

Toward an Optimal Procedure for Variable Selection and QSAR Model Building

A. Yasri* and D. Hartsough

Computational Design Group, ArQule Inc., 19 Presidential Way, Woburn Massachusetts 01801

Received March 12, 2001

In this work, we report the development of a novel QSAR technique combining genetic algorithms and neural networks for selecting a subset of relevant descriptors and building the optimal neural network architecture for QSAR studies. This technique uses a neural network to map the dependent property of interest with the descriptors preselected by the genetic algorithm. This technique differs from other variable selection techniques combining genetic algorithms to neural networks by two main features: (1) The variable selection search performed by the genetic algorithm is not constrained to a defined number of descriptors. (2) The optimal neural network architecture is explored in parallel with the variable selection by dynamically modifying the size of the hidden layer. By using both artificial data and real biological data, we show that this technique can be used to build both classification and regression models and outperforms simpler variable selection techniques mainly for nonlinear data sets. The results obtained on real data are compared to previous work using other modeling techniques. We also discuss some important issues in building QSAR models and good practices for QSAR studies.

1. INTRODUCTION

At present, one cannot talk about drug design without mentioning quantitative structure–activity relationships (QSAR). The pioneering work introducing the QSAR concept was the performed by Hansch et al. (1962).¹ This work demonstrated that the biological activity of a molecule can be quantitatively linked to some physicochemical parameters (i.e. molecular descriptors), and it also introduced the idea that the activity can be described by more than one parameter (i.e. multiple regression). These structural descriptors also provided useful information for quantitative structure–property relationships (QSPR).

The QSAR approach attempts to find consistent relationships between the variations in the values of molecular properties and the biological activity for a series of compounds, so that these “rules” can be used to evaluate new chemical entities. Like other data mining approaches, QSAR is performed in successive steps including data preparation, data reduction, and data modeling and prediction. The success of a QSAR study relies heavily on how each of these steps is conducted and how the analysis is performed.

Building a QSAR model for a set of compounds begins by collecting and organizing experimental and theoretical data for these compounds. The experimental information may concern biological properties (activity, toxicity or bio-availability, metabolism), which correspond to the dependent variables in the data mining process. The other data are physicochemical calculated descriptors including parameters to account for hydrophobicity,^{2,3} topology,⁴ electronic,⁵ and steric⁶ effects which have been determined empirically or, more recently, by computational methods. These structural attributes represent the independent variables.

The use of numerous descriptors that are indicative of molecular structure and topology is becoming more common

in QSAR.⁷ These types of descriptors are easily calculated from molecular structures and potentially number in the thousands. Most QSAR methods deal with descriptors of molecular structure derived from either a two-dimensional (2D) representation of molecular structure based on molecular connectivity⁸ or three-dimensional (3D) structural information such as shape descriptors (volume, surface area, moment of inertia, ...), steric or electrostatic fields.⁹ The clear physicochemical sense of steric and electrostatic descriptors have made the CoMFA approach one of the most popular methods for QSAR/QSPR. However, several problems have persisted in the use of 3D approaches: (1) In the CoMFA approach, the initial compound alignment remains a difficult task and dramatically affects the results. (2) Any 3D descriptor depends on the conformation used to represent the molecule. Because it is difficult to define the conformation of a molecule responsible for its activity, this type of descriptor will introduce some noise during the modeling step and consequently affect the QSAR model. (3) The use of QSAR models including 3D descriptors for predicting large sets of compounds such as a virtual or real combinatorial library is very costly. On the other hand, computing 2D descriptors is very fast and does not suffer from any alignment or conformational problem. Thus, such descriptors can be easily used for representing a large set of compounds from combinatorial libraries. For these reasons, we have considered the use of various 2D descriptors for developing QSAR.

The data preparation step ends by organizing the data in a spreadsheet where the molecules and the descriptors correspond to the rows and the columns of data, respectively. This data representation is almost acceptable to techniques of predictive data mining. The second step is to reduce the data by selecting pertinent descriptors from a large set that faithfully describes the activity of interest. Choosing the adequate descriptors for QSAR studies is difficult because

* Corresponding author phone: (781)994-0547; fax: (781)994-0679; e-mail: ayasri@arqule.com.

there are no absolute rules that govern this choice. To deal with this issue, variable selection techniques were introduced where optimization search algorithms such as stepwise forward, stepwise backward, and simulated annealing are used. These search methods were combined to scoring functions such as multiple linear regression (MLR),¹⁰ linear or quadratic discriminant analyses (LDA, QDA),¹¹ and partial least-squares (PLS) analysis.^{12–14} More recently, evolutionary algorithms and specifically genetic algorithms were used for the variable selection problem when combined to PLS, MLR, k-nearest neighbors (kNN),¹⁵ and artificial neural networks (ANNs).^{16,45}

It has been shown that these variable selection techniques improved the QSAR models compared to those without variable selection. However, there are some recurrent problems in all of these methods: (1) The increase in the number of descriptors with regard to the number of compounds results in some chance correlation between the descriptors,¹⁷ and these correlations are sometimes treated intuitively by removing the unfavorable descriptors. (2) Linear statistical techniques such as MLR depend on an assumed linear relationship between the dependent variable and one or more descriptors. (3) In some regression techniques (MLR, PLS, PCR), there is an a priori assumption of the model form (quadratic, cubic, cross-terms, ...). (4) When nonlinear techniques such as ANNs are used as a scoring function, the topology of the network remains fixed during the optimization and no exploration of the network topology is performed. It is known that the performance of ANNs is fundamentally tied to their topology. Furthermore, an ANN with fixed size may lead to the problem of under/overfitting. In this paper, we choose to use genetic algorithms to optimize both the variable selection and the model form (neural networks architecture). The derived models are compared with those obtained (1) from much simpler variable selection techniques such as stepwise forward selection (FS) combined to LDA or MLR and (2) from models without variable selection step.

The last steps in a QSAR study are the data modeling and prediction. Generally, when a variable selection procedure is performed, the data modeling step is a part of the variable selection. Many modeling and prediction techniques are available and accessible; most can be readily applied to QSAR data. These techniques can be grouped into three categories: (i) linear and nonlinear mathematics; (ii) distance-based techniques; and (iii) logic techniques. The application of these techniques for QSAR model searching is performed on a training set, a subset representing either the whole or a part of the data (about 70–80%). The quality of the model is estimated by how well it performs the mapping between the descriptors and the activity in the training set. This mapping is expressed by the correlation coefficient R^2 or the root-mean-square error (RMSE) between the experimental activity and the predicted one

$$R^2 = \frac{\sum(\hat{Y}_i - \bar{Y})^2}{\sum(Y_i - \bar{Y})^2} \quad \text{RMSE} = \sqrt{\frac{\sum(\hat{Y}_i - \bar{Y})^2}{N}} \quad (1)$$

where Y_i and \hat{Y}_i are the observed and predicted activities for compound i , and \bar{Y} is the mean activity over the N compounds in the training set.

The validation of the model involves internal and/or external validation. Internal validation is performed by cross-validation techniques and/or a randomization test.¹⁸ Cross-validation consists of the following: (a) removing one (leave-one-out) or groups (leave-group-out) of compounds in a systematic or random way, (b) generating a model from the remaining compounds, and (c) predicting the removed compounds. The quality of the model is expressed by the cross-validated correlation coefficient q^2 :

$$q^2 = 1 - \frac{\sum(\hat{Y}_i - Y_i)^2}{\sum(Y_i - \bar{Y})^2} \quad (2)$$

The randomization technique consists of giving random values to the dependent variable and constructing a model with the real input descriptors. This randomization is repeated several times, and the resulting fit scores are compared to the fit with the real nonrandomized model. External validation is the most important in measuring the robustness of a model. It consists of making predictions for an independent set of compounds (about 20–30% of the whole data) not used in model training. In this work, we used the internal leave-group-out cross-validation technique during the variable selection step, and the randomization technique as well as test set prediction to test the validity of our QSAR models.

Overall, in this paper we report the development and validation of a variable selection and model building process using a careful analysis during all the QSAR model development phases. This approach combines genetic algorithms with artificial neural networks to preselect descriptors. It also identifies the optimal network topology to map the activity of interest to the selected descriptors. It differs from other variable selection techniques combining genetic algorithms to neural networks by two main features: (1) the variable selection search procedure is not constrained to a defined number of descriptors and (2) we perform an exploration of the optimal neural network architecture. Background on genetic algorithms and artificial neural networks is given below.

1.1. Background. Genetic Algorithms. Genetic algorithms (GAs) are perhaps the best known of all evolution-based search algorithms.¹⁹ GAs were developed by John Holland in an attempt to explain the adaptive processes of natural systems and to design artificial systems based upon these natural systems.^{20,21} Having become widely used for a broad range of optimization problems in the last 15 years,²² the GA has been described as being a search algorithm with some of the innovative flair of human search.²³

Traditionally, GAs have used binary strings to encode the features that compose an individual in the population; the binary segments of an individual that represent a specific feature are known as chromosomes. Binary strings are convenient to use because they can be easily manipulated by GA operators such as crossover and mutation. Binary chromosomes can also be used to represent nonbinary numbers that have integer and floating point types. Given a problem and a population of individuals, a GA will evaluate each individual as a potential solution according to a predefined evaluation function (i.e. scoring function). The evaluation function assigns a value of goodness to each

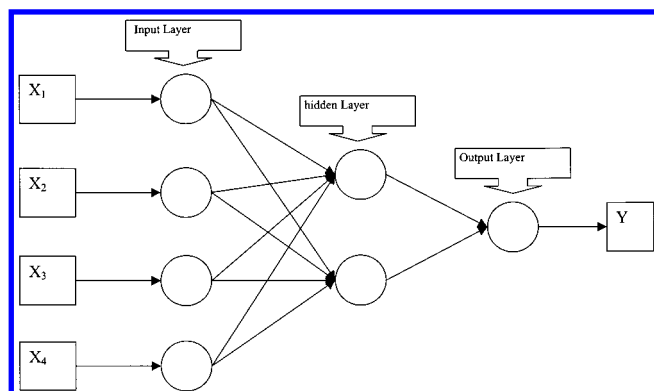


Figure 1. Schematic representation of a fully connected feedforward neural network with three layers. Each node in the input layer receive signal from an input dependent variable X . The hidden nodes role is to intervene between the external input (X) and the network output (Y).

individual based on how well the individual solves a given problem. This metric is then used by a fitness function to determine which individuals will breed to produce the next generation. Breeding is done by crossing existing solutions with each other in order to produce new solutions. In addition, a mutation factor is present which will randomly modify existing solutions. Mutations helps the GA break out of local minima. While there is no guarantee that GAs will find an optimal solution, their method of selection and breeding candidate solutions can converge to good solutions provided that enough generations are given.

GAs have been used for variable and descriptor selection in QSAR studies.^{12,16} However, these GAs were constrained by a fixed number of variables represented in each chromosome. In the present work, we employed an unconstrained GA, and the number of descriptors used is searched during the optimization process.

Neural Networks. The neural network approach to data analysis has received much attention of late. Neural networks have overcome the theoretical limitations of *perceptrons* and early linear networks by the introduction of “hidden layers” to represent intermediate processing and to compute non-linear recognition functions.²⁴ The rapid advancement of computing systems in the past decade has also contributed significantly to the success of this approach in various engineering, business, and medical applications. Neural network systems learn to discriminate among classes of patterns within an input domain in a holistic manner. They are presented with training sets of representative instances of each class, correctly classified, and they learn to recognize and predict other new instances of these classes. Learning consists of adjusting weights in a fixed-topology network via different learning algorithms.

Back-propagation networks, commonly implemented as layered, fully connected feedforward networks (Figure 1), have been used quite successfully in various engineering and business applications. Recently they have also been adopted in biomedical areas such as analysis of appendicitis and cancer patient data;²⁵ cancer image extraction and classification;²⁶ studies of soybean diseases;²⁷ and in pharmaceutical applications such as pharmaceutical production development,²⁸ pharmacodynamic modeling,²⁹ and pharmacological effects of drug concentrations.³⁰

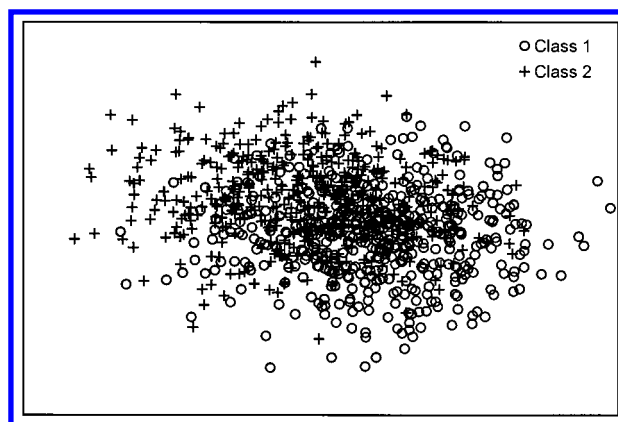


Figure 2. Distribution of the simulated classification data in the plan defined by the two informative variables. The two classes show overlapping at the middle of the plan and separation at each side of the plan.

2. MATERIALS AND METHODS

2.1. Simulated Data. Data Set 1 (Classification Problem 1: Sim_linear). For this data set, we used 1000 samples, 98 random variables from standard normal distributions (mean = 0; stdev = 1), and two informative variables dividing the samples into two classes with roughly a 25% overlap (Figure 2A). In this way, we know that the two last variables are the informative features, and the rest are noisy variables. The samples are distributed uniformly within the two classes (50–50%). The data are divided into a training set (700 samples) and a prediction set (300 samples).

Data Set 2 (Regression Problem 1, Sim_quadratic). Ninety-seven random distributions were generated for 1000 samples and three quadratic informative variables (V_1, V_2, V_3). The output dependent variable Y is generated according to the following expression with some noise represented by ϵ :

$$Y = V_1^2 + V_2^2 + V_3^2 + \epsilon \quad (3)$$

The data are divided into a training set (700 samples), a validation set (150 samples), and a prediction set (150 samples).

Data Set 3 (Regression Problem 2, Sim_sinusoidal). Ninety-seven distributions were randomly generated for 1000 samples and three sinusoidal informative variables (V_1, V_2, V_3). The output dependent variable Y is generated according to the following expression:

$$Y = \sin(V_1) + \sin(V_2) + \sin(V_3) + \epsilon \quad (4)$$

The data are divided into a training set (700 samples), a validation set (150 samples), and a prediction set (300 samples).

2.2. Experimental Data. Data Set 4 (Regression Problem 1: Benzodiazepines Data). Benzodiazepines are well-known as anxiolytics, tranquilizers, and anticonvulsants in epilepsy treatment. For this data set, we used 55 compounds whose biological activity (IC_{50}) was reported in the work of Haefely et al. (1985).³¹ We randomly divided the data into three subsets: a training set of 40 compounds, a validation set of 6 compounds, and a prediction set of 8 compounds.

Data Set 5: (Regression Problem 2: Tubulin Polymerization Inhibitors). Data on inhibition of tubulin polymerization in vitro were taken from Weigt and co-workers.^{32–35}

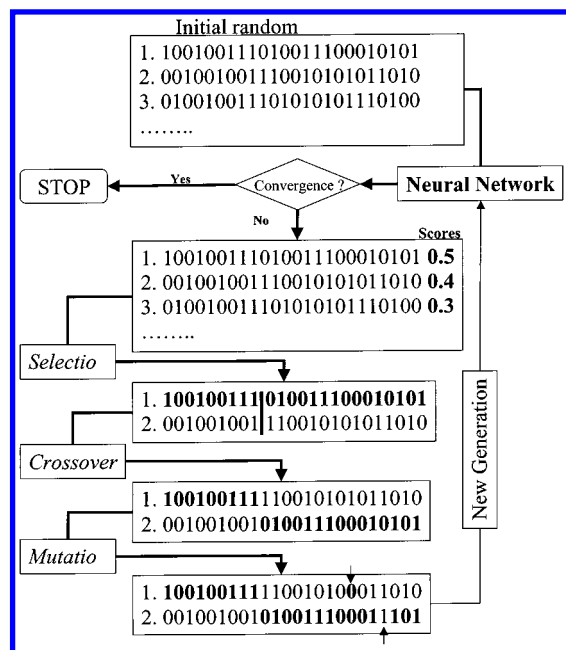


Figure 3. Flow diagram describing the steps used in selecting the best subset of descriptors and building a neural network model by ARQeDES program.

These data was divided into two groups: a training set of 40 compounds and a validation set of eight compounds. A set of 10 compounds, whose IC₅₀ values are more than 40 μ M, was left as a prediction set. This prediction set was also used in the work of Weigt and Wiese to validate the derived CoMFA models.

Data Set 6: (Classification Problem 1: CNS versus Non CNS Drugs). A set of 268 compounds with known CNS activity were collected from the recent work of Ajay et al. (1999).³⁶ Many of these compounds are reported in the works of Waterbeemd,³⁷ Fisher,³⁸ and Norinder.³⁹ Here, we adopted the classification used by Ajay and co-workers using as the CNS-active class compounds with BBB ≥ 1.0 and as CNS-inactive to those with BBB < 1.0 ; here BB = $C_{\text{brain}}/C_{\text{blood}}$. The compounds were randomly divided into three subsets: training, validation, and test sets containing 130, 38, and 107 compounds, respectively.

2.3. Computational Details. All calculations were done on R10000 SGI Iris workstations. For structural descriptor generation, we used Molconnz,⁴⁰ DiverseSolutions,⁴¹ and MOE,⁴² software with structure input in MDL SD file format. The program combining genetic algorithms and neural networks called ARQeDES was implemented in C code and executed on an SGI workstation. This program uses the PGAPack library for the genetic algorithm⁴³ and the SNNS package for the artificial neural networks.⁴⁴

2.4. Genetic Algorithms. Figure 3 shows a schematic representation of the GA optimization strategy. The first step in the GA is to create a gene pool of individuals (initial random population). Each individual or chromosome encodes a different subset of descriptors by a binary string representation. The value of one implies that the corresponding descriptor is included in the subset, and zero means that the descriptor is excluded. The length of each chromosome is the same and is equal to the total number of descriptors. Using each parent combination of descriptors associated with the training set, a neural network (NN) is trained and cross-

validated; thus for each chromosome a value corresponding to the fit is returned as a fitness value. Then, two parents are selected randomly based on the roulette wheel selection method. The population of individuals (descriptor combinations) is evolved by performing a crossover between two selected parents who produce two offspring. Each offspring is subjected to a random single-point mutation, i.e., a randomly selected one (or zero) is changed to zero (or one). The fitness of each offspring is evaluated by training of a neural network model. When the resulting offspring are characterized by a better value of the fitness function, they replace the parents; otherwise, the parents are kept. This process is repeated until a predefined convergence criterion is achieved.

In summary, each parent in our optimization problem represents a randomly preselected descriptor subset, and the objective of the optimization is to evolve the initial population of descriptor subsets to the population with the best fit of the training data using the neural network.

2.4. Neural Networks. During the GA optimization procedure, NNs are used as the evaluation function for mapping the molecular descriptors to the activity of interest (dependent variable) by using the Quickprop learning function. Quickprop is a back-propagation technique which speeds up the learning by using information about the curvature of the error surface.³³ This requires the computation of the second derivatives of the error function. Quickprop assumes the error surface to be locally quadratic and attempts to move in one step from the current position directly to the minimum of the parabola.

For each descriptor subset, the best topology of the neural network is searched by using the training set and the validation set. The validation set is used to monitor over-training of the network. Once the best topology of the network is found and the convergence criteria is reached, the model is cross-validated with the leave-group-out technique. This procedure returns a fitting score to the GA expressed in the following equation

$$\text{Score} = \text{Fit}_T - 0.4 \cdot |\text{Fit}_T - \text{Fit}_{CV}| \quad (5)$$

where Fit_T is the fit result on the training set, and Fit_{CV} is the cross-validated fit. Here, the fit is expressed by the correlation coefficient R^2 between experimental and predicted dependent variable. The 0.4 weight factor of the equation's second term was used because it has been found to produce the best behavior of the GA in previous studies.^{45,46}

3. RESULTS

To demonstrate the effectiveness of the method, two type of data were used: (1) three sets of simulated data concerning regression and classification problems with known solutions and (2) three sets of real data covering different biological domains. The results obtained for these data sets will be discussed in relation to the selected descriptors, the statistical significance of the models judged by q^2 , and the results for the prediction sets.

For all of the optimization procedures performed, the GA used 200 members as a population size (i.e. 200 variable combinations) and evolved for 300 generations. Since the GA optimization procedure uses the elitism concept (the best combinations are kept in the following generation), only 10%

of the whole combinations are evaluated by the NN in each generation. All of the NN models were cross-validated using the leave-group-out technique and internally validated by a randomization test. The randomization test consists of randomizing the values of the dependent variable (the activity) N times, and the resulting data are trained against real independent variables (the input descriptors). In this work, the values of the dependent variable were randomized 100 times, and the q^2 of each resulting model was compared to the real model. The rationale behind this test is that the significance of the real QSAR model decreases if there is a significant chance correlation between the selected descriptors and the randomized response variable.¹⁶

3.1. Simulated Data. Artificial data are very useful to validate variable selection techniques since we know the form of the response surface, the informative variables fitting the data set, and the noisy variables. In this work we studied artificial classification and regression data where only a few variables are necessary to fit the dependent variable.

In the first simulated data set *sim_linear*, 1000 samples are distributed equally between two classes and described by 100 variables. The samples within the two classes can be classified at 75% by only two variables, and the remaining 98 variables are independent of the classes. This data set was split into a training set (700 samples), a validation set (150 samples), and a prediction set (150 samples).

Using different starting random seeds for the genetic algorithm, three different runs of ARQeDES variable selection and model building were performed. To validate the results of ARQeDES variable selection, a classification model was built with only the two informative variables and compared to those obtained from the variable selection procedures. All of the neural network models required only one hidden neuron to perform classification of the training set. The cross-validated results and classification predictions for the training set and prediction set are summarized in Table 1. All three trials of the variable selection procedure resulted in a subset of variables including the two informative variables and some noisy variables. The neural network models obtained with these descriptors performed classification at $76.89 \pm 1.93\%$ for the training set and $66.78 \pm 2.23\%$ for the prediction set. This classification result is similar to that obtained with the model without noisy variables (74.13% and 65.90 for training and prediction sets, respectively) and with the model from FS-LDA which selects the two informative and seven noisy variables (73.15% and 67% for training and prediction sets, respectively). The good classification performance of the neural networks with true and noisy variables illustrates (1) that the variable selection by ARQeDES was able to catch the relevant variables for the classification problem and (2) how neural networks are robust in the sense that they respond with the right output even when presented with input patterns containing noise by attributing small weights to links coming from these noisy inputs.

The second artificial data set is for a quadratic regression problem (*sim_quadratic*). It was shown previously that neural networks are able to discover the quadratic nature of analyzed data even without square terms.^{47,48} As with the first data set, a neural network model with only the informative variables was built and compared to the model

Table 1. Classification Results of the Training and Prediction Sets by ARQeDES Variable Selection^c

		training			prediction			CPU time ^a
		class 1	class 2	%WC	class 1	class 2	%WC	
2V _{inf}	class 1	266	88	75.14	52	24	68.42	<i>b</i>
	class 2	93	253	73.12	27	47	63.51	
	mean			74.13			65.90	
2V _{inf} and 6V _{nois}	class 1	281	78	78.27	48	28	63.15	2.5
	class 2	79	262	76.83	23	51	68.91	
	mean			77.55			66.03	
2V _{inf} and 7V _{nois}	class 1	291	79	78.64	49	27	64.47	2.8
	class 2	72	258	78.18	25	49	66.21	
	mean			78.41			65.34	
2V _{inf} and 5V _{nois}	class 1	261	93	73.72	54	22	71.05	2.6
	class 2	84	262	75.72	24	50	67.56	
	mean			74.72			69.30	
FS-LDA	class 1	253	93	73.12	49	25	66	0.33
	class 2	91	263	74.29	24	52	68	
	mean			73.15			67	

^a Hours of CPU time for the variable selection procedures. ^b The CPU time for the model with only the informative variables is not provided since it is negligible. ^c Comparison with the model using only informative variables for normal distribution data set (*sim_lin*) and FS-LDA model. V_{inf}: informative variable; V_{nois}: noisy variable; WC: well classified; FS-LDA: stepwise forward selection combined to linear discriminant analysis.

Table 2. Comparison of Neural Network Models by ARQeDES with the Informative and Noisy Variables for the Quadratic, Sinusoidal Response Functions and the FS-MLR Procedure^c

	variables	q^2 ^a	R^2 ^b	CPU time ^c
quadratic	3V _{inf}	0.88	0.91	<i>d</i>
	3V _{inf} and 1V _{nois}	0.92	0.91	4.2
FS-MLR	1V _{inf} and 2V _{nois}	-0.025	0.007	0.25
sinusoidal	3V _{inf}	0.74	0.68	<i>d</i>
	3V _{inf} and 1V _{nois}	0.75	0.65	3.8
FS-MLR	3V _{inf} and 3V _{nois}	0.45	0.154	0.3

^a Cross-validated q^2 for the training set. ^b Correlation coefficient between the actual and predicted values for the test set. ^c Hours of CPU time for the variable selection procedures. ^d The CPU time for the model with only the informative variables is not provided since it is negligible. ^e V_{inf}: informative variable; V_{nois}: noisy variable; FS-MLR: stepwise forward selection combined to multiple linear regression.

obtained with variables selected by ARQeDES and by FS-MLR. The results are summarized in Table 2.

The variable selection search resulted in three variables corresponding to the informative variables and one noisy variable. Using these independent variables a NN model was trained with a cross-validated q^2 of 0.88 for the training set and a correlation coefficient R^2 of 0.91 between the actual and predicted values for the test set (Figure 4). The NN model with only the three informative variables showed quite similar q^2 and R^2 (0.92 and 0.91, respectively). This demonstrates the ability of the ARQeDES variable selection procedure to extract the pertinent variables for a quantitative dependent variable and to build a quantitative model for a quadratic response surface. As for the variable selection by FS-MLR, it resulted in one informative and two noisy variables giving very poor predictive performances for both the training ($q^2 = -0.025$) and the test sets ($R^2 = 0.007$). For this nonlinear data the FS-MLR variable selection was significantly outperformed by the GA-NN procedure.

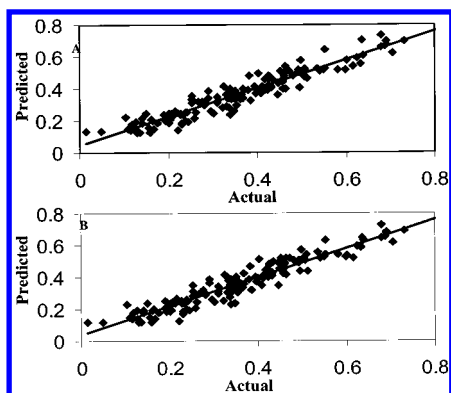


Figure 4. Actual versus predicted values for the test set obtained by the NN model with only three informative variables (A) and with ARQeDES selected variables (B).

The last simulated problem also uses a quantitative dependent variable and is a function of three sinusoidal distributions. The NN model built with only the three true independent variables resulted in a q^2 for the training set of 0.74 and an R^2 for the test set of 0.68. The variable selection by ARQeDES resulted in the three informative variables and one noisy variable being selected out of the pool of 100 variables and a NN model with 0.75 and 0.65 for q^2 and R^2 , respectively (Table 2). The variable selection by FS-MLR also selected the three informative variables and three noisy variables, but the derived linear model presented poor predictive performances for the training data ($q^2 = 0.45$) and test data ($R^2 = 0.15$). Again, the variable selection by ARQeDES identified the informative variables and was able to derive a NN model with good prediction power for the test set and outperformed the model derived from FS-MLR procedure.

As for the required time for the GA-NN procedure, the variable selection and model searching steps do not take more than 4 CPU hours at the average (Tables 1 and 2). The simpler variable selection techniques take much less time, but the derived models from these selection procedures present very poor predictive performances.

3.2. Experimental Biological Data. Three experimental data sets were collected from the literature and used in this work. The results of variable selection and model building are compared to prior QSAR models obtained for the same data set. The models obtained in our work are used to predict a set of compounds not used during the training phase.

3.2.1. Benzodiazepine Data. The first data set is a collection of 54 benzodiazepine derivatives with antipentylenetetrazole activity (data set 4, methodology section). The activity was expressed as $-\log(\text{IC}_{50})$ and scaled between 0.1 and 0.9. A training set of 40 compounds and a validation set of six compounds were represented by 404 descriptors submitted to ARQeDES to select a pertinent subset of descriptors and to build a NN regression model. Two types of optimization searches were performed: (1) four unconstrained procedures where the number of optimal descriptors and the topology of the NN are searched and (2) two constrained procedures where the number of descriptors is fixed to 6 and the topology of the NN is fixed to 6-2-1 and 6-4-1. In Table 3, we summarize the results for the selected variables and models obtained for each search procedure as well as the model derived from the FS-MLR selection.

Table 3. ARQeDES Results for Variable Selection and Model Building for Benzodiazepine Data Set^f

model name	NN topology	q^2	R_p^2
BZ_dynsc_1	8-2-1	0.85	0.58
BZ_dyn_2	6-2-1	0.90	0.57
BZ_dyn_3	6-2-1	0.81	0.58
BZ_dyn_4	8-2-1	0.88	0.54
BZ_fix_1	6-2-1	0.88	0.58
BZ_fix_2	6-4-1	0.84	0.60
FS-MLR	7-0-1	0.80	0.59
NN-0 ^a	404-2-1	0.9	0.1
Maddalena_NN ^b	10-3-1	0.89	
So_NN ^c	6-2-1	0.93	
Winkler_NN ^d	11-6-1	0.82	
Bono_MLR ^e		0.86-0.88	

^a Neural network model obtained without any variable selection procedure. ^b Results by back-propagation NN reported in ref 49. ^c Results by back-propagation NN reported in ref 16. ^d Results by back-propagation NN reported in ref 55. ^e Results by multilinear regression (MLR) reported in ref 50. ^f Comparison with other NN and MLR models reported in previous works. q^2 : cross-validated correlation coefficient for the training set; R_p^2 : correlation coefficient between experimental and predicted IC_{50} values for the test set.

The unconstrained procedures resulted in six to eight selected descriptors and NN topologies with only two hidden neurons (BZ_dyn1-4). These models all show good cross-validated correlation coefficients (q^2) ranging from 0.81 to 0.9 which are comparable to the NN models reported by Maddalena and Johnston⁴⁹ and So and Karplus,¹⁶ both using back-propagation. These models are also in good agreement with the model recently reported by Bono and Nedad using multilinear regression (MLR) analysis.⁵⁰ However, the previously reported works on benzodiazepine derivatives data set presented linear or nonlinear models which were trained on the whole data set, and no testing was presented. One can argue that the cross-validated q^2 is a good indicator of the predictivity of a model, but since the model is not tested externally, it is not obvious that it can truly generalize.

The constrained model BZ_fix_2 whose NN topology is 6-2-1 reproduces the results from the nonconstrained procedure. When the NN topology is increased to four hidden neurons in model BZ_fix_4, the predictivity of the model is not significantly increased. This suggests that a simple NN topology is sufficient to fit the data and avoid the risk of overfitting.

The FS-MLR variable selection resulted in seven descriptors and a linear model showing similar prediction performance to that of the GA-NN model for the training and test data (Table 3). When no descriptor selection is performed and all of the 404 descriptors are used, the obtained NN model presented overtraining and very poor prediction performance (Table 3).

Figure 5 shows the randomization results expressed as R^2 for the training against q^2 for one of the nonconstrained models (BZ_dyn_3). One can see that there is a clear separation between the real model and the randomized ones suggesting that the correlations found are significant.

3.2.2. Tubulin Polymerization Data. The tubulin polymerization activity of 48 compounds was contained in this data set.³²⁻³⁵ The activity data were expressed as $-\log(\text{IC}_{50})$ and scaled between 0.1 and 0.9. For this data set, constrained and nonconstrained variable selection and model building were performed. This resulted in a good nonconstrained

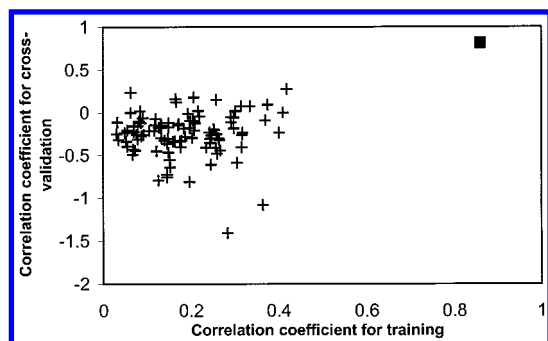


Figure 5. Scatter plot for R^2 (correlation coefficient for training) versus q^2 (cross-validated correlation coefficient) for the real model (plain square) and randomized models (star).

Table 4. Comparison for Prediction by NN Model Obtained by ARQeDES and CoMFA Model from Weigt and Weise Work³²

compound name	NN ARQeDES model	CoMFA model
q^2 ^a	0.80	0.87
19	21.45	9
39	42.09	25
42	40.45	32
43	41.89	29
51	36.20	12
54	35.06	14
55	35.83	25
56	23.84	17
59	23.31	38
60	23.74	33

^a q^2 : cross-validated correlation coefficient for the training set.

solution with seven descriptors giving a NN model with 7–2–1 architecture and a cross-validated q^2 for the training set of 0.75. The constrained search procedures with 4–1–1, 4–2–1, 6–1–1, 6–2–1, 7–1–1, 7–2–1, 8–1–1, and 8–2–1 architectures resulted in values of q^2 ranging from 0.34 to 0.62, except for 7–2–1 which gave a q^2 of 0.78 as the nonconstrained solution. To validate the best model, we predicted 10 compounds which were not seen by the model and whose IC₅₀ values are more than 40 μ M, and we compared the results to predictions by the work of Weigt and Wiese using a CoMFA model (Table 4).

As reported by Weigt and Weise, the prediction of these 10 compounds by the CoMFA model is not obvious since it requires an extrapolation of activities out of the activity range of the training set. In our NN model, the activity of the training set was scaled in a way to have a minimum IC₅₀ value of 0.05 and a maximum of 150 μ M. All of the 10 compounds are predicted with low activity; namely compounds 39, 42, and 43 whose predicted IC₅₀ values are more than 40 μ M. The other compounds have predicted IC₅₀ values ranging from 21 to 36 μ M. Similarly, the predicted activities for the 10 compounds from Weigt and Weise are also very low, but the large predicted values are below 40 μ M.

3.3.3. CNS versus Non CNS Drugs. Three compound subsets were used to build and validate a classification model discriminating between drugs with or without CNS activity: (1) a training set of 130 compounds with 58 CNS active and 72 CNS inactive compounds, (2) a validation set of 38 compounds with 21 CNS active and 17 CNS inactive compounds to monitor neural network over-training, and (3) a test set of 100 compounds to validate the final model.

ARQeDES variable selection and modeling building resulted in a subset of six topological descriptors and a NN

model having 6–3–2 as the network topology. This model achieved 100% classification accuracy for CNS active and 93% for CNS inactive compounds in the training set. In the test set, this model accurately predicted 91.25% of CNS active and 70% of CNS inactive compounds. This prediction accuracy is in good agreement with prediction on the same compounds in the work of Ajay and co-workers (1999).³⁶

The variable selection procedure by FS-LDA procedure for the CNS data set resulted in 14 descriptors and presented a classification accuracy of 88.10% for the training set (96.5% for GA-NN) and 16% for the test set (80.62% for GA-NN). In these data, the GA-NN procedure of ARQeDES significantly outperformed the FS-LDA selection.

4. DISCUSSION

This work employed a genetic algorithm combined to an artificial neural network to find a set of descriptors and to build a QSAR model for biological activity. The most important issue dealt with is how each step for QSAR model building is performed from data collection to model validation and prediction. These QSAR model steps were illustrated by using both artificial and real biological data.

The first question stressed in this work is the importance of the molecular descriptors for QSAR studies. The three-dimensional descriptors such as shape descriptors present the disadvantage of depending on the conformation of the compounds. When these types of descriptors are used, we must make the assumption that the conformation used to derive the descriptor is responsible for the biological activity of the compound. This assumption remains true in the case of rigid compounds. However, when there are a large number of rotatable bonds within the compound, one conformation is inadequate for representing the entire conformational space covered by the compound, and consequently the descriptor may be useless and introduce noise. Conversely, 2D information derived from the molecular graph of the molecule does not suffer any conformational problem and can be used to describe both rigid and flexible compounds. In fact, we have shown that 2D descriptors are sufficient to capture pertinent information for the compounds under study. We applied the 2D description to data sets concerning three biological activities: benzodiazepines,³¹ tubuline polymerization inhibitors,³² and CNS drugs.³⁶ In each case, the selected descriptors describe essential interactions such as steric, electronic, hydrogen bonding, and hydrophobic parameters. They also served to build QSAR models able to predict biological activity of new compounds.

The most important and limiting step in any QSAR study is variable selection because the method of choosing descriptors can have a great influence on all subsequent steps in drug design from QSAR model building to use of models in screening new compounds. Ideally, the best solution could be obtained by an exhaustive search of all possible combinations of the initial descriptors, which is possible with small data set or with a small number of descriptors. But in most cases, and because of numerous and diverse descriptors used in QSAR studies, this search is prohibitively time-consuming. In this paper, we used genetic algorithms, a stochastic optimization technique which is the method of choice for searching spaces of high dimension with multiple optima such as a QSAR function linking the chemical space (set of descriptors) to the biological space (biological activity).

In this work, we showed that constrained search by fixing the number of desired descriptors influences the solution of the genetic algorithm and may result in a nonoptimal solution. On the other hand, the unconstrained search can better explore the solution space and results in a subset of descriptors that faithfully fit the training set and allow construction of a predictive NN model.

The scoring function evaluating the solutions proposed by the genetic algorithm must allow the comparison of solutions by guaranteeing that attributed scores do not depend solely on either the number of descriptors or on the function expression. To this purpose, artificial neural networks were used to fit the preselected descriptors to the activity of interest. The way we used the neural networks is original in terms of exploring the model space. In fact, for each preselected descriptor subset, increasing the size of hidden layer step-by-step, the optimal number of hidden nodes is searched. For each hidden neuron added to the network, the fit is calculated, and the prediction error on the validation set is determined. The addition of hidden neurons is stopped when the neural network prediction error on the validation set starts increasing, even if the error on the training set continues to decrease. This approach allows us to use the minimum number of hidden neurons to fit the training set and minimize the risk of over-fitting. Once the optimal neural network model is set, a cross-validation procedure is performed in order to test the stability of the model, and a combined score from the training set and cross-validation is returned to the genetic algorithm.

Further internal validation is performed on each final model by performing a randomization test.¹⁶ We showed that all the descriptors selected by the genetic algorithm combined with the neural networks implemented in ARQeDES, produced models that presented a significant separation from models generated by randomized activities. This provides compelling evidence that ARQeDES's variable selection and model building procedure is most likely a result of genuine correlation between the chosen descriptors and activity.

The good performance of our variable selection and NN model construction is supported by the results from the simulated data. In fact, artificial data are very important to validate any data mining technique because we know the correct solution to the problem, and because we can draw a clear conclusion about the limits of the technique. In our case, we used different simulated data for quantitative and qualitative dependent variables having different response surfaces (linear, quadratic, and sinusoidal). For the classification as well as for the quantitation problem, the variable selection technique of ARQeDES was able to capture the pertinent and informative variables from a large set of noisy variables. Some noisy variables comprised part of the solution provided by the technique, but they did not adversely affect the predictivity of the neural network model. Indeed, Figure 6 shows the attributed normalized weights to the selected variables for simulated data. All of the weights assigned to the noisy variables are very small compared to those attributed to the informative variables. This finding suggests that the way the neural network is trained is optimal in terms of discriminating between true and noisy signals by the attribution of the appropriate weight to each input node.

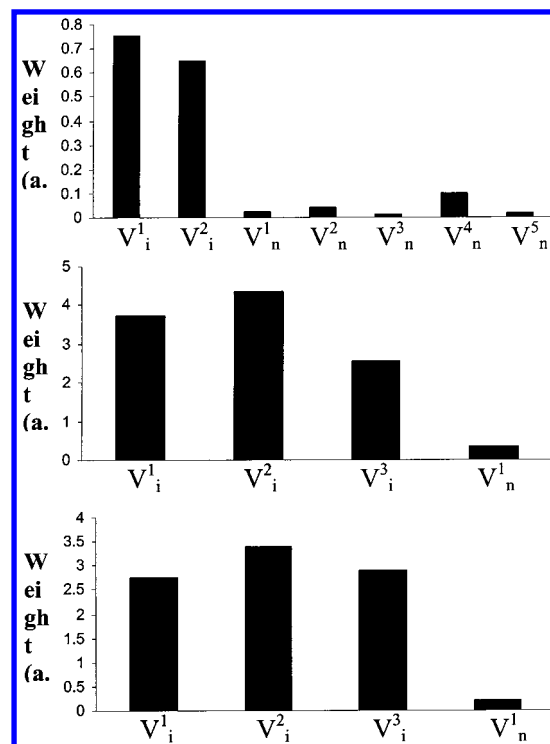


Figure 6. Normalized weights of the selected variables for linear classification and regression artificial data: V_i (Informative variable) and V_n (noisy variable).

Real biological data such as the anticonvulsant activity of benzodiazepines has been widely used to validate regression and variable selection techniques because a wealth of chemical and pharmacological information is available for this chemical series.^{49,51–53} All of the previous studies using a variety of multivariate statistical techniques, linear or nonlinear, reported very good models for benzodiazepines activity. A few of these studies reported real predictions for the activity of these molecules;⁵⁴ the others claimed that their QSAR models exhibited considerable predictive power based on the values of cross-validated q^2 (0.7–0.9).^{7,16,48,50,55,56} q^2 criteria is a good indicator of the predictivity of the model on the training set, and since no real prediction on a test set is performed, there is no guarantee that these models can generalize. In this work, even if the number of compounds in the benzodiazepines data is not large, we kept a test set to validate our NN models derived from the variable selection procedure. Furthermore, any new data mining technique presenting good performances on benzodiazepines data must be validated on other real data to prove its power and to show the limits. The application of the variable selection and neural network model building technique implemented in ARQeDES was successful in modeling tubulin polymerization inhibition data and in the discrimination between CNS active and CNS inactive drugs.

5. CONCLUSIONS

In this paper, we have reported the development of a novel variable selection and neural network model building technique called ARQeDES. The method was successfully applied to both artificial and real biological data sets. The effectiveness of the method is demonstrated by the selection of the pertinent and informative variables for each problem treated. This selection is shown to be optimal to derive a

predictive neural network model. The method combines genetic algorithms, a stochastic optimization method, to artificial neural networks. The genetic algorithm is used to preselect descriptor subsets, and the neural network is used to map the input variables (descriptors) to the dependent variable (activity).

The good performance of the technique on simulated data allowed us to conclude that (1) the variable selection by ARQeDES is able to capture the relevant independent variables explaining the dependent one (activity of interest), and (2) the derived neural network models are robust in the sense that they respond with the right output even when presented with input patterns containing some noise; this is achieved by attributing "unfavorable" weights to links coming from these noisy inputs.

By using real biological data, we also showed that nonconstrained variable selection procedures, where the number of descriptors and the network architecture are not fixed but searched during the optimization, provide better solutions to the problem than the constrained ones. These neural network models are derived from a careful analysis involving a deep exploration of the model space.

The molecular descriptors used in this work are only 2D descriptors, and no 3D information was involved in modeling the data. These 2D descriptors were applied to data sets concerning three biological activities: benzodiazepines, tubulin polymerization inhibitors, and CNS drugs. In each case, the selected descriptors describe essential interactions such as steric, electronic, hydrogen bonding, and hydrophobic parameters. They also served to build QSAR models able to predict the biological activity of new compounds.

The vast amount of information about structure–activity relationships produced by combinatorial chemistry and high-throughput screening in pharmaceutical companies for lead identification and optimization has created a great need for fast, accurate, and fully automated data mining tools. We have shown that the ARQeDES method presented in this work is an elaborated QSAR technique which can be successfully used to optimize the set of descriptors for neural networks trained on biological data. This technique could be employed to rapidly model the activity for a large set of compounds such as primary high-throughput screening results and to build nonlinear classification models able to discriminate between active and nonactive compounds. In the lead optimization domain, the technique can be used to optimize the set of descriptors for neural networks trained on quantitative activity for the chemical series containing the lead compound and then to derive a regression model which can be used to optimize the activity of that lead.

ACKNOWLEDGMENT

The authors would like to thank Dr. Barry Wythoff and Dr. Ryszard Czerminski from ArQule Inc. for their fruitful discussion, critical reading, and support of this work.

REFERENCES AND NOTES

- (1) Hansch, C.; Muir, R. M.; Fujita, T.; Maloney, P. P.; Geiger, F.; Streich, M. The correlation of biological activity of plant growth regulators and chloromycetin derivatives with Hammett constants and partition coefficients. *J. Am. Chem. Soc.* **1963**, *85*, 2817–2824.
- (2) Hansh C.; Leo A. In *Exploring QSAR: Fundamentals and applications in chemistry and biology*; Heller S. R., Eds; American Chemical Society: Washington, DC, 1995.
- (3) Rekker, R. F. The Hydrophobic Fragmental Constants. Its Derivation and Application. A Means of Characterizing Membrane Systems. In *Pharmacochemistry Library*; Nauta, W. T., Rekker, R. F., Eds.; Elsevier: New York, 1977; Vol. 1.
- (4) Hall, L. H.; Kier, L. B. The molecular connectivity Chi Indexes and Kappa Shape Indexes in structure–Property Modeling. In *Reviews in computational chemistry II*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH Publishers: 1991; pp 367–422.
- (5) Exner, O. In *Advances in Free Energy Relationships*; Chapman, N. B., Shorter, J., Eds.; Plenum: New York, 1972; p 1.
- (6) Verloop, A.; Hoogenstraaten, W.; Tipker, J. In *Drug Design*; Ariens, E. J., Ed.; Academic Press: New York, 1976; Vol. VII, p 165.
- (7) Waller, C. L.; Bradley, M. P. Development and Validation of a novel variable selection technique with application to multidimensional quantitative structure–activity relationship studies. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 345–355.
- (8) Hall, L. H.; Kier, L. B. The molecular connectivity Chi Indexes and Kappa Shape Indexes in Structure–Property Modeling. In *Review in Computational Chemistry II*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH Publishers: 1991; pp 367–422.
- (9) Cramer, R. D.; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect of shape on binding of steroids to carrier proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (10) Shi, L. M.; Fan, Y.; Myers, T. G.; O'Connor, P. M.; Paull, K. D.; Friend, S. H.; Weinstein, J. N. Mining the NCI anticancer drug discovery databases: genetic function approximation for the QSAR study of anticancer ellipticine analogues. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 2 189–99.
- (11) Partek Pro 2000. Partek Incorporated, Partek Analysis and Recognition technologies, Copyright 1993–1999.
- (12) Hasegawa, K.; Miyashita, Y.; Funatsu, K. GA strategy for variable selection in QSAR studies: GA-based PLS analysis of calcium channel antagonists. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 2, 306–310.
- (13) Hasegawa, K.; Miyashita, Y.; Funatsu, K. GA strategy for variable selection in QSAR studies: GA-PLS and D-optimal Designs for predictive QSAR model. *J. Mol. Struct. (THEOCHEM)* **1998**, *425*, 255–262.
- (14) Cho SJ, Zheng W, Tropsha A. Rational combinatorial library design. 2. Rational design of targeted combinatorial peptide libraries using chemical similarity probe and the inverse QSAR approaches. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 2, 259–268.
- (15) Zheng W.; Tropsha, A. Novel Variable Selection Quantitative Structure–Property relationships approach based on the k-nearest neighbor principle. *J. Inf. Comput. Sci.* **2000**, *40*, 185–194.
- (16) So, S. S.; Karplus, M. Genetic Neural Networks for quantitative structure–activity relationships: improvement and application of Benzodiazepine affinity for Benzodiazepine/GABA A receptor. *J. Med. Chem.* **1996**, *39*, 5246–5256.
- (17) Topliss, J. G.; Edwards, R. P. Chance factors in Studies of quantitative Structure–Activity Relationships. *J. Med. Chem.* **1979**, *22*, 1238–1244.
- (18) Van de Waterbeemd, H. *Chemometrics Methods in Molecular Design; Methods and principles in medicinal chemistry*; VCH: Weinheim, 1995; Vol. 2.
- (19) Bentley, P. J. *Evolutionary Design by computers*; Morgan K. Publishers, Inc.: San Francisco, CA, 1999.
- (20) Holland J. H. Genetic algorithms and Optimal allocations of trails. *SIAM J. Computing* **1973**, *2*, 2, 88–105.
- (21) Holland J. H. *Adaptation in Natural and artificial systems*; University of Michigan Press: Ann Arbor, 1975.
- (22) Holland J. H. *Genetic algorithms*; Scientific American: **1992**, 66–72.
- (23) Goldberg, D. E. *Genetic algorithms in search, optimization and machine learning*; Addison-Wesley: 1989.
- (24) Rumelhart, D. E.; Hinton G. E.; Williams R. J. Learning internal representations by error propagation. In *Parallel Distributed Processing*; Rumelhart, D. E., McClelland, J. L., and the PDP Research Group, Eds.; The MIT Press: Cambridge, MA, 1986; pp 318–362.
- (25) Weiss, S. M.; Kapouleas, I. An empirical comparison of pattern recognition, neural nets, and machine learning classification methods. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)*; Detroit, MI, 1989; pp 781–787.
- (26) Moallemi, C. Classifying cells for cancer diagnosis using neural networks. *IEEE EXPERT*; 1991; pp 8–12.
- (27) Mooney, R.; Shavlik, J.; Towell, G.; Gove, A. An experimental comparison of symbolic and connectionist learning algorithms. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)*; Detroit, MI, 1989; pp 775–780.
- (28) Hussain, A. S.; Yu, X.; Johnson R. D. Application of neural network computing in pharmaceutical product development. *Pharm. Res.* **1991**, *8*, 1248–1252.

- (29) Veng-Pedersen, P.; Modi, N. B. Neural networks in pharmacodynamic modeling. Is current modeling practice of complex kinetic systems at a dead end? *J. Pharmacokinet. Biopharm.* **1992**, *20*, 397–412.
- (30) Erb, R. J. Neural network computation of nonlinear pharmacological effects from concentration of drug in the central compartment independent of equilibrium status. *Pharm. Res.* **1992**, *9*, 293.
- (31) Haefely, W.; Kyburz, E.; Gegecke, M.; Möhler, H. Recent advances in the molecular pharmacology of benzodiazepine receptors and the structure–activity relationships of their agonist and antagonists. *Adv. Drug Res.* **1985**, *14*, 165–322.
- (32) Weigt, M.; Weise, M. A comparative Molecular Field Analysis of Inhibitors of Tubulin Polymerization. *QSAR J.* **19(2)**, **2000**, 142–148.
- (33) Fahlman, S. E. In *Faster-learning variations on back-propagation: An empirical study*; Sejnowski, T. J., Hinton, G. E., Touretzky, D. S., Eds.; Connectionist Model Summer School, Morgan Kaufmann: San Mateo, CA, 1988.
- (34) Li, L.; Wang, H. K.; Kuo, S. C.; Wu, T. S.; Lednicer, D.; Lin, C. M.; Hamel, E.; Lee, K. H. Antitumor agents 155. Synthesis and biological evaluation of 3',6,7-substituted 2-phenyl-4-quinolones as antimicrotubule agents. *J. Med. Chem.* **1994**, *37*, 3400–3407.
- (35) Chen, K.; Kuo, S. C.; Hsieh, M. C.; Mauger, A.; Lin, C. M.; Hamel, E.; Lee, K.-H. Antitumor agents 174. 2',3',4',5,6,7-substituted 2-phenyl-1,8-naphthyridin-4-ones: their synthesis, cytotoxicity, and inhibition of tubulin polymerization. *J. Med. Chem.* **1997**, *40*, 2266–2275.
- (36) Ajay, B. W. B.; Murcko, M. A. Designing libraries with CNS activity. *J. Med. Chem.* **1999**, *42*, 24 4942–51.
- (37) van de Waterbeemd, H.; Camenisch, G.; Folkers, G.; Chretien, J. R.; Raevsky, O. A. Estimation of blood-brain barrier crossing of drugs using molecular size and shape, and H-bonding descriptors. *J. Drug. Target.* **1998**, *6*, 2, 151–65.
- (38) Fisher, H.; Gottschlich, R.; Seeling, A. Blood-brain barrier permeation: molecular parameters governing passive diffusion. *J. Membr. Biol.* **1998**, *165*, 3, 201–11.
- (39) Norinder, U.; Sjöberg, P.; Osterberg, T. Theoretical calculation and prediction of brain-blood partitioning of organic solutes using MolSurf parametrization and PLS statistics. *J. Pharm. Sci.* **1998**, *87*, 8, 952–959.
- (40) Molconn-Z, Molconn software 3.5; Lowell H. Hall, copyright 1998.
- (41) DiverseSolutions v 3.0.2, Distributed by Tripos, Inc. on behalf of the Laboratory of Molecular Graphics and theoretical Modeling, College of Pharmacy, University of Texas at Austin; Austin, TX, 1997.
- (42) MOE 2000.02. Molecular Operating Environment, Chemical Computing Group Inc.: Copyright 1997–2000.
- (43) PGAPack Parallel, Genetic Algorithm library, Nevine D., Mathematics and Computer Science Division at Argonne National Laboratory: Argonne, IL.
- (44) Malmgren, H. For classification and predictive purposes, simulated neural networks (SNNs; more often called artificial neural networks, ANNs) offer a powerful alternative to traditional statistical analyses [letter]. *Epilepsia* 1999, *40*, 9, 1323–4.
- (45) Wessel, M. D.; Jurs, P. C.; Tolan, J. W.; Muskal, S. M. Prediction of Human Intestinal Absorption of Drugs Compounds from Molecular Structure. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 726–735.
- (46) Bakken, G. A.; Jurs, P. C. Prediction of Hydroxyl radical rate constants from molecular structure. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 1064–1075.
- (47) Tetko, I. V.; Livingstone, D. J.; Luil, A. I. Neural Network studies. 1. Comparison of overfitting and overtraining. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 826–833.
- (48) Tetko, I. V.; Vila, A. E. P.; Livingstone, D. J.; Luil, A. I. Neural Network studies. 2. Variable Selection. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 794–803.
- (49) Maddalena, D. J.; Johnston, G. A. Prediction of receptor properties and binding affinity of ligands to benzodiazepine/GABA_A receptors using artificial neural networks. *J. Med. Chem.* **1995**, *38*, 4, 715–724.
- (50) Bono, L.; Nenad, T. Multivariate Regression Outperforms Several Robust Architectures of Neural Networks in QSAR Modeling. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 121–132.
- (51) Loew, G. H.; Nienow, J. R.; Poulsen, M. Theoretical structure–activity studies of benzodiazepine analogues. *Mol. Pharmacol.* **1984**, *26*, 19–34.
- (52) Greco, G.; Novellino, E.; Silipo, C.; Vittoria, A. Study of benzodiazepines receptor sites using a combined QSAR–CoMFA approach. *Quant. Struct.-Act. Relat.* **1992**, *11*, 461–477.
- (53) Ghose, A. K.; Crippen, G. M. Modeling the benzodiazepine receptor binding site by the general three-dimensional structure-directed quantitative structure–activity relationship method REMOTEDISC. *Mol. Pharmacol.* **1990**, *37*, 725–734.
- (54) Winkler, D. A.; Burden, F. R. Holographic QSAR of Benzodiazepines. *QSAR J.* **1998**, *17*, 224–231.
- (55) Winkler, D. A.; Burden, F. R.; Watkins, A. J. R. Atomistic Topological Indices Applied to Benzodiazepines using Various Regression Methods. *QSAR J.* **1998**, *17*, 14–19.
- (56) Kovalishyn, V. V.; Tetko, I. V.; Luik, A. I.; Kholodovych, V. V.; Villa, A. E. P.; Livingstone, D. J. Neural network studies. 3. Variable Selection in the cascade-correlation learning architecture. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 651–659.

CI010291A