# Heuristic Extraction of Rules in Pruned Artificial Neural Networks Models Used for Quantifying Highly Overlapping Chromatographic Peaks

César Hervás

Department of Computer Science, Edificio C2, Rabanales Campus, University of Cordoba,
E-14071 Cordoba, Spain

Manuel Silva,* Juan Manuel Serrano, and Eva Orejuela

Department of Analytical Chemistry, Edificio C3-Anexo, Rabanales Campus, University of Cordoba,
E-14071 Cordoba, Spain

The suitability of an approach for extracting heuristic rules from trained artificial neural networks (ANNs) pruned by a regularization method and with architectures designed by evolutionary computation for quantifying highly overlapping chromatographic peaks is demonstrated. The ANN input data are estimated by the Levenberg−Marquardt method in the form of a four-parameter Weibull curve associated with the profile of the chromatographic band. To test this approach, two *N*-methylcarbamate pesticides, carbofuran and propoxur, were quantified using a classic peroxyoxalate chemiluminescence reaction as a detection system for chromatographic analysis. Straightforward network topologies (one and two outputs models) allow the analytes to be quantified in concentration ratios ranging from 1:7 to 5:1 with an average standard error of prediction for the generalization test of 2.7 and 2.3% for carbofuran and propoxur, respectively. The reduced dimensions of the selected ANN architectures, especially those obtained after using heuristic rules, allowed simple quantification equations to be developed that transform the input variables into output variables. These equations can be easily interpreted from a chemical point of view to attain quantitative analytical information regarding the effect of both analytes on the characteristics of chromatographic bands, namely profile, dispersion, peak height, and residence time.

## INTRODUCTION

Chemometrics is a discipline of chemistry concerned with the application of mathematical methods to design or select optimal measurement procedures and experiments and to provide maximum chemical information by analyzing chemical data.[1] Resolution is one of the principal aims of all chromatographic techniques. Thus, a major goal of chromatographic method development is to obtain the adequate separation of all components of interest within a reasonable elution window. One of the most powerful choices is liquid chromatography (LC) coupled with some sort of multichannel detection method such as a diode array detector (DAD). This hyphenated technique generates data matrices where the columns correspond to chromatograms and the rows to spectra. Despite its high separation efficiency, peaks are often poorly resolved, and measurements must be made using data analysis methods.[2,3] A variety of chemometric techniques have been used for this purpose considering that the data matrix obtained from LC/DAD is bilinear. As a result, iterative target transformation factor analysis,[4] evolving factor analysis,[4−6] window factor analysis,[6] generalized rank annihilation method,[7,8] and heuristic evolving latent projections,[4,9] among others, have been developed and applied to quantify the components that provide overlapping chromatographic bands. The majority of these methods do not work

well for highly overlapping peaks, and, also, they are normally applied to cases where peak shapes are symmetrical, which is often not the case in many real-world analyses. Moreover, pure reference standards are commonly required, and some hypothesis about the data must be assumed. Recently, the use of ANNs for the resolution of overlapping chromatographic analytical signals has gained popularity, although few publications have dealt with this subject.[10−12] ANNs provide better results than ordinary least-squares (OLS) and partial least-squares (PLS) methods by using a three-dimensional data matrix for each sample obtained by the LC/DAD technique.[10] Despite these efforts, the problem remains to be completely solved, and further research is therefore required, especially that focused on methods for resolving highly overlapping chromatographic bands which are not restricted to two-dimensional arrays of data, such as those obtained from LC/DAD.

One of the most outstanding features of ANNs, in contrast to other chemometric tools, is that "no prior knowledge" is needed about the relation among variables. ANNs, however, are considered "black boxes" as the models obtained are difficult to interpret. Nevertheless, this viewpoint has changed in recent years owing to a recent trend in this field to develop ANNs of minimal size to solve real-world problems, i.e. networks having a small number of weights, but possessing a reasonably high generalization capacity. Moreover, a small-sized neural network will be less prone to overtraining noise or the structure of the data in the training set, thus increasing

* Corresponding author phone: +34-957-212099; e-mail: qa1sirom@uco.es.

PRUNED ARTIFICIAL NEURAL NETWORKS MODELS

*J. Chem. Inf. Comput. Sci., Vol. 44, No. 5, 2004* **1577**

its generalization capacity over a new data set and allowing users to gain knowledge from the trained neural network in order to achieve a better understanding of how the network solves the problem. In general, most of the studies related to the extraction of symbolic rules have been aimed at solving classification problems,[13] and there are few methods related to extracting rules from trained neural networks for regression purposes.[14,15] One of these is the artificial neural network-decision trees (ANN-DT) method.[15] This algorithm is capable of extracting rules from function approximating networks in such a way that decision trees are constructed based on the network inputs and its respective outputs without analyzing the activation values for the hidden units or the connection weights of the networks. Given that this treatment precludes interpretability due to the cause–effect relationship between inputs and outputs, one of the most important aims of this work is to determine the explicative ability of the proposed models. Rules for function approximation normally take the following form: if (a condition or restriction in the input variables, x, is satisfied), then the output predicts y = f(x), where f(x) is a constant or a linear or nonlinear simple function of x. This kind of rule is acceptable if we take into account its similarity with nonlinear classification and statistical regression methods. A method for rule extraction from function approximating neural networks algorithm has been recently reported.[16]

In this work we use pruned ANN models for quantifying analytes that provide highly overlapping chromatographic peaks obtained from a single detector instrument, such as the chemiluminescence (CL) detector. The William regularization method[17] was used for the pruned process and a genetic algorithm,[18] in which the connection weights are codified by real numbers, to design the ANN architecture. Heuristic rules were used to simplify the best ANN models obtained by removing the sigmoidal functions with a negligible value with respect to the output value (ca. 1–5%). In this way, it is possible to remove several addends in certain subsets of the input space, thereby improving model interpretability. The proposed methodology was validated by determining two *N*-methylcarbamate pesticides, carbofuran and propoxur, whose chromatographic peaks exhibit a high degree of overlapping according to the composition of the mobile phase used for their separation by LC. The single analytical response provided by the chemiluminescence detector was fitted to a four-parameter Weibull function by using least-squares regression, and the estimates thus obtained were used as inputs to the ANNs. Further details on the analytical methodology are described elsewhere.[19]

## THEORETICAL BACKGROUND

The proposed approach for the quantification of overlapping chromatographic peaks is based on a two-step procedure to construct ANN models for predicting the contribution of each component to the overall chromatographic signal. The first step of this approach consists of extracting the information from the analytical response (chromatographic peak) in order to select the inputs to the ANNs. Upon examining this response, it can be observed that the chromatographic peaks $(t_i, S_{t_i})$ can be accurately fitted by least-squares regression to a four-parameter Weibull curve defined by $S_m$ (peak height), $t_m$ (residence time), $B$ (dispersion of the analytical signal

values from $S_m$), and $C$ (related to the function profile, which is associated to the inflection points of the curve), when the time domain ranges from $t_i > t_m - B((C-1)/C)^{(1/C)}$. If $t_i$ is standardized by subtracting $t_m$, dividing by the dispersion parameter $B$, and displacing the standardized variable $((C-1)/C)^{(1/C)}$ units, a new temporal variable $t_i'$ can be defined as follows

$$t_i' = \frac{t_i - t_m}{B} + \left(\frac{C-1}{C}\right)^{1/C} \tag{1}$$

which is characterized by the location parameter $((C-1)/C)^{(1/C)}$ and a dispersion parameter equal to one. In this case, $S_{t_i}$ is the response variable (analytical signal), which is proportional to the contribution (concentration) of each component to the chromatographic peak, and $t_i'$ is the independent variable. If it is assumed that the change of $S_{t_i}$ with time is proportional to the inverse of the transformation time, $t_i'$ and the parameter, $C$, associated with the convexity of the function, the following differential equation can be obtained:

$$\frac{\partial S_{t_i}}{\partial(t_i')} = \left(\frac{C-1}{t_i'} - C\right) S_{t_i} \tag{2}$$

Considering that $S_{t_i}$ at time $t_i' = 0$ is given by $S_m \times C \times \exp(-(C-1/C))$, where $S_m = S_{t_m}$, the integration of this equation provides

$$S_{t_i} = S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C} \tag{3}$$

for $t_i' \geq 0$, that is $t_i > t_m - B((C-1)/C)^{(1/C)}$, and $C > 1$ which corresponds to a four-parameter Weibull function. If additive errors $(\epsilon_i)$ are assumed, the nonlinear model is given by

$$S_{t_i} = S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C} + \epsilon_i \tag{4}$$

The least-squares principle was used to estimate the parameters in this nonlinear model. The least-squared estimates of $S_m$, $B$, $C$, and $t_m$, labeled as $\hat{S}_m$, $\hat{B}$, $\hat{C}$, and $\hat{t}_m$, are those that minimize the sum of squared residuals:

$$SS(\text{Re}s) = \sum_{i=1}^{n} [S_{t_i} - S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C}]^2 \tag{5}$$

To obtain these minimums, the partial derivatives of $SS(Res)$ with respect to each parameter were calculated and subsequently set at zero to obtain the four normal equations:

$$\frac{\partial SS(\text{Re}s)}{\partial S_m} = \sum_{i=1}^{n} (S_{t_i} - S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C}) \times$$
$$(Ce^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C}) = 0 \tag{6}$$

$$\frac{\partial SS(\text{Re}s)}{\partial B} = \sum_{i=1}^{n} (S_{t_i} - S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C}) \times$$
$$\frac{\partial (S_m C e^{-(C-1/C)} (t_i')^{C-1} e^{-t_i' C})}{\partial B} = 0 \tag{7}$$

$$\frac{\partial SS(\mathrm{Re}s)}{\partial C} = \sum_{i=1}^{n}(S_{t_i} - S_m Ce^{-(C-1/C)}(t_i')^{C-1}e^{-t_i'C}) \times$$

$$\frac{\partial(S_m Ce^{-(C-1/C)}(t_i')^{C-1}e^{-t_i'C})}{\partial C} = 0 \quad (8)$$

$$\frac{\partial SS(\mathrm{Re}s)}{\partial t_m} = \sum_{i=1}^{n}(S_{t_i} - S_m Ce^{-(C-1/C)}(t_i')^{C-1}e^{-t_i'C}) \times$$

$$\frac{\partial(S_m Ce^{-(C-1/C)}(t_i')^{C-1}e^{-t_i'C})}{\partial t_m} = 0 \quad (9)$$

There are no explicit solutions for these equations, thus calling for an iterative numerical procedure. We used a Levenberg−Marquardt method,[20] which tends toward the Gauss−Newton fitting[21] if the residual sum of the squares is reduced at each step or toward the steepest descent fitting[22] if the residual sum of squares increases in any step. The estimates obtained, $\hat{S}_m$, $\hat{B}$, $\hat{C}$, and $\hat{t}_m$, were used as inputs to the ANN. Hence, all the neural network models proposed in this work had four variables in the input layer.

In the second step, a procedure was used to construct various regularization-pruning ANN models in conjunction with genetic algorithms based on the learning of the network from the pattern of the training set. Genetic algorithms have been successfully used to optimize network topologies as well as to compute connection weights while avoiding minimum locals derived from overtraining.[23−25] Reproduction, crossover, and mutation operators are essential to the evolutionary process because they provide both selective pressure and diversity in such a way that the searching algorithm which seeks the best network presents an equilibrium between the exploitation of the best solutions and the exploration of other parts of the search space, avoiding local optima problems. In this work we use a new version of a real coded genetic algorithm (RCGA) which we recently developed[18,26,27] to optimize the network topology and a backpropagation learning procedure EDBD[28] for the learning process.

Binary representations for individuals (ANN models in this work) have been traditionally used in genetic algorithms, which evenly discretize the real domain of the ANN weights. The main problem of binary representation is that if the number of variables is large, the length of the chromosome is too long. A solution to this problem is the real codification of the chromosomes, where each gene is a floating-point number that represents a variable of the function to be optimized. In real coded genetic algorithms (RCGAs) a chromosome is coded as a finite-length string of the real numbers corresponding to the independent variables or design variables. The floating-point representation is robust, accurate, and efficient because it is closest to the real domain space, and moreover, the string length reduces to the number of independent variables or the number of genes, and these RCGAs have outperformed binary-representations in function optimization. Several RCGAs for nonlinear function optimization of continuous variables, using floating-point representation, have been early reported.[29−32]

In the RCGA used in this work, based on the one proposed by Beasley,[33,34] reproduction and crossover are applied in

**Table 1.** Parametric Values Used by the Algorithm

| William pruning | genetic operators | fitness function | genetic algorithm convergence |
|---|---|---|---|
| $ghiw = 1.1$ | $P_{cross} = 0.8$ | $\lambda_{0GA} = 0.75$ | $\gamma = 0.9$ |
| $ghider = 1.1$ | $P_{mut} = 0.01$ | $\beta_{GA} = 2$ | |
| | $P_{fro} = 0.9$ | $tolerance = 0.05$ | |

each generation to a small part of the population (typically just two individuals). Initially, a random population of 100 neural networks with an identical number of neurons in one hidden layer and a different connectivity (different weights) are created following a Laplace distribution.[17] A percentage of the connections are frozen, with a probability $P_{fro}$, so as to start the process with nets which are not fully connected, thus assisting the pruning process. The crossover operator removes the two worst individuals in the population and replaces them with two other individuals obtained by crossing two of the N-2 remaining ones, which are randomly selected. The crossing process consists of exchanging the input weights for a node in the hidden layer. The biological foundation behind this method is that species with a long life (i.e. parents and children) coexist for a certain period of time so that the latter can learn from the former, while at the same time there is competition among them. The mutation operator randomly selects a neuron in the hidden layer and then adds to each of its weights a random value (following a uniform distribution in $[-1,1]$) multiplied by a parameter called the mutation range $R_{mut}$. Formally

$$w^*_{jl} = w_{jl} + \xi \times R_{mut} \qquad \xi \in U[-1, 1] \qquad (10)$$

where $w_{jl}$ is the weight associated to the connection between the $l$ neuron of the hidden layer and the $j$ input variable. The mutation probability of each individual is defined by the parameter $P_{mut}$. This search method ameliorates the changes between generations. Therefore the number of generations increases until the stop criterion is met (maxgen). Nevertheless, the method yields excellent results for the generalization error.

The fitness function used in the genetic algorithm has two objectives: first, to minimize the residual sum of squares, and second, to reduce the unnecessary weights. The latter, weighted by a $\lambda$ parameter, penalizes more complex models (models with a larger number of weights) as compared to less complex models, albeit with the same generalization capacity. Thus, models with a smaller number of weights to be estimated that maintain the mean square error of the generalization set are favored. This is known as weight decay. Some techniques known as pruning reduce the network size by modifying not only the weights but also the network structure during training, beginning with a network design with an excessive number of nodes and gradually eliminating the unnecessary nodes or connections.[35,36] The algorithm performs the connection pruning twice, and the best individual network is selected from the new population. Table 1 shows the parameter values used to ensure the optimal application of the algorithm.

Sigmoidal and linear functions were used for hidden and output nodes, respectively. To start processing data and to avoid saturation problems with the sigmoidal functions in the ANN model, the input and output variables were scaled on a range from 0.1 to 0.9. Thus, the new scaled variables

Pruned Artificial Neural Networks Models

*J. Chem. Inf. Comput. Sci., Vol. 44, No. 5, 2004* **1579**

are as follows: $\hat{S}_{gn}^*$, $\hat{B}_m^*$, $\hat{C}^*$, and $\hat{t}_m^*$ for the input variables and $[\hat{CF}]^*$ and $[\hat{P}]^*$ for the output variables, which correspond to the contribution of each analyte [carbofuran (*CF*) and propoxur (*P*)] to the chromatographic peak. After optimizing the ANN model, estimations should be descaled following the same equation.

## EXPERIMENTAL SECTION

The Levenberg−Marquardt algorithm was used to solve eqs 6−9, obtaining the four estimated coefficients. Convergence of the iterative process was achieved with a tolerance of 0.0001 and a maximum number of 100 iterations. The algorithm software for ANN in C language was run on a PC Pentium IV compatible computer. An overall set of 25 synthetic samples (in triplicate) containing uniformly distributed concentrations of carbofuran (30−150 ng/mL) and propoxur (30−210 ng/mL) was prepared as described elsewhere.[19] Of the three replicates, two randomly chosen replicates were used to design the training set, while the other was used for the generalization set. The performance of the algorithm was tested using various network topologies that were run 30 times. The accuracy of each model was assessed in terms of the SEP for the results obtained for both data sets, that is, SEP$_T$ for the training set, and SEP$_G$ for the generalization set

$$SEP = \frac{100}{\bar{A}_i}\sqrt{\frac{\sum_{i=1}^{n}(A_i - \hat{A}_i)^2}{n}} \qquad (11)$$

where $A_i$ and $\hat{A}_i$ are the experimental and expected values for the analyte concentration in the mixture, $\bar{A}_i$ is the mean of the experimental values of the training set, or of the generalization set, and $n$ is the number of patterns used. The nonparametric Kolmogorov−Smirnov (K−S) and Levene tests were performed using SPSS 12.0 statistical software[37] and used to evaluate the performance of the different models in selecting the most suitable network topology.

## RESULTS AND DISCUSSION

The quantitative resolution of overlapping chromatographic peaks is an area of growing interest for analytical chemists on the grounds of the high discrimination power of current chemometric tools. Despite the fact that a large number of approaches for this purpose have been described in the literature, two hypotheses should be generally assumed prior to their use: first, the aid of spectral discrimination is required, such as LC-DAD systems, and second, the components should be moderately overlapped to obtain fairly good results as strong overlapping poses difficulties for many chemometric methods. To overcome these drawbacks, the goal of this work was to develop a general methodology based on pruned ANNs for the quantification of unresolved chromatographic bands with a high degree of overlapping by using analytical data provided by a single detector, that is, in the absence of spectral discrimination. The approach was tested on the simultaneous determination of *CF* and *P*, two *N*-methylcarbamate pesticides that cannot be resolved by LC even with multistep gradient elution.[38] After hydroly-
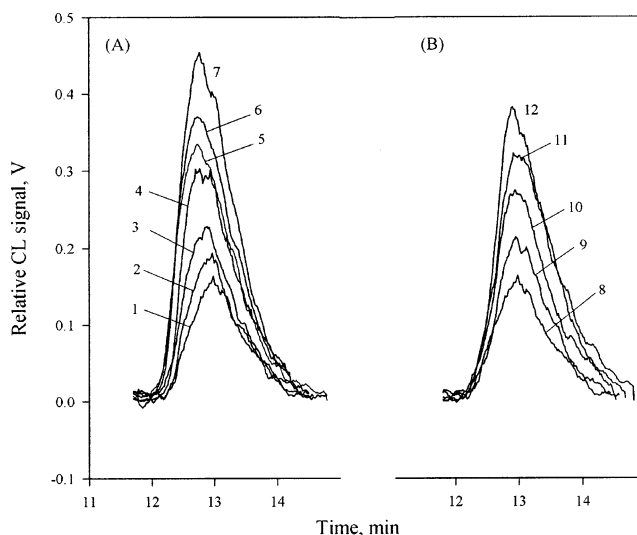


**Figure 1.** Typical chromatograms of samples containing different ratios of carbofuran and propoxur. (A) Curves 1−7 correspond to samples with a fixed amount of carbofuran (30 ng/mL) and variable amounts of propoxur (30−210 ng/mL). (B) Curves 8−12 correspond to samples in which the amount of propoxur was fixed at 30 ng/mL and carbofuran was changed from 30 to 150 ng/mL.

sis of the pesticides and derivatization of their hydrolytic metabolites with dansyl chloride (DNS−Cl), chromatographic data were obtained by monitoring the single CL signal provided by the reaction of the coeluted dansylated metabolites with the classical peroxyoxalate chemiluminescence (PO−CL) reaction based on the bis(2,4,6-trichlorophenyl) oxalate−hydrogen peroxide system.

Figure 1 shows the chromatograms achieved using different mixtures of these pesticides following their hydrolysis and derivatization with DNS−Cl. The samples were separated on a Nova-Pack C$_{18}$ 150 × 3.9 mm (4 $\mu$m) column under isocratic elution with an acetonitrile−water (70:30, v/v) mobile phase flowing at 0.5 mL/min. As can be seen, the two compounds exhibit a rather different chromatographic behavior; in fact, chromatograms with the same concentration of *CF* (30 ng/mL) and a variable concentration of *P* over the range 30−150 ng/mL exhibit a gradual increase in the peak height, $S_m$, and a slight decrease in the residence time, $t_m$, with an increase in the concentration of the latter (Figure 1A). Distinct behavior is observed for chromatograms containing mixtures with a constant concentration of *P* (30 ng/mL) and a variable concentration of *CF* from 30 to 210 ng/mL (Figure 1B). In this case, $t_m$ remains virtually constant, while $S_m$ increases with an increase in the concentration of *CF*, showing a slight increase in sensitivity in comparison to the previous case (compare curves 5 and 12 in Figure 1). Additional conclusions can be derived from these chromatograms with respect to the variation of its profile and dispersion as a function of the pesticide concentration in the mixture. According to the chromatograms in Figure 1B, it is clear that these parameters are closely related to the presence of *CF* in the mixture. In fact, the chromatograms reached the baseline at higher times as the concentration of *CF* increased, in contrast to the behavior shown in Figure 1A for *P*.

The above qualitative study regarding the influence of relative concentrations of *CF* and *P* on the shape of the chromatographic peak suggests that it should be modeled
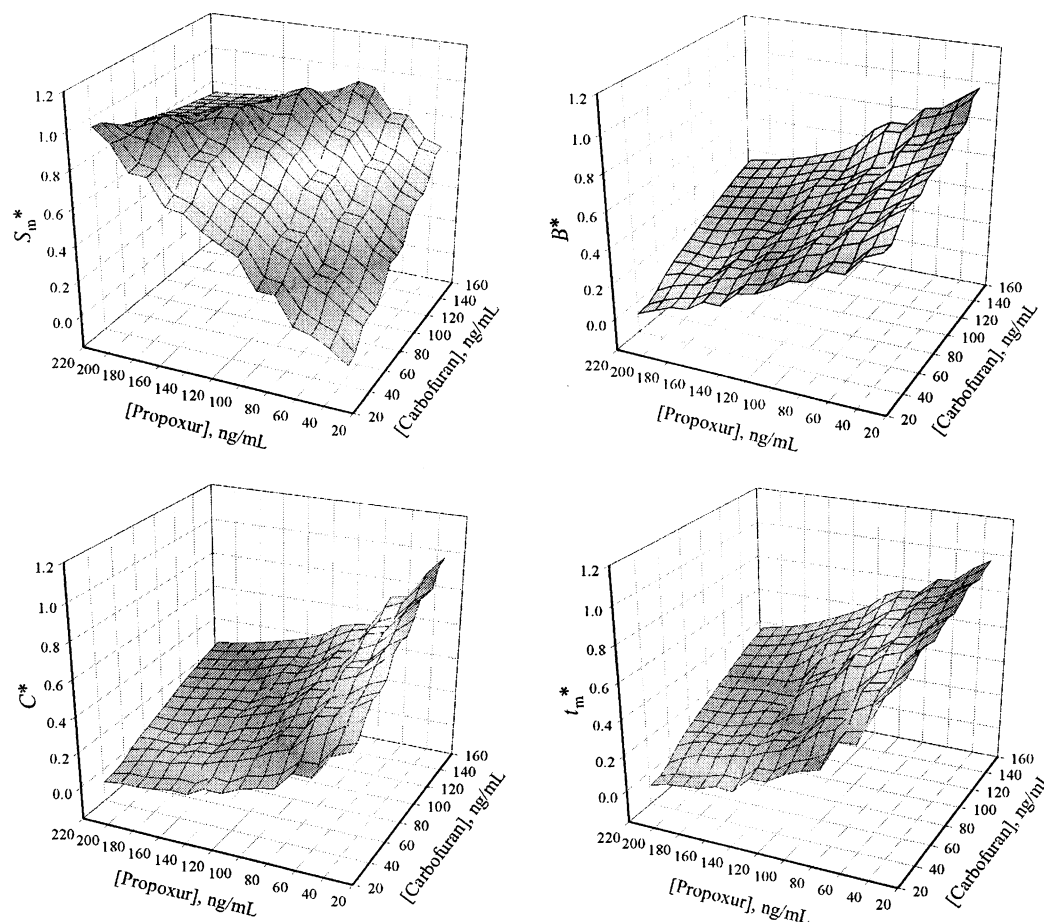
**Figure 2.** Three-dimensional plots showing the dependence of carbofuran and propoxur concentrations on the four parameters estimated by the Weibull curve fitted to the chromatographic peaks.

**Table 2.** Accuracy and Statistical Results of the Algorithm Used with Various Network Topologies (Over 30 Runs)

| analyte | starting topology | connections | | $SEP_T$ | | | | $SEP_G$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mean | SD | mean | SD | best | worst | mean | SD | best | worst |
| carbofuran | 4:4:1 | 18.3 | 3.71 | 3.18 | 0.29 | 2.67 | 3.71 | 3.33 | 0.31 | 2.69 | 3.91 |
| | 4:5:2 | 25.7 | 7.04 | 3.15 | 0.52 | 2.28 | 4.17 | 3.23 | 0.34 | 2.67 | 3.86 |
| propoxur | 4:4:1 | 16.2 | 4.72 | 3.65 | 0.75 | 1.99 | 5.42 | 3.58 | 0.69 | 2.29 | 5.44 |
| | 4:5:2 | 25.7 | 7.04 | 3.15 | 0.52 | 2.28 | 4.17 | 3.13 | 0.66 | 2.05 | 4.76 |

with a predetermined function in order to use its characteristics as inputs to ANNs. This situation reduces ANN complexity and learning time as well as the number of patterns needed to train the network, which is of great practical interest. Taking into account that chromatographic bands are nonsymmetric, the four-parameter Weibull function was found to be the best choice for modeling chromatographic data (see Theoretical Background section). Figure 2 shows the three-dimensional plots of each parameter (the values were scaled between 0 and 1 to improve clarity in comparison) as a function of the *CF* and *P* concentrations in the sample subjected to chromatographic separation. As can be seen, $S_m^*$ and $t_m^*$ are mainly related to the concentration of *P* in the mixture, whereas *B\** and *C\** depend to a higher extent on the concentration of *CF*. This is in agreement with the above conclusions reached solely from observing the chromatograms, thereby confirming that the four-parameter Weibull function can be readily used for modeling the chromatographic data in order to select the inputs to the ANNs.

**Optimization of the Network Topology.** To further investigate the prediction ability of pruning ANNs, in addition to the two outputs models, corresponding to the concentration of both analytes, one output neural network models were also made considering the output layer corresponding to the only analyte to be analyzed. Twenty-five chromatograms provided by samples containing the analytes in [*CF*]/[*P*] ratios from 5:1 to 1:7 were analyzed in triplicate. Fifty chromatographic peaks (two randomly selected from each sample) were used for the training set, while the remaining 25 (one of the three replicated measurements for each sample) were retained for the generalization test. Although this procedure affords lower SEP values, it is appropriate enough to evaluate the ability of the proposed ANN models to solve the addressed analytical problem and to extract significant conclusions.

Table 2 shows the statistical results obtained over 30 runs using one or two outputs networks to quantify both analytes. The 4:4:1 and 4:5:2 network topologies were chosen to start the computation pruned process for one and two outputs

PRUNED ARTIFICIAL NEURAL NETWORKS MODELS

*J. Chem. Inf. Comput. Sci., Vol. 44, No. 5, 2004* **1581**

**Table 3.** Statistical Comparison of the Generalization Ability (SEP$_G$) and Number of Connections ($n$) as Applied to the Determination of Carbofuran and Propoxur Using Overlapping Chromatographic Data

Kolmogorov−Smirnov Test

|  | statistic parameter | carbofuran | | propoxur | |
|---|---|---|---|---|---|
|  |  | 4:4:1 model | 4:5:2 model | 4:4:1 model | 4:5:2 model |
| SEP$_G$ | $Z$ | 0.688 | 0.587 | 0.897 | 0.591 |
|  | $p$ | 0.731 | 0.880 | 0.396 | 0.875 |
| $n$ | $Z$ | 0.919 | 0.716 | 0.676 | 0.716 |
|  | $p$ | 0.367 | 0.684 | 0.751 | 0.684 |

Levene and Student's *t*-Tests

| 4:4:1 vs 4:5:2 | SEP$_G$ | | | | number of connections | | | |
|---|---|---|---|---|---|---|---|---|
|  | Levene test | | student's *t*-test | | Levene test | | student's *t*-test | |
|  | $F$ | $p$ | $t$ | $p$ | $F$ | $p$ | $t$ | $p$ |
| carbofuran | 1.220 | 0.274 | 1.222 | 0.227 | 21.34 | 0.000 | −5.067 | 0.000 |
| propoxur | 0.000 | 0.989 | 2.590 | 0.012 | 9.414 | 0.003 | −6.114 | 0.000 |

topologies, respectively. As can be seen, both models provided quite good results (in terms of accuracy and precision) for determining the concentration of each pesticide from strongly overlapping chromatographic bands: the SEP mean values for the training and generalization tests ranged from 3.13 to 3.65%, whereas the standard deviation (SD) varied between 0.29 and 0.75%.

To ascertain the statistical significance of the observed differences between the means (SEP$_G$ for each compound and the number of connections, $n$) for the proposed ANN models, the nonparametric Kolmogorov−Smirnov test (KS-test) with a signification level, $\alpha$, equal to 0.05 was used to evaluate if the SEP$_G$ and $n$ values follow a normal distribution. As can be seen from the results in Table 3, a normal distribution can be assumed because the critical levels, $p$, were higher than 0.05 in all cases with smaller values of 0.367 and 0.396. At this point, the proposed ANN models for each compound were compared under the hypothesis of normal distribution for SEP$_G$ and $n$, by using the student's *t*-test for differences between means with equal or different variance according to the results previously obtained by the Levene test, in both cases with $\alpha = 0.05$. The statistical results obtained are also shown in Table 3.

When determining *CF*, no significant differences were found between the SEP$_G$ mean values and variances for both ANN models according to the $p$-values provided by the tests: 0.274 and 0.277. The opposite behavior was found when the statistical tests were applied to the number of connections. In short, both ANN models show a similar homogeneity and efficacy for determining *CF* from the overlapping chromatographic peaks: there are no significant differences between their SEP$_G$ mean values. Regarding the determination of *P*, although both models prove to have the same efficacy for determining the SEP$_G$ mean values, the student's *t*-test shows that there are significant differences between the means, providing the two outputs model with smaller values. Similar claims can be made for *CF* as for *P* regarding the statistical parameter values found for the number of connections (see Table 3). According to this study, the following final optimal models can be chosen: 4:3:1 and 4:2:1 network topologies for the single determination of *CF* and *P*, respectively, and a 4:4:2 architecture for the simul-

taneous quantification of both pesticides in the sample.

The simplicity of the proposed ANN architectures permits us to derive straightforward quantitative equation systems for the direct determination of the contribution of each pesticide to the overlapping chromatographic bands using (a) the parameters estimated by the Weibull regression of the peak; (b) the optimized network weights; and (c) the sigmoidal transfer functions. Table 4 shows the quantitative equations corresponding to each network topology, together with its number of connections and the SEP$_T$ and SEP$_G$ values. According to the SEP values, both topologies can be readily used for the quantification of the analytes in the sample, although the two outputs topology provides slightly better results. Table 5 shows the results for various synthetic samples containing variable amounts of *CF* and *P*. As can be seen, samples containing [*CF*]/[*P*] ratios ranging from 5:1 to 1:7 can be accurately resolved with relative errors of less than ca. $\pm 5\%$ in both components. In some cases, relative errors increased for samples with a lower content of any pesticide; however, they decreased when using the two outputs topology.

**Refining Optimized Topologies by Using Heuristic Rules: Chemical Interpretation.** It is a widely accepted fact in the literature that ANNs are effective chemometric tools for solving a great variety of nonlinear analytical problems, although the models used for these purposes are difficult to interpret, and ANNs are therefore considered "black boxes". The current introduction of a pruned process in conjunction with genetic algorithms allows ANN architectures with reduced dimensions to be obtained in such a way that the mathematical transformation of the input vector into the output vector can be easily implemented via software based on easily solved equations, such as those shown in Table 4. Because these models are simple, they can be more easily approached and explained. However, using heuristic rules, simpler models could be obtained from which quality chemical information can be derived to explain the analytical problem at hand. In this study, heuristic rules were used based on the following hypotheses: (a) the values of the sigmoidal functions involved in their respective quantification equation used to determine pesticide concentration are calculated over the range studied for the input variables, $\hat{S}_m^*$, $\hat{B}^*$, $\hat{C}^*$, and $\hat{t}_m^*$, to establish their relative contribution for the determination of *CF* and *P*; (b) according to these results, the base sigmoidal function is chosen, while the others are associated to it in a positive or negative form; and (c) a predetermined range for the input variables can be established to neglect some associated sigmoidal term with a view to reducing the network topology. As a criteria, the term is neglected if its value is smaller than 5% of that corresponding to the output variables $[\hat{CF}]^*$ or $[\hat{P}]^*$. In light of the equations shown in Table 4, the most notable features of the proposed models are discussed below from a computational and chemical point of view as are the heuristic rules used.

*Model for Carbofuran (4:3:1 ANN).* Based on the equations shown in Table 4 and after evaluating the relative contribution of each sigmoidal function, it follows that $\hat{h}_1$ is the sigmoidal base term and bears a direct relationship to the concentration of *CF* in the sample. This function essentially depends on $\hat{B}^*$ and $\hat{C}^*$ parameters, in this order, and also on $\hat{S}_m^*$ and $\hat{t}_m^*$, albeit to a smaller extent. These

**Table 4.** Quantification Equations and Accuracy Provided by the Optimized Network Topologies as Applied to the Determination of Carbofuran and Propoxur from Overlapping Chromatographic Peaks

|  | 4:3:1 network topology | 4:2:1 network topology | 4:4:2 network topology |
|---|---|---|---|
| quantification equations | $[\hat{CF}]^* = -0.49 + 0.54\hat{h}_1 + 1.21\hat{h}_2 - 0.42\hat{h}_3$ | $[\hat{P}]^* = -0.23 + 1.03\hat{h}_1 + 0.80\hat{h}_2$ | $[\hat{CF}]^* = 0.03 + 0.99\hat{h}_1 + 0.76\hat{h}_2 - 0.65\hat{h}_3 - 1.66\hat{h}_4$ $[\hat{P}]^* = -0.72\hat{h}_1 + 0.99\hat{h}_2 + 1.00\hat{h}_4$ |
| sigmoidal functions | $\hat{h}_1 = 1/(1 + \exp(-0.57 + 1.73\hat{S}_m^* - 6.94\hat{B}^* - 3.27\hat{C}^* + 1.23\hat{t}_m^*))$ $\hat{h}_2 = 1/(1 + \exp(3.60 - 2.80\hat{S}_m^* - 1.49\hat{B}^* - 1.17\hat{C}^* + 0.80\hat{t}_m^*))$ $\hat{h}_3 = 1/(1 + \exp(2.22\hat{S}_m^*))$ | $\hat{h}_1 = 1/(1 + \exp(-1.65\hat{S}_m^* + 1.52\hat{B}^* + 0.58\hat{t}_m^*))$ $\hat{h}_2 = 1/(1 + \exp(8.93\hat{B}^* - 4.58\hat{C}^* - 0.60\hat{t}_m^*))$ | $\hat{h}_1 = 1/(1 + \exp(-3.43 - 1.47\hat{S}_m^* - 1.29\hat{B}^* - 1.61\hat{C}^* - 0.76\hat{t}_m^*))$ $\hat{h}_2 = 1/(1 + \exp(1.24 - 2.66\hat{S}_m^*))$ $\hat{h}_3 = 1/(1 + \exp(1.03\hat{S}_m^*))$ $\hat{h}_4 = 1/(1 + \exp(1.59 - 1.41\hat{S}_m^* + 4.13\hat{C}^* + 2.27\hat{t}_m^*))$ |
| connections | 15 | 9 | 20 |
| SEP$_T$ | 2.68% | 2.00% | 2.28% |
| SEP$_G$ | 2.69% | 2.29% | 2.67% (*CF*), 2.06% (*P*) |

**Table 5.** Comparison of the Quality Achieved for the Quantification of Carbofuran and Propoxur Using One and Two Outputs Network Topologies

| [carbofuran]/ [propoxur] | carbofuran (ng/mL) | | | propoxur (ng/mL) | | |
|---|---|---|---|---|---|---|
|  | added | estimated | | added | estimated | |
|  |  | 4:3:1 model | 4:4:2 model |  | 4:2:1 model | 4:4:2 model |
| 1:1 | 30 | 33.3 | 32.3 | 30 | 31.0 | 32.2 |
| 1:2 | 30 | 28.3 | 28.6 | 60 | 59.7 | 57.9 |
| 1:3 | 30 | 29.3 | 27.9 | 90 | 89.8 | 89.7 |
| 1:4 | 30 | 33.5 | 30.1 | 120 | 123.6 | 126.2 |
| 1:5 | 30 | 30.7 | 28.1 | 150 | 150.8 | 153.4 |
| 1:6 | 30 | 29.4 | 31.2 | 180 | 178.7 | 177.2 |
| 1:7 | 30 | 31.1 | 31.6 | 210 | 211.0 | 213.6 |
| 2:1 | 60 | 56.5 | 56.5 | 30 | 28.8 | 30.4 |
| 2:2 | 60 | 58.0 | 59.8 | 60 | 59.3 | 59.1 |
| 2:3 | 60 | 61.0 | 63.0 | 90 | 90.3 | 91.8 |
| 2:4 | 60 | 60.0 | 61.2 | 120 | 120.9 | 122.5 |
| 2:5 | 60 | 56.1 | 57.4 | 150 | 151.5 | 150.3 |
| 2:6 | 60 | 58.5 | 57.4 | 180 | 173.9 | 177.5 |
| 3:1 | 90 | 89.8 | 88.4 | 30 | 29.4 | 28.4 |
| 3:2 | 90 | 90.0 | 90.4 | 60 | 60.4 | 61.8 |
| 3:3 | 90 | 88.1 | 89.9 | 90 | 89.2 | 89.1 |
| 3:4 | 90 | 90.5 | 92.7 | 120 | 120.4 | 118.6 |
| 3:5 | 90 | 93.6 | 91.6 | 150 | 152.3 | 151.3 |
| 4:1 | 120 | 123.2 | 121.2 | 30 | 30.9 | 30.0 |
| 4:2 | 120 | 117.6 | 116.3 | 60 | 57.2 | 59.7 |
| 4:3 | 120 | 117.9 | 117.8 | 90 | 90.4 | 90.6 |
| 4:4 | 120 | 122.8 | 122.3 | 120 | 122.9 | 118.9 |
| 5:1 | 150 | 152.2 | 153.8 | 30 | 32.9 | 33.1 |
| 5:2 | 150 | 148.0 | 149.5 | 60 | 61.8 | 62.5 |
| 5:3 | 150 | 149.9 | 150.3 | 90 | 88.1 | 87.0 |

dependencies are direct for $\hat{B}^*$ and $\hat{C}^*$ and inverse for $\hat{S}_m^*$ and $\hat{t}_m^*$. Regarding the other sigmoidal functions, $\hat{h}_2$ (related in a positive manner to $[\hat{CF}]^*$) depends on the input variables in the following order of importance: $\hat{S}_m^*$, $\hat{B}^*$, $\hat{C}^*$, and $\hat{t}_m^*$; whereas $\hat{h}_3$ (bearing an inverse linear relationship with $[\hat{CF}]^*$) only depends on the $\hat{S}_m^*$ value. From the foregoing it follows that (a) $\hat{B}^*$ and $\hat{C}^*$ are the two parameters of the Weibull curve which most strongly influence the quantification of *CF* in the sample; (b) the peak height value, $\hat{S}_m^*$, also contributes to this, although to a lesser extent; and (c) the influence of the residence time, $\hat{t}_m^*$, is practically negligible. This quantitative information, provided by the pruned ANN model, is in agreement with the qualitative conclusions extracted from Figures 1 and 2. In sum, the dispersion and shape of the chromatographic bands are closely related to the concentration of *CF* in the analyzed sample.

*Model for Propoxur (4:2:1 ANN).* In this case, both sigmoidal functions exhibit an appreciable influence on the determination of $[\hat{P}]^*$, although $\hat{h}_1$ to a higher extent. The respective dependencies are as follows: $\hat{h}_1$ is closely related to $\hat{S}_m^*$ and $\hat{B}^*$, and to a lesser extent to $\hat{t}_m^*$; the dependence being direct with the peak height and inverse with the other variables; $\hat{h}_2$ depends on $\hat{B}^*$, $\hat{C}^*$, and $\hat{t}_m^*$ in this order of importance. According to these results we can infer that the $\hat{S}_m^*$ variable is key to determining $[\hat{P}]^*$ in the sample (see also Figures 1 and 2), whereas the contribution of the other variables provides a better fitting in its quantification. This asseveration can be confirmed if a heuristic rule is assumed in order to neglect the $\hat{h}_2$ term in the corresponding quantification equation. In fact, at high values of $\hat{B}^*$, $\hat{C}^*$, and $\hat{t}_m^*$, such as $\hat{B}^* \geq 0.68 \wedge \hat{C}^* \geq 0.54 \wedge \hat{t}_m^* \geq 0.7$ ($\hat{B} \geq 1.03 \wedge \hat{C} \geq 2.12 \wedge \hat{t}_m \geq 13$), $\hat{h}_2$ can be neglected because its value is smaller than 5% of that corresponding to the $[\hat{P}]^*$. Note that this heuristic rule is useful for any $\hat{S}_m^*$ value. Chemically, the range of the input variables involved in the heuristic rule corresponds to samples with a high content of *CF* and [*CF*]/[*P*] ratios. In sum, for these kind of samples, the content of *P* can be easily determined by using the following equation:

$$[\hat{P}]^* = -0.23 + \frac{1.03}{1 + \exp(-1.65\hat{S}_m^* - 1.52\hat{B}^* - 0.58\hat{t}_m^*)} \quad (12)$$

Note that, in this case, $[\hat{P}]^*$ is independent of the $\hat{C}^*$ parameter associated with the shape of the chromatographic band, which, as stated above, is closely related to the *CF* concentration in the sample.

*Model for Both Pesticides (4:4:2 ANN).* When this ANN model is used, the most relevant addend for both quantification equations is the $\hat{h}_2$ function, which only depends on the $\hat{S}_m^*$ value, exhibiting a more pronounced effect on the determination of propoxur (see coefficients for $\hat{h}_2$ in these quantification equations). The other sigmoidal functions are associated to $\hat{h}_2$ to a different extent in order to obtain the best fitting for the concentration of both analytes. Thus, for $[\hat{CF}]^*$, $\hat{h}_1$ is added in a positive form and depends directly on the input variables $\hat{C}^*$, $\hat{S}_m^*$, $\hat{B}^*$, and $\hat{t}_m^*$ (in that order), whereas $\hat{h}_3$ and $\hat{h}_4$ provided a negative dependence, the first with respect to $\hat{S}_m^*$ and the latter basically on $\hat{C}^*$. Regarding $[\hat{P}]^*$, $\hat{h}_1$ and $\hat{h}_4$ exhibit opposite effects. This ANN model

PRUNED ARTIFICIAL NEURAL NETWORKS MODELS

*J. Chem. Inf. Comput. Sci., Vol. 44, No. 5, 2004* **1583**

can be simplified by using a similar heuristic rule to that used above for $[\hat{P}]^*$. In fact, the $\hat{h}_4$ term can be neglected (errors < 4%) if $\hat{C}^* \geq 0.55 \wedge \hat{t}_m^* \geq 0.65$ ($\hat{C} \geq 2.13 \wedge \hat{t}_m \geq 12.97$). From the foregoing, the following conclusions can be drawn: (a) $\hat{S}_m^*$ is the parameter of the Weibull curve which most strongly influences the quantification of the pesticides in the sample, especially for $[\hat{P}]^*$; (b) the $\hat{B}^*$ and $\hat{C}^*$ parameters exhibit a positive influence on the determination of $[\hat{CF}]^*$, as in the one output ANN model proposed for this compound; and (c) a simplified model can be obtained for samples with a high content of *CF* and $[CF]/[P]$ ratios, similar to the results found for *P* based on its one output ANN model.

## CONCLUSIONS

As shown in this study, neural networks in conjunction with genetic algorithms and heuristic rules provide straightforward ANN models for the quantification of chromatographic peaks with a high overlapping degree by using a single detector; that is, without the aid of spectral discrimination. The proposed pruned ANN models are quite simple and easy to interpret from both a computational and a chemical point of view, in contrast to the classical assertion that ANNs act like "black boxes". Although the proposed approach was tested for asymmetric chromatographic bands (the four parameters, estimated by NLR, of the Weibull curve fitted to the chromatographic peaks were used as input variables to the ANNs), it can also be useful for symmetric chromatographic bands; in this case, the characteristic parameters of the Gaussian curve could be used as inputs. Finally, taking into account that the proposed methodology is practically independent of the degree of overlapping, in theory any composition of the mobile phase could be used in the chromatographic separation, thus reducing residence time (increasing the sample throughout) and increasing sensitivity if the analytical signal is related to its composition, as is the case of the PO−CL detection.

## ACKNOWLEDGMENT

## REFERENCES AND NOTES

(1) Massart, D. L.; Vandeginste, B. G. M.; Buydens, L. M. C.; DeJong S.; Lewi, P. J.; Smeyers-Verbeke, J. *Handbook of chemometrics and qualimetrics. Part A. Data handling in science and technology;* Elsevier: Amsterdam, 1997; Vol. 20 A.

(2) Caselli, M.; de-Gennaro, G.; Ielpo, P. Application of a fast deconvolution method to high performance ion chromatography peaks. *LC-GC-Europe* **2002**, *15*, 104, 106, 108.

(3) Dunkerley, S.; Brereton, R. G.; Crosby, J. A comparison of deconvolution methods as applied to high-performance liquid chromatography-diode array detector-electrospray mass spectrometry of 2- and 3-hydroxypyridine at varying pH in the presence of severely tailing peak shapes. *Chemom. Intell. Lab. Syst.* **1999**, *48*, 99−119.

(4) van-Zomeren, P. V.; Darwinkel, H.; Coenegracht, P. M. J.; de Jong, G. J. Comparison of several curve resolution methods for drug impurity profiling using high-performance liquid chromatography with diode-array detection. *Anal. Chim. Acta* **2003**, *487*, 155−170.

(5) Pasadakis, N.; Gaganis, V.; Varotsis, N. Accurate determination of aromatic groups in heavy petroleum fractions using HPLC−UV-DAD. *Fuel* **2001**, *80*, 147−153.

(6) Sanchez, F. C.; Rutan, S. C.; Garcia, M. D. G.; Massart, D. L. Resolution of multicomponent overlapped peaks by the orthogonal projection approach, evolving factor analysis and window factor analysis. *Chemom. Intell. Lab. Syst.* **1997**, *36*, 153−164.

(7) Gross, G. M.; Prazen, B. J.; Synovec, R. E. Parallel column liquid chromatography with a single multiwavelength absorbance detector for enhanced selectivity using chemometric analysis. *Anal. Chim. Acta* **2003**, *490*, 197−210.

(8) Fraga, C. G.; Bruckner, C. A.; Synovec, R. E. Increasing the number of analyzable peaks in comprehensive two-dimensional separations through chemometrics. *Anal. Chem.* **2001**, *73*, 675−683.

(9) Liang, Y. Z.; Hamalainen, M. D.; Kvalheim, O. M.; Andersson, R. Assessment of peak origin and purity in one-dimensional chromatography by experimental design and heuristic evolving latent projections. *J. Chromatogr. A* **1994**, *662*, 113−122.

(10) Garrido-Frenich, A.; Martinez-Galera, M.; Gil-Garcia, M. D.; Martinez-Vidal, J. L.; Catasus, M.; Marti, L.; Mederos, M. V. Resolution of HPLC-DAD highly overlapping analytical signals for quantitation of pesticide mixtures in groundwater and soil using multicomponent analysis and neural networks. *J. Liq. Chromatogr. Relat. Technol.* **2001**, *24*, 651−668.

(11) Li, Y. B.; Huang, X. Y.; Sha, M.; Meng, X. S. Resolution of overlapping chromatographic peaks by radial basis function neural network. *Sepu* **2001**, *19*, 112−115.

(12) Galeano-Diaz, T.; Guiberteau, A.; Ortiz, J. M.; Lopez, M. D.; Salinas, F. Use of neural networks and diode-array detection to develop an isocratic HPLC method for the analysis of nitrophenol pesticides and related compounds. *Chromatographia* **2001**, *53*, 40−46.

(13) Andrews R.; Diederich, J.; Ticle, A. Survey and critique of techniques for extracting rules from trained artificial neural networks. *Knowledge Based Syst.* **1998**, *8*, 373−389.

(14) Diederich, A. J.; Tickle, B.; Andrews, R.; Golea, M. The truth will come to light: Directions and challenges in extracting the knowledge embedded within trained artificial neural networks. *IEEE Trans. Neural Networks* **1998**, *9*, 1057−1068.

(15) Schmitz, G. P. J.; Aldrich, C.; Gouws, F. S. ANN-DT: An algorithm for extraction of decision trees from artificial neural networks. *IEEE Trans. Neural Networks* **1999**, *10*, 1392−1402.

(16) Setieno, R.; Leow, W. K.; Zurada, J. M. Extraction of rules from artificial neural networks for nonlinear regression. *IEEE Trans. Neural Networks* **2002**, *13*, 564−577.

(17) Williams, P. M. Bayesian regularization and pruning using a Laplace prior. *Neural Comput.* **1995**, *7*, 117−143.

(18) Hervás, C.; Toledo, R.; Silva, M. Use of pruned computational neural networks for processing the response of oscillating chemical reactions with a view to analyzing nonlinear multicomponent mixtures. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1083−1092.

(19) Orejuela, E.; Silva, M. Monitoring some phenoxy-type *N*-methylcarbamate pesticide residues in fruit juices using high-performance liquid chromatography with peroxyoxalate-chemiluminescence detection. *J. Chromatogr. A* **2003**, *1007*, 197−201.

(20) Marquardt, D. W. An algorithm for least-squares estimation of nonlinear parameters. *J. Soc. Indust. Appl. Mathem.* **1963**, *11*, 431−441.

(21) Hartley, H. O. The modified Gauss−Newton method for the fitting of nonlinear regression functions by least-squares. *Technometrics* **1961**, *3*, 269−280.

(22) Rawlings, J. O.; Pantula, S. G.; Dickey, D. *Applied regression analysis: A research tool*; Springer-Verlag: New York, 1998.

(23) Angeline, P. J.; Saunders, G. M.; Pollack, J. B. An evolutionary algorithm that constructs recurrent neural networks. *IEEE Trans. Neural Networks* **1994**, *5*, 54−65.

(24) Yao, X. Evolving Artificial Neural Networks. *IEEE Trans. Neural Networks* **1999**, *9*, 1423−1447.

(25) García, N.; Hervás, C.; Muñoz, J. Covnet: Cooperative coevolution of neural networks. *IEEE Trans. Neural Networks* **2003**, *14*, 575−596.

(26) Hervás, C.; Zurera, G.; García, R. M.; Martínez, J. A. Optimisation of a computational neural network for its application to the prediction of microbial growth in foods. *Food Sci. Technol.* **2001**, *7*, 1−5.

(27) Hervás, C.; Algar, J. A.; Silva, M. Correction of temperature variations in kinetic-based determinations by use of pruning computational neural networks in conjunction with genetic algorithms. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 724−731.

(28) Minai, A. A.; Williams, R. J. Back-propagation heuristics: A study of the extended delta-bar-delta, IEEE International Joint Conference on Neural Networks, San Diego, CA, 1990; pp 595−600.

(29) Goldberg, D. E. Real-coded genetic algorithms, virtual alphabets, and blocking. *Complex Systems* **1991**, *5*, 139−167.

(30) Eshelman, L. J.; Schaffer, J. D. *Real-coded genetic algorithms and interval-schemata. Foundations of Genetic Algorithms 2;* Morgan Kaufmann: San Mateo, CA, 1993; pp 187−202.

(31) Ono, I.; Kita, H.; Kobayashi, S. A robust real-coded genetic algorithm using unimodal normal distribution crossover augmented by uniform crossover: effects of self-adaptation of crossover probabilities. Genetic and Evolutionary Computation Conference (GECCO-99), Martin Kaufmann: San Mateo, CA, 1999; pp 496−503.

(32) Herrera, F.; Lozano, M. Gradual distributed real-coded genetic algorithms. *IEEE Trans. Evolut. Comput.* **2000**, *4*, 43−63.

(33) Whitley, D.; Kauth, J. GENITOR: A different genetic algorithm. Rocky Mountain Conference on Artificial Intelligence, Denver, CO, 1988; pp 118−130.

(34) Beasley, D.; Bull, D.; Martin, R. An overview of genetic algorithms: Part 1, Fundamentals. *University Computing* **1993**, *15*, 58−69.

(35) Cun, L. Y.; Denker, J. S.; Solla, S. A. *Optimal brain damage advances in neural information processing systems;* Morgan Kaufmann: San Mateo, CA, 1990; pp 598−605.

(36) Hassibi, B.; Stork, D. G.; Wolff, G. J. Optimal brain surgeon and general networks pruning. IEEE International Conference on Neural Networks, San Francisco, CA, 1993; pp 291−299.

(37) SPSS, Advanced Models. Copyright 12.0 SPSS Inc., 2003, Chicago, IL.

(38) DeKok, A.; Hiemstra, M.; Brinkman, U. A. T. Low ng $L^{-1}$ level determination of twenty *N*-methylcarbamate pesticides and twelve of their polar metabolites in surface water via offline solid-phase extraction and high-performance liquid chromatography with post-column reaction and fluorescence detection. *J. Chromatogr. A* **1992**, *623*, 265−276.