
118. Schroeter TS, Schwaighofer A, Mika S, ter Laak A, Suelzle D, Ganzer U, Heinrich N, Müller K-R (2007) Estimating the domain of applicability for machine learning QSAR models: a study on aqueous solubility of drug discovery molecules. *J Comput-Aided Mol Des* 21:485–498 2298

**CE3** Please provide author initials.

- 2300 119. Shawe-Taylor J, Cristianini N (2004) Kernel methods for pat- 2361  
 2301 tern analysis. Cambridge University Press, Cambridge 2362
- 2302 120. Shen M, LeTiran A, Xiao Y, Golbraikh A, Kohn H, Tropsha A 2363  
 2303 (2002) Quantitative structure-activity relationship analysis 2364  
 2304 of functionalized amino acid anticonvulsant agents using 2365  
 2305 *k*-nearest neighbor and simulated annealing PLS methods. 2366  
 2306 *J Med Chem* 45:2811–2823 2367
- 2307 121. Shen M, Xiao Y, Golbraikh A, Gombar VK, Tropsha A (2003) De- 2368  
 2308 velopment and validation of *k*-nearest-neighbor QSPR mod- 2369  
 2309 els of metabolic stability of drug candidates. *J Med Chem* 2370  
 2310 46:3013–3020 2371
- 2311 122. Smola AJ, Schölkopf B (2004) A tutorial on support vector re- 2372  
 2312 gression. *Stat Comput* 14:199–222 2373
- 2313 123. Sommer S, Kramer S (2007) Three data mining techniques 2374  
 2314 to improve lazy structure-activity relationships for noncon- 2375  
 2315 generic compounds. *J Chem Inf Model* 47:2035–2043 2376
- 2316 124. Sorich MJ, McKinnon RA, Miners JO, Smith PA (2006) The im- 2377  
 2317 portance of local chemical structure for chemical metabolism 2378  
 2318 by human uridine 5'-diphosphate-glucuronosyltransferase. 2379  
 2319 *J Chem Inf Model* 46:2692–2697 2380
- 2320 125. Sun H (2005) A naive Bayes classifier for prediction of mul- 2381  
 2321 tidrug resistance reversal activity on the basis of atom typing. 2382  
 2322 *J Med Chem* 48:4031–4039 2383
- 2323 126. Suykens JAK (2001) Support vector machines: a nonlinear 2384  
 2324 modelling and control perspective. *Eur J Control* 7:311–327 2385
- 2325 127. Suykens JAK, Van Gestel T, De Brabanter J, De Moor B, Vande- 2386  
 2326 walle J (2002) Least squares support vector machines. World 2387  
 2327 Scientific, Singapore 2388
- 2328 128. Svetnik V, Liaw A, Tong C, Culbertson JC, Sheridan RP, Feuston 2389  
 2329 BP (2003) Random forest: a classification and regression tool 2390  
 2330 for compound classification and QSAR modeling. *J Chem Inf 2391*  
 2331 *Comput Sci* 43:1947–1958 2392
- 2332 129. Svetnik V, Wang T, Tong C, A Liaw, Sheridan RP, Song Q (2005) 2393  
 2333 Boosting: an ensemble learning tool for compound classifica- 2394  
 2334 tion and QSAR modeling. *J Chem Inf Model* 45:786–799 2395
- 2335 130. Swamidass SJ, Chen J, Phung P, Ralaivola L, Baldi P (2005) Ker- 2396  
 2336 nels for small molecules and the prediction of mutagenicity, 2397  
 2337 toxicity and anti-cancer activity. *Bioinformatics* 21[S1]:i359– 2398  
 2338 i368 2399
- 2339 131. Terfloth L, Bienfait B, Gasteiger J (2007) Ligand-based models 2400  
 2340 for the isoform specificity of cytochrome P450 3A4, 2D6, and 2401  
 2341 2C9 substrates. *J Chem Inf Model* 47:1688–1701 2402
- 2342 132. Tobita M, Nishikawa T, Nagashima R (2005) A discriminant 2403  
 2343 model constructed by the support vector machine method 2404  
 2344 for HERG potassium channel inhibitors. *Bioorg Med Chem 2405*  
 2345 *Lett* 15:2886–2890 2406
- 2346 133. Todeschini R, Consonni V (2003) Descriptors from molecular 2407  
 2347 geometry. In: Gasteiger J (ed) *Handbook of chemoinformatics*, 2408  
 2348 vol 3. Wiley, Weinheim, pp 1004–1033 2409
- 2349 134. Tong W, Hong H, Fang H, Xie Q, Perkins R (2003) Decision for- 2410  
 2350 est: Combining the predictions of multiple independent de- 2411  
 2351 cision tree models. *J Chem Inf Comput Sci* 43:525–531 2412
- 2352 135. Tong W, Xie Q, Hong H, Shi L, Fang H, Perkins R (2004) As- 2413  
 2353 sessment of prediction confidence and domain extrapola- 2414  
 2354 tion of two structure-activity relationship models for predict- 2415  
 2355 ing estrogen receptor binding activity. *Env Health Perspect 2416*  
 2356 112:1249–1254 2417
- 2357 136. Trinajstić N (1992) *Chemical graph theory*. CRC Press, Boca Ra- 2418  
 2358 ton 2419
- 2359 137. Urrestarazu Ramos E, Vaes WHJ, Verhaar HJM, Hermens JLM 2420  
 2360 (1998) Quantitative structure-activity relationships for the 2421  
 aquatic toxicity of polar and nonpolar narcotic pollutants. 2422  
*J Chem Inf Comput Sci* 38:845–852 2423
138. Vapnik VN (1979) Estimation of dependencies based on em- 2424  
 pirical data. Nauka, Moscow 2425
139. Vapnik VN (1995) *The nature of statistical learning theory*. 2426  
 Springer, New York 2427
140. Vapnik VN (1998) *Statistical learning theory*. Wiley, New York 2428
141. Vapnik VN (1999) An overview of statistical learning theory. 2429  
*IEEE Trans Neural Netw* 10:988–999 2430
142. Vapnik V, Chapelle O (2000) Bounds on error expectation for 2431  
 support vector machines. *Neural Comput* 12:2013–2036 2432
143. Vapnik VN, Chervonenkis AY (1974) *Theory of pattern recog- 2433*  
 nition. Nauka, Moscow 2434
144. Vapnik V, Lerner A (1963) Pattern recognition using general- 2435  
 ized portrait method. *Automat Remote Control* 24:774–780 2436
145. Varnek A, Kireeva N, Tetko IV, Baskin II, Solov'ev VP (2007) Ex- 2437  
 haustive QSPR studies of a large diverse set of ionic liquids: 2438  
 how accurately can we predict melting points? *J Chem Inf 2439*  
*Model* 47:1111–1122 2440
146. Vogt M, Bajorath J (2008) Bayesian similarity searching in 2441  
 high-dimensional descriptor spaces combined with Kull- 2442  
 back–Leibler descriptor divergence analysis. *J Chem Inf 2443*  
*Model* 48:247–255 2444
147. von Korff M, Sander T (2006) Toxicity-indicating structural 2445  
 patterns. *J Chem Inf Model* 46:536–544 2446
148. Votano JR, Parham M, Hall LM, Hall LH, Kier LB, Oloff S, Trop- 2447  
 sha A (2006) QSAR modeling of human serum protein bind- 2448  
 ing with several modeling techniques utilizing structure-in- 2449  
 formation representation. *J Med Chem* 49:7169–7181 2450
149. Wang J, Du H, Yao X, Hu Z (2007) Using classification structure 2451  
 pharmacokinetic relationship (SCPR) method to predict drug 2452  
 bioavailability based on grid-search support vector machine. 2453  
*Anal Chim Acta* 601:156–163 2454
150. Watson P (2008) Naïve Bayes classification using 2D pharma- 2455  
 cophore feature triplet vectors. *J Chem Inf Model* 48:166–178 2456
151. Witten IH, Frank E (2005) *Data mining: practical machine 2457*  
 learning tools and techniques, 2nd edn. Morgan Kaufmann, 2458  
 San Francisco, p 525 2459
152. Xiao Z, Xiao Y-D, Feng J, Golbraikh A, Tropsha A, Lee K-H 2460  
 (2002) Antitumor agents. 213. Modeling of epipodophyllo- 2461  
 toxin derivatives using variable selection *k*-nearest neighbor 2462  
 QSAR method. *J Med Chem* 45:2294–2309 2463
153. Xue Y, Li ZR, Yap CW, Sun LZ, Chen X, Chen YZ (2004) Ef- 2464  
 fect of molecular descriptor feature selection in support vec- 2465  
 tor machine classification of pharmacokinetic and toxicolog- 2466  
 ical properties of chemical agents. *J Chem Inf Comput Sci 2467*  
 44:1630–1638 2468
154. Yamashita F, Hara H, Ito T, Hashida M (2008) Novel hierarchi- 2469  
 cal classification and visualization method for multiobjective 2470  
 optimization of drug properties: application to structure-acti- 2471  
 vity relationship analysis of cytochrome P450 metabolism. 2472  
*J Chem Inf Model* 48:364–369 2473
155. Yap CW, Chen YZ (2005) Prediction of cytochrome P450 3A4, 2474  
 2D6, and 2C9 inhibitors and substrates by using support vec- 2475  
 tor machines. *J Chem Inf Model* 45:982–992 2476
156. Yap CW, Cai CZ, Xue Y, Chen YZ (2004) Prediction of torsade- 2477  
 causing potential of drugs by support vector machine ap- 2478  
 proach. *Toxicol Sci* 79:170–177 2479
157. Yu G-X, Park B-H, Chandramohan P, Munavalli R, Geist A, Sam- 2480  
 atova NF (2005) In silico discovery of enzyme-substrate speci- 2481  
 ficity-determining residue clusters. *J Mol Biol* 352:1105–1117 2482

- 2422 158. Yue P, Li Z, Moulton J (2005) Loss of protein structure stability  
2423 as a major causative factor in monogenic disease. *J Mol Biol*  
2424 353:459–473
- 2425 159. Zhang S, Golbraikh A, Oloff S, Kohn H, Tropsha A (2006)  
2426 A novel automated lazy learning QSAR (ALL-QSAR) approach:  
2427 method development, applications, and virtual screening  
2428 of chemical databases using validated ALL-QSAR models.  
2429 *J Chem Inf Model* 46:1984–1995
- 2430 160. Zhang S, Golbraikh A, Tropsha A (2006) Development of  
2431 quantitative structure-binding affinity relationship models  
2432 based on novel geometrical chemical descriptors of the pro-  
2433 tein-ligand interfaces. *J Med Chem* 49:2713–2724
- 2434 161. Zheng WF, Tropsha A (2000) Novel variable selection quanti-  
2435 tative structure-property relationship approach based on the  
2436 *k*-nearest-neighbor principle. *J Chem Inf Comput Sci* 40:185–  
2437 194

Uncorrected Proof  
2008-09-04