



Quantitative structure–electrochemistry relationship for variously-substituted 9, 10-anthraquinones using both an heuristic method and a radial basis function neural network

Huitao Liu^{a,*}, Ping Han^a, Yingying Wen^a, Feng Luan^a, Yuan Gao^a, Xiuyong Li^b

^a Department of Applied Chemistry, Yantai University, Yantai 264005, PR China

^b Yantai Inspection and Quarantine Bureau, Yantai 264000, PR China

ARTICLE INFO

Article history:

Received 9 March 2009

Received in revised form

26 July 2009

Accepted 27 July 2009

Available online 14 August 2009

Keywords:

Dye intermediate

9,10-Anthraquinone

Quantitative structure–property

relationship (QSPR)

Heuristic method (HM)

Radial basis function neural network

(RBFNN)

Half-wave potential ($E_{1/2}$)

ABSTRACT

Quantitative structure–property relationship models correlating the half-wave potentials ($E_{1/2}$) of the dye intermediate 9, 10-anthraquinone and its derivatives were developed using both linear and non-linear modelling approaches. Descriptors calculated from molecular structures alone were used to represent the $E_{1/2}$ of the 9, 10-anthraquinones. An heuristic method was used to select the most appropriate molecular descriptors whilst a linear, quantitative structure–property relationship model was developed; using the selected descriptors, a radial basis function neural network was employed for the non-linear model development. The statistical parameters provided by the heuristic model were $R^2 = 0.945$; $F = 98.04$; $RMS = 0.0360$ for the training set and $R^2 = 0.863$; $F = 18.903$; $RMS = 0.0600$ for the test set. The radial basis function neural network model gave better results with $R^2 = 0.960$; $F = 610.336$; $RMS = 0.0299$ for the training set and $R^2 = 0.863$; $F = 18.946$; $RMS = 0.0476$ for the test set.

© 2009 Elsevier Ltd. All rights reserved.

1. Introduction

Anthraquinone (Aq) derivatives, which constitute the largest and most important group of natural quinones, enjoy widespread usage in various areas of chemistry, biochemistry and pharmacology [1]. The hydroxy-derivatives of 9, 10-anthraquinone provide both natural and synthetic dyes that are mostly of a red hue. Aq derivatives are also employed as redox catalysts in various industrial processes, this having led to extensive studies of their electrochemical behaviour, namely half-wave potential ($E_{1/2}$), which is a characteristic constant for a reversible oxidation–reduction system and which can be useful in predicting other electrochemical properties of organic compounds [2].

Various methods can be used to determine the half-wave potential of organic, inorganic and organometallic compounds [3,4]. Quantitative structure–activity/property relationship (QSAR/QSPR) studies, which are one of the most important areas in chemometrics, are

mathematical equations relating chemical structure to a wide range of physical, chemical, biological and technological properties [5]. It is well known that various chemical characteristics, which are closely related to molecular structure, can be calculated or predicted using various methods; QSAR/QSPR analysis quantifies the relationship between property and structure. The method comprises a reliable statistical model for the prediction of property/behaviour for novel chemicals and analytical systems. These relationships also seek to identify and isolate the most important structural descriptors that affect given physicochemical properties [5]. The advantage of this method over others lies in the fact that it requires only a knowledge of chemical structure and is not dependent on the experimental process. In recent years, numerous quantitative QSAR/QSPR models have been introduced for calculating the physicochemical properties of molecules from chemical structure; the applications of QSAR/QSPR in electrochemistry are described [4]. Bahram Hemmateenejad and Mahdiah Yazdani established a QSPR model to predict the $E_{1/2}$ of steroids by multiple linear regression (MLR) and principle component regression (PCR) analysis [5]. Nesmerak et al. employed QSPR approach to investigate the electrooxidation of new benzoxazines as a model of metabolic degradation [6]. Shamsipur and Hemmateenejad employed PCR and principle component-artificial neural network

* Corresponding author at: Department of Applied Chemistry, 32 Qingquan Road, Yantai University, Yantai 264005, PR China. Tel./fax: +86 535 6902401.

E-mail address: liuht-ytu@163.com (H. Liu).

(PC-ANN) models in a QSPR study of some organic compounds [7]. Mojtaba Shamsipur et al. had shown the application of the QSPR model in predicting $E_{1/2}$ of 9, 10-anthraquinones [1]. In ref. 1, a linear five-parametric equation was obtained which consisted of three constitutional descriptors, one topological descriptor and one electronic descriptor. However, the R^2 of the equation was 0.888, which was smaller than that of our four-descriptor linear model including two electrostatic descriptors and two quantum-chemical descriptors. In addition, no external test sets were selected to validate the model in the ref. 1 and only linear model was employed.

The Heuristic method (HM) usually produces correlations between 2 and 5 times faster than other methods, with comparable accuracy [8]. The rapidity of calculations using the heuristic method render it the first method of choice in practical research; thus, in this work HM was used. The main aim of the present study was to develop linear (HM) and non-linear (RBFNN) QSPR models to predict $E_{1/2}$ values of Aq derivatives based on the descriptors calculated from their molecular structures alone. According to the best of our knowledge, this is the first report involving prediction of $E_{1/2}$ for Aq derivatives using RBFNN.

2. Experimental section

2.1. Dataset

The dataset was available from the literature reported by Shamsipur et al. [1]. The general structure of the Aq derivatives is shown in Fig. 1, which also shows both the substituents and the experimental $E_{1/2}$ values for the Aq derivatives. The dataset was randomly divided into a training set of 25 compounds together with a test set of 8 compounds. The training set was applied to adjust the parameters and establish the QSPR models and the test set was applied to evaluate their prediction ability.

2.2. Descriptor generation

The structures of the Aq derivatives were drawn with ISIS Draw 2.4 program. All molecules were preoptimized using molecular mechanics MM + method in the HyperChem program [9]. Then a more precise optimization is done with semiempirical PM3 method in the MOPAC 6.0 program [10]. The output files exported from MOPAC were transferred into the software CODESSA [11] for generating descriptors. In this work, constitutional, topological, geometrical, electrostatic and quantum-chemical descriptors were calculated, and 420 descriptors were calculated.

2.3. Heuristic method

It was impossible to use many descriptors to build a QSPR model, so it was necessary to select the most important descriptors which played major roles in the feature of $E_{1/2}$ of these Aq derivatives. In this study the HM method in CODESSA was employed to select descriptors based on the training set and establish the linear

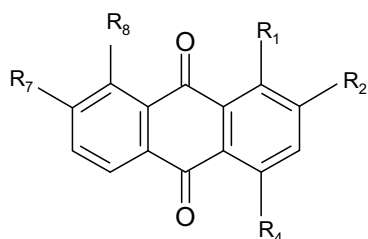


Fig. 1. General structure of the Aq derivatives used in this study.

model. After calculating a large number of descriptors, a feature selection step was applied to reduce the large set of descriptors to a suitable number without losing any significant information. The HM of the descriptor selection proceeds with a pre-selection by eliminating descriptors that (i) are not available for each structure; (ii) have a small variation in magnitude for all structure; (iii) have a Fisher F-criterion below 1.0; and (iv) have t -values less than the user-specified value (by default 0.1), etc. This procedure ordered the descriptors by the decreasing correlation coefficient when used in one-parameter correlations. As a next step, the program calculated the pair correlation matrix of descriptors and further reduced the descriptor pool by eliminating highly correlated descriptors.

After the pre-selection of descriptors, multi-linear regression models were developed in a stepwise procedure. Then, descriptors and correlations were ranked in the light of the values of the F -test and the correlation coefficient. Starting with the top descriptor from the list, two-parameter correlations are calculated. In the following steps new descriptors were added one by one until the pre-selected number of descriptors in the model is achieved. The final result is a list of the 10 best models according to the values of the F -test and correlation coefficient. The goodness of correlation is tested by the coefficient regression (R^2), the F -test (F), and the standard deviation (s^2) [12].

2.4. Radial basis function neural networks

RBFNN, one type of neural networks, has been widely employed to establish non-linear model. The theory of RBFNN has been adequately described elsewhere [13,14]. Here, only a brief description of the RBFNN principle is shown. The architecture of RBFNN is shown in Fig. 2, which consists of three layers: input layer, hidden layer and output layer. The input layer does not process the information; it only distributes the input vectors to the hidden layer. The hidden layer of RBFNN consists of a number of RBF units (n_h) and bias (b_k). Each neuron on the hidden layer employs a radial basis function as non-linear transfer function to operate on the input data. The more often used RBF is a Gaussian function that is characterized by a center (c_j) and width (r_j). The RBF functions by measuring the Euclidean distance between input vector (X) and the radial basis function center (c_j) and performs the non-linear transformation with RBF in the hidden layer as given below:

$$h_j = \exp(-\|X - c_j\|^2 / r_j^2) \quad (1)$$

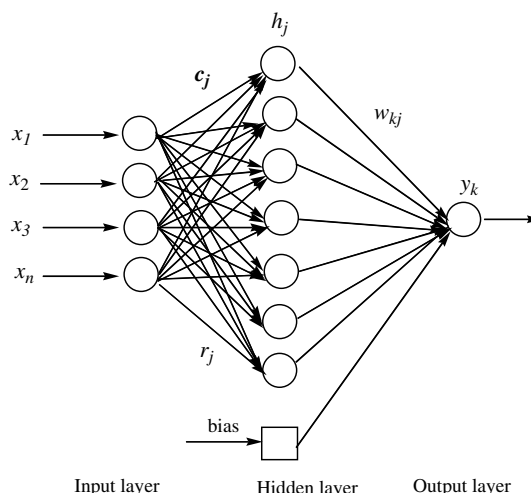


Fig. 2. The architecture of RBFNN.

In which h_j is the notation for the output of the j_{th} RBF unit. For the j_{th} RBF, c_j and r_j are the center and the width, respectively. The operation of the output layer is linear, which is given in Eq. (2):

$$y_k(X) = \sum_{j=1}^{n_h} w_{kj} h_j(X) + b_k \quad (2)$$

where y_k is the k_{th} output unit for the input vector X , w_{kj} is the weight connection between the k_{th} output unit and the j_{th} hidden layer unit, and b_k is the bias.

From Eqs. (1) and (2), we can see that designing RBFNN involves selecting centers, number of hidden layer units, width and weights. There are various ways for selecting the centers, such as random subset selection, K-means clustering, orthogonal least squares learning algorithm, RBF-PLS, etc. In the present paper, a forward subset selection routine was used to select the centers from training set samples. The widths of the radial basis function can either be chosen the same for all the units or can be chosen different for each unit. In this paper, considerations were limited to the Gaussian functions with a constant width, which was the same for all units. The adjustment of the connection weight between the hidden layer and output layer is performed using a least squares solution after the selection of centers and width of radial basis functions. The overall performance of RBFNN is evaluated in terms of root mean squared error (RMS) according to the equation below:

$$RMS = \sqrt{\frac{\sum_{i=1}^{n_s} (y_k - \hat{y}_k)^2}{n_s}} \quad (3)$$

Where y_k is the desired output, \hat{y}_k is the actual output of the network, and n_s is the number of compounds in analyzed set.

All calculation programs implementing RBFNN were written in M-file based on a basis MATLAB script for RBFNN. RBFNN toolbox in MATLAB 7.0 was used to develop RBFNN. The scripts were run on a personal computer.

3. Results and discussion

3.1. Results of the HM

HM was used to select the descriptors and establish the linear model for the prediction of $E_{1/2}$ of Aq derivatives using all the descriptors. After the heuristic reduction, the pool of descriptors was reduced from 420 to 120. At the same time, to avoid the “overparameterization” of the model, an increase of the R^2 value of less than 0.02 was chosen as the breakpoint criterion. We can see that four descriptors were eventually selected according to the principles of HM. The statistical data were given in Table 2. Detailed explanations about these descriptors were found in CODESSA [11]. The predicted data by HM were shown in Table 1. The statistical parameters of the best four descriptors of the linear model were listed in Table 3. From Table 3, we can see that four descriptors have the major effect on $E_{1/2}$. Among them, only RNCS Relative negative charged surface area (SAMNEG*RNCG)[Zefirov's PC] (RNCS) has positive effect on the $E_{1/2}$, which means that the $E_{1/2}$ value will increase with the increasing of this descriptor. While the other three descriptors, Min valency of a C atom (MVCA), PNSA-1 partial negative surface area [Zefirov's PC] (PNSA-1) and Min e–e repulsion for a C–O bond (MEERCOB) have the negative effect on the $E_{1/2}$, which means that the $E_{1/2}$ value will reduce with the increasing of these descriptors. The statistical parameters for the training set are $R^2 = 0.945$, $F = 98.04$, and $RMS = 0.0360$. With the test set, the

Table 1
The substituents and experimental values of anthraquinone derivatives.

Aq	R ₁	R ₂	R ₄	R ₇	R ₈	The experimental $E_{1/2}$	The predicted $E_{1/2}$	
							HM	RBFNN
Aq ₁ *	H	H	H	H	H	−0.911	−0.9541	−0.9129
Aq ₂	H	CH ₃	H	H	H	−0.925	−0.9266	−0.9194
Aq ₃	H	C ₂ H ₅	H	H	H	−0.924	−0.9294	−0.9377
Aq ₄	OH	H	H	H	H	−0.728	−0.6548	−0.6647
Aq ₅	OH	CH ₃	H	H	H	−0.639	−0.6499	−0.6500
Aq ₆	OH	CH ₂ OCH ₃	H	H	H	−0.679	−0.6951	−0.6907
Aq ₇	OH	CH ₂ OC ₂ H ₅	H	H	H	−0.728	−0.7362	−0.7463
Aq ₈ *	OH	CH ₂ O-n-Pr	H	H	H	−0.732	−0.7217	−0.7367
Aq ₉	OH	CH ₂ O-n-But	H	H	H	−0.736	−0.7529	−0.7584
Aq ₁₀	OH	CH ₂ O-isoBut	H	H	H	−0.736	−0.7527	−0.7601
Aq ₁₁	OH	CH ₂ O(CH ₃) ₂ OH	H	H	H	−0.730	−0.7546	−0.7410
Aq ₁₂	OH	CH ₂ OCH ₂ CH ₂ OH	H	H	H	−0.843	−0.8105	−0.7967
Aq ₁₃	OH	(CH ₂ OCH ₂) ₂ CH ₂ OH	H	H	H	−0.847	−0.8354	−0.8175
Aq ₁₄	OH	(CH ₂ OCH ₂) ₅ CH ₂ OH	H	H	H	−0.844	−0.9156	−0.8943
Aq ₁₅ *	OH	(CH ₂ OCH ₂) ₆ CH ₂ OH	H	H	H	−0.849	−0.9524	−0.9351
Aq ₁₆	OH	CH ₂ CHCH ₂	H	H	H	−0.779	−0.753	−0.7379
Aq ₁₇	OH	CH ₂ Br	H	H	H	−0.797	−0.8123	−0.8067
Aq ₁₈	OH	CHBr ₂	H	H	H	−1.120	−1.0711	−1.1138
Aq ₁₉	OH	CH ₂ OCOCH ₂ Br	H	H	H	−0.830	−0.8024	−0.8265
Aq ₂₀	OH	OCH ₂ CHCH ₂	H	H	H	−0.752	−0.7346	−0.7234
Aq ₂₁	OH	H	OCH ₃	H	H	−0.750	−0.7098	−0.7520
Aq ₂₂ *	OH	CH ₂ CHCH ₂	OCH ₂ CHCH ₂	H	H	−0.745	−0.7429	−0.7304
Aq ₂₃	OH	H	H	OCH ₂ CHCH ₂	H	−0.785	−0.7872	−0.7694
Aq ₂₄	OH	H	H	OH	H	−0.635	−0.6875	−0.6880
Aq ₂₅	OH	H	OH	H	H	−0.757	−0.708	−0.7068
Aq ₂₆	OH	CH ₂ CHCH ₂	H	OH	H	−0.640	−0.6863	−0.6788
Aq ₂₇	OH	CHCHCH ₃	H	OH	H	−0.639	−0.6658	−0.6620
Aq ₂₈	OH	CH ₂ CHCH ₂	H	OH	CH ₂ CHCH ₂	−0.619	−0.6184	−0.6212
Aq ₂₉ *	OH	CHCHCH ₃	H	OH	CHCHCH ₃	−0.620	−0.693	−0.6806
Aq ₃₀	OCH ₂ CHCH ₂	H	H	H	H	−1.077	−1.053	−1.0566
Aq ₃₁	OCH ₂ CHCH ₂	CH ₃	H	H	H	−1.021	−1.0361	−1.0293
Aq ₃₂	OCHCHCH ₂	H	H	H	H	−1.017	−0.9959	−1.0235
Aq ₃₃	OCH ₂ CHCH ₂	H	H	H	OCH ₂ CHCH ₂	−1.031	−1.0729	−1.0348

* Compound belonged to test set.

Table 2

The statistical parameters of the selection of descriptors.

The number of descriptors	The squared correlation coefficient (R^2)	The value of F	Root mean square error (RMS)
3	0.8923	66.27	0.0500
4	0.9446	98.04	0.0360
5	0.9641	118.06	0.0300

prediction results are obtained, the statistical parameters are $R^2 = 0.863$, $F = 18.90$, and $RMS = 0.0600$. Fig. 3 shows the predicted vs. observed values for all of the 33 compounds studied.

3.2. Results of the RBFNN

Using the four selected descriptors, an RBFNN non-linear model was established. To obtain better results, the parameters that influence the performance of RBFNN were optimized. The selection of the optimal width value for RBFNN was performed by systemically changing its value in the training step. The value that gives the best leave one out (LOO) cross-validation result was used in the model. Based on the above optimization, the value of the optimal width is 3.5 and the corresponding number of centers (hidden layer nodes) of RBFNN is 12. The predicted data of the non-linear models are also shown in Table 1 and the plot of predicted and experimental values of RBFNN is shown in Fig. 4. The R^2 , F , and RMS for the training set are 0.960, 610.336 and 0.0299, respectively. The statistical parameters of the test set are $R^2 = 0.863$; $F = 18.946$; and $RMS = 0.0476$.

Comparing the results of the HM with those of the RBFNN, it can be concluded that the performance of RBFNN is a little bit better than that of HM. And from Figs. 3 and 4, we can see that both the two methods show satisfactory predicted results, indicating that both the HM and RBFNN models give satisfactory results.

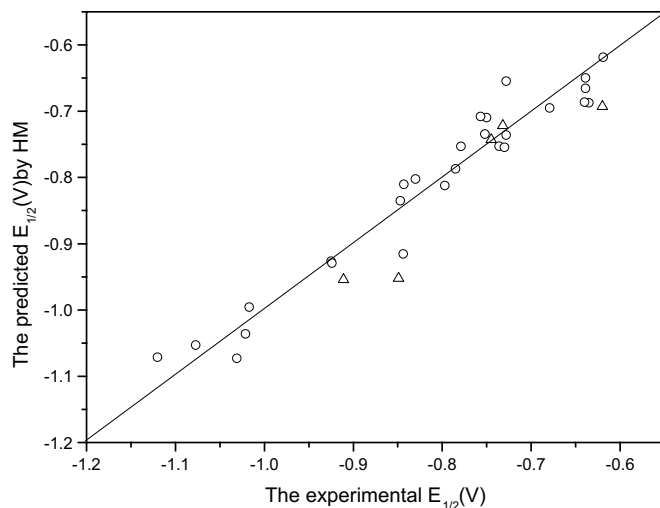
3.3. Discussions of the input parameters

By interpreting the descriptors in the model, it is possible to gain some insight into factors affecting the half-wave potential value and find out which structural factor plays an important role during the reduction reaction. In the linear model, four descriptors were found to be important for these compounds studied.

Both PNSA-1 and RNCS belong to electrostatic descriptors. They are kinds of charged partial surface area descriptors, which reflect characteristics of the charge distribution of the molecules. They are calculated in terms of the whole surface area of the molecule and functional group portions. From Table 3 we can see both the descriptors encode features responsible for redox process of molecules. The negative sign of the coefficient of PNSA-1 indicates that the 9, 10-anthraquinones with larger negative partial surface area will reduce at more negative potentials. While the positive sign of the coefficient of RNCS means that increasing the relative negative charged surface area will result in a more positive potential needed for the molecular to be reduced. Therefore we can conclude that

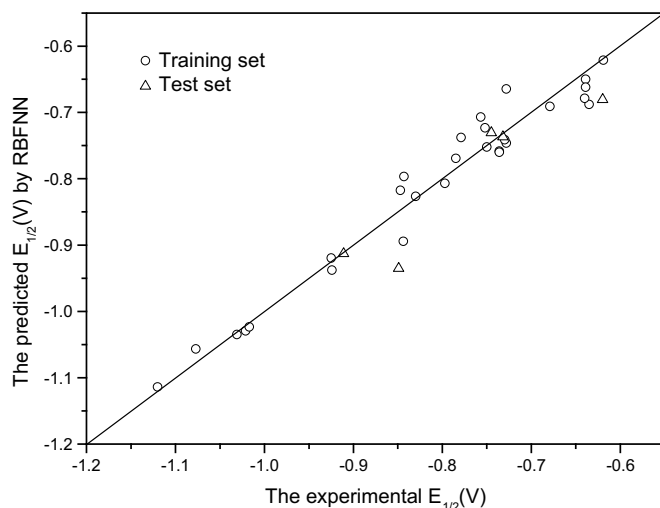
Table 3Descriptors, coefficients, standard error, and t -test values for the HM model.

Descriptors	Coefficients	Standard error	t -test
Intercept	8.0334e+01	4.5612e+00	17.6124
Min valency of a C atom	-1.9703e+01	1.1478e+00	-17.1662
PNSA-1 partial negative surface area [Zefirov's PC]	-2.7541e-03	2.9164e-04	-9.4435
RNCS Relative negative charged surface area (SAMNEG*RNCS) [Zefirov's PC]	1.9379e-02	2.6227e-03	7.3891
Min e-e repulsion for a C-O bond	-2.5378e-03	5.445e-03	-4.6604

**Fig. 3.** Plot of observed vs predicted Aq derivatives by HM.

there is a competition between PNSA-1 and RNCS, which influences the redox process and $E_{1/2}$ values of compounds rightabout.

The other two descriptors, MVCA and MEERCOB are quantum-chemical descriptors. MVCA is a valency-related descriptor, which relates to the strength of intramolecular bonding interactions and characterizes the stability of the molecules, their conformational flexibility and other valency-related properties. Carbon is the major element of organic compounds, a molecule with higher value of MVCA has increased tendency to be reduced, which results in a decrease of $E_{1/2}$ value. Hence the coefficient of this descriptor has negative sign. MEERCOB is a kind of quantum mechanical energy-related descriptors. It characterizes the total energy of the molecule in different energy scales and the intramolecular energy distribution. Meanwhile it reflects minimum values of the electron–electron repulsion energy for a given atomic species in the molecule. The electron–electron repulsion energy describes the electron repulsion driven processes in the molecule, and may be related to the conformational (rotational, inversional) changes or atomic reactivity in the molecule. As can be seen from the MLR model, MEERCOB also has a negative influence on $E_{1/2}$, indicating that the molecule with higher minimum e-e repulsion energy for a C–O bond has higher tendency to obtain an electron, and results a decrease in $E_{1/2}$.

**Fig. 4.** Plot of observed vs predicted Aq derivatives by RBFNN.

4. Conclusions

In this work, non-linear model has been compared with the linear model. The results demonstrated that both linear and non-linear QSPR models can be used for the prediction of $E_{1/2}$ of Aq derivatives based on descriptors. The high R^2 , low RMS obtained from the models suggest both of the models have good predictability. Therefore the HM method and RBFNN method can be employed to predict the $E_{1/2}$ with satisfactory results independently. The QSPR models are useful because they rationalize experimental observations and save time and money.

Acknowledgements

We are grateful for financial support from the postdoctoral foundation of Yantai University x(HY03B12) and Shandong Provincial International Cooperation Project for Excellent Teachers in Chinese Universities.

References

- [1] Shamsipur M, Sirouejinejad A, Hemmateenejad B, Abaspour A, Sharghi H, Alizadeh K, et al. Cyclic voltammetric, computational, and quantitative structure–electrochemistry relationship studies of the reduction of several 9,10-anthraquinone derivatives. *J Electroanal Chem* 2007;600:345–58.
- [2] McNaughton RL, Tipton AA, Rubie ND, Conry RR, Kirk ML. Electronic structure studies of oxomolybdenum tetrathiolate complexes: origin of reduction potential differences and relationship to cysteine–molybdenum bonding in sulfite oxidase. *Inorg Chem* 2000;39(25):5697–706.
- [3] Niu S, Wang XB, Nichols JA, Wang LS, Ichiye T. Combined quantum chemistry and photoelectron spectroscopy study of the electronic structure and reduction potentials of rubredoxin redox site analogues. *J Phys Chem A* 2003;107(16):2898–907.
- [4] Fatemi MH, Hadjmohammadi MR, Kamel K, Biparva P. Quantitative structure–property relationship prediction of the half-wave potential for substituted nitrobenzenes in five nonaqueous solvents. *Bull Chem Soc Jpn* 2007;80(2):303.
- [5] Hemmateenejad B, Yazdani M. QSPR models for half-wave reduction potential of steroids: a comparative study between feature selection and feature extraction from subsets of or entire set of descriptors. *Anal Chim Acta* 2009;634:27–35.
- [6] Nesmerak K, Nemece I, Sticha M, Waisser K, Palat K. Quantitative structure–property relationships of new benzoxazines and their electrooxidation as a model of metabolic degradation. *Electrochim Acta* 2005;50(6):1431.
- [7] Hemmateenejad B, Shamsipur M. Quantitative structure–electrochemistry relationship study of some organic compounds using PC-ANN and PCR. *Internet Electron J Mol Des*, http://biochempress.com/Files/IECMD_2003/IECMD_2003_051.pdf 2003.
- [8] Katritzky AR, Lobanov VS, Karelson M. CODESSA: reference manual. Gainesville, FL: University of Florida; 1994.
- [9] HyperChem 6.01. Hypercube, Inc.; 2000.
- [10] Stewart JPP. MOPAC 6.0, quantum chemistry program exchange, QCPE, No. 455. Bloomington, IN: Indiana University; 1989.
- [11] Katritzky AR, Lobanov VS, Karelson M. CODESSA: training manual. Gainesville, FL: University of Florida; 1995.
- [12] Xia BB, Ma WP, Zheng B, Zhang XY, Fan BT. Quantitative structure–activity relationship studies of a series of non-benzodiazepine structural ligands binding to benzodiazepine receptor. *Eur J Med Chem* 2008;43(7):1491.
- [13] Yao XJ, Panaye A, Doucet JP, Zhang RS, Chen HF, Liu MC, et al. Comparative study of QSAR/QSPR correlations using support vector machines, radial basis function neural networks, and multiple linear regression. *J Chem Inform Comput Sci* 2004;44(4):1257–66.
- [14] Luan F, Zhang XY, Zhang HX, Zhang RS, Liu MC, Hu ZD, et al. QSPR study of permeability coefficients through low-density polyethylene based on radial basis function neural networks and the heuristic method. *Comput Mater Sci* 2006;37(4):454–61.