



FEATURE EXTRACTION USING LINEAR AND NONLINEAR QSAR STUDY ON SEVERAL TAXOL DERIVATIVES AS ANTICANCER DRUGS

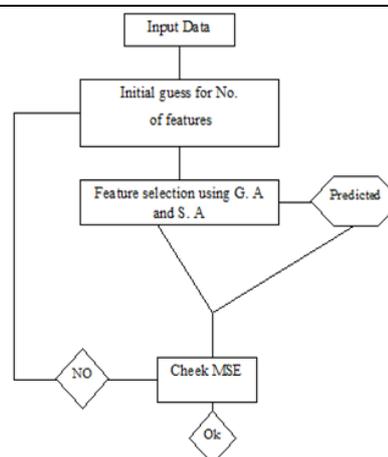
Robabeh SAYYADIKORDABADI^{a*} and Asghar ALIZADEHDAKHEL^b

^a Department of Chemistry, Rasht Branch, Islamic Azad University, Rasht, Iran

^b Department of Chemical Engineering, Rasht Branch, Islamic Azad University, Rasht, Iran

Received August 4, 2017

The activity of the Paclitaxel derivatives was estimated using multiple linear regression (MLR), artificial neural network (ANN) as modelling tools, and genetic algorithm (GA) and simulated annealing algorithm (SA) as optimization techniques. These models were employed to choose the best set of descriptors in a cross-validation procedure for non-linear-log (IC_{50}) (the empirical negative logarithm half maximal inhibitory concentration) prediction. A high predictive ability was observed for the MLR-MLR, GA-MLR, GA-ANN models, with root mean sum square errors (RMSE) of 0.421, 0.0712, 0.160, 0.0534 in gas phase and 0.910, 0.965, 0.922, 0.976 in solvent, respectively. The results obtained using the GA-ANN method indicated that the activity of the derivatives of Paclitaxel depends on different parameters such as E2u, BELP1, HATS6p, piPC05, Mor14u, BELv8, RDF120m, RDF025p descriptors in gas phase including BEHe8, Mor07u, H5u, Eig11r in the solvent phase.



INTRODUCTION

Paclitaxel sold under the brand name Taxol among others is an anti-cancer (“antineoplastic” or “cytotoxic”) chemotherapy drug for treating ovarian, breast, lung, pancreatic and other cancers.¹ It is classified as a “plant alkaloid”, a “taxane” and an “antimicrotubule agent” and its mechanism of action involves interference with the normal breakdown of microtubules during cell division.^{2,3}

QSAR approaches are mathematical equations relating chemical structure to their biological activity and present information that is useful for drug design and medicinal chemistry.⁴⁻⁶ A major step in constructing the QSAR method is to find a set of molecular descriptors representing the higher impact on the biological activity of interest.⁷⁻¹⁰

In the current study, multiple linear regressions (MLR), and artificial neural networks (ANN) as linear and nonlinear modeling tools and simulated annealing (SA) and genetic algorithm (GA)²¹⁻²⁴ as optimization model were applied to investigate the QSAR in Paclitaxel derivatives. Various QSAR models have been utilized to select the best descriptors for the important prediction of inhibitory activity of paclitaxel derivatives, and then these models were compared.

COMPUTATIONAL METHODS

The geometric optimizations of Paclitaxel derivatives were performed using Gaussian 03W¹¹ at B3lyp/6-31g. Polarized continuum model (PCM)

* Corresponding author: sayyadi@iaurasht.ac.ir

was applied to consider the non-specific solvent effect, and all Paclitaxel derivatives were optimized in H₂O solvent. Three thousand, two hundred and twenty six (3226) molecular descriptors in topological, geometrical, MoRSE,^{12,13} RDF,^{13,14} GETAWAY,^{15,16} auto-correlations¹⁷ and WHIM^{18,19} groups were calculated employing the Dragon program²⁰ and then in three steps, the number of descriptors was decreased via an objective feature selection.

Initially, in the dataset of Paclitaxel derivatives, the descriptors that had the same value of at least 70% were removed and, thereafter, the descriptors with correlation coefficient less than 0.25 with the dependent variable (-log IC₅₀) were considered redundant and removed.²¹ After these two steps, the number of descriptors was decreased to 1 047 and 1 110 in the gas and the solvent phase respectively, and then stepwise multiple linear regression procedure was utilized for rejection of descriptors. The QSAR approach with high correlation coefficient (R), low standard deviation, least numbers of independent variables, high ability to predict and high F statistic value is an ideal method.²²

In SA-ANN and GA-ANN models, 1047 and 1 110 descriptors in the gas and solvent phase were considered as possible input of the ANN and fed into the input layer of the ANNs (Fig. 1). The neural networks utilized in this work were all three-layer feed-forward network and they were

trained using the TSET members with Levenberg-Marquart algorithm.⁷ Modelling and optimization calculations were carried out using Matlab. 7.12.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - y_o)^2}{n}} \quad (1)$$

where y_i is the desired output, y_o is the predicted value by method, and n is the number of molecules in this study's data set.

RESULTS AND DISCUSSION

Geometry of the Thirty four different Paclitaxel derivatives was optimized using Gaussian 09W at B3LYP/6-31 g. All the optimized Paclitaxel compounds are illustrated in Fig. 2. All studied Paclitaxel compounds with the calculated fundamental vibration values are shown in Fig. 2. For the whole compounds, it was obtained that the values of NImag are zero and the values of the fundamental vibrations are positive. Therefore, all of the considered compounds are stable.

The best selected descriptors utilizing MLR-MLR, SA-ANN, MLR-GA and GA-ANN methods in the gas and solvent phases are shown in Tables 1 to 4.

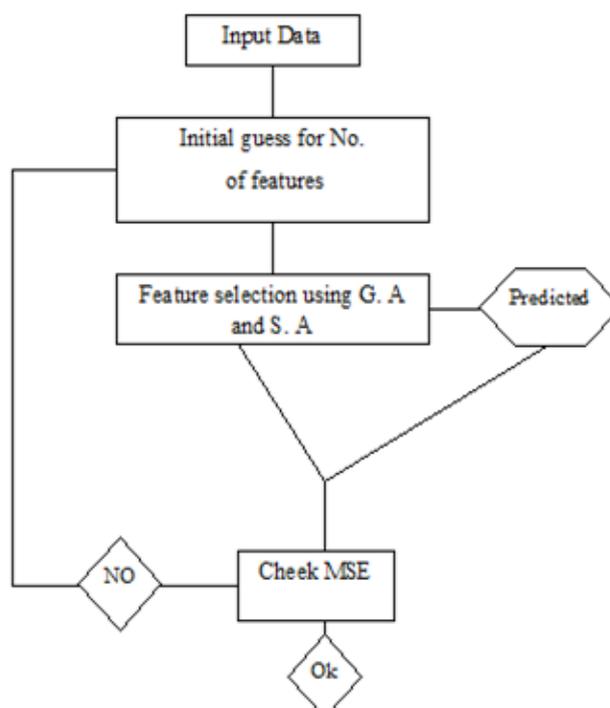


Fig. 1 –The employed procedure for finding optimum descriptors of the nonlinear approaches.

The mean square error is defined as follows:

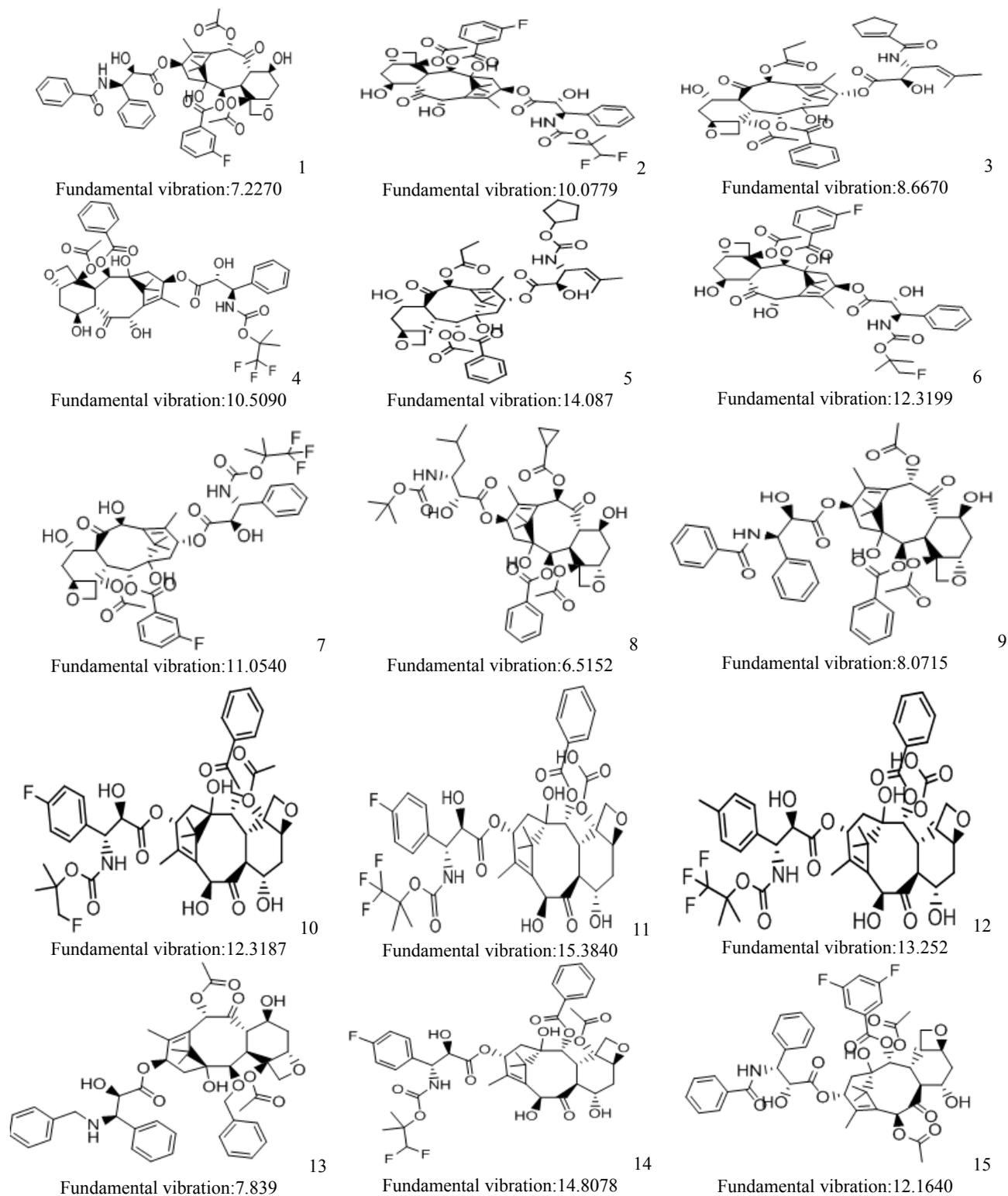


Fig. 2 – Optimized structure of the Paclitaxel derivatives used to build QSAR methods with B3lyp/6-31g in gas phase.

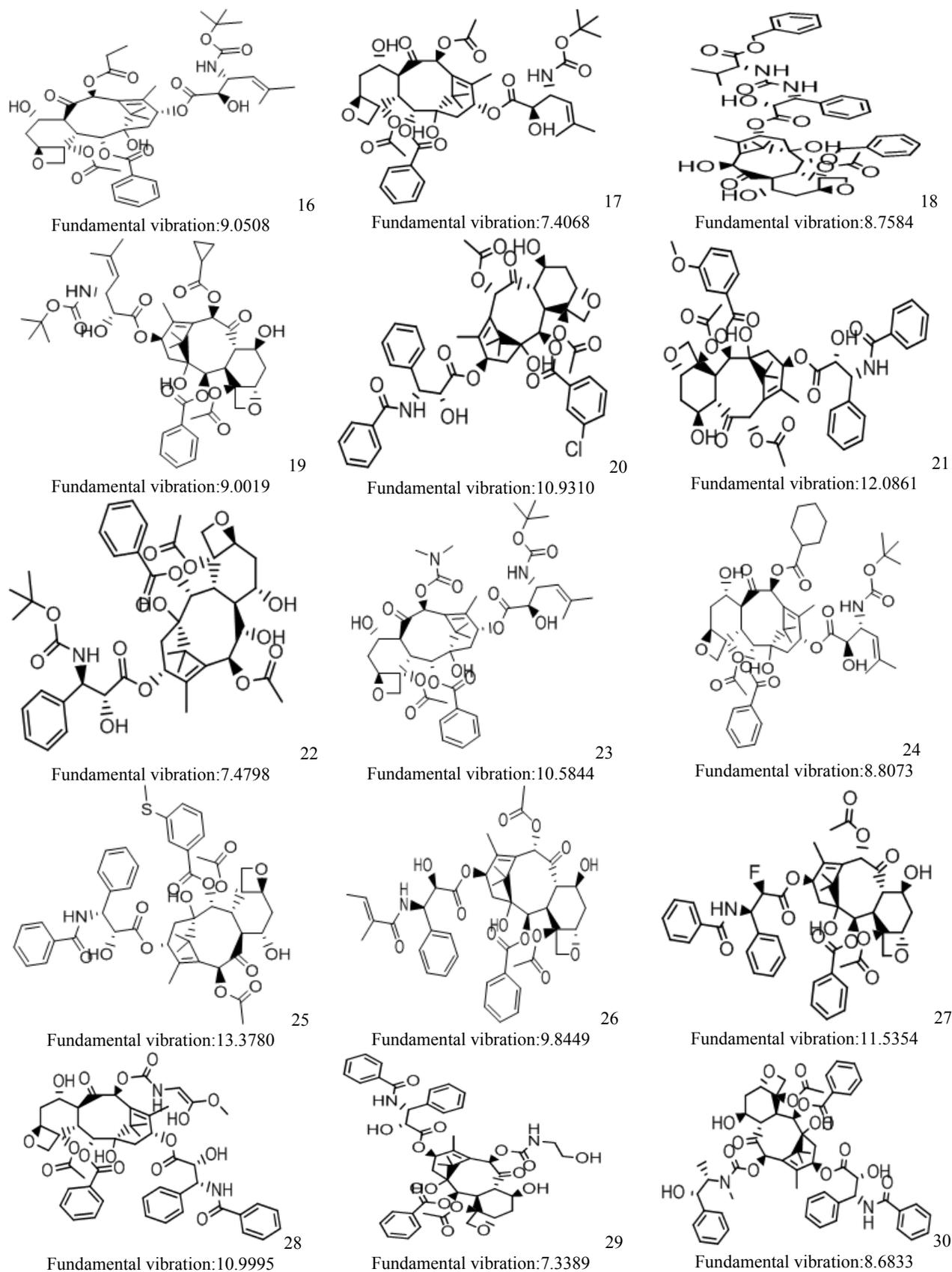


Fig. 2 – Optimized structure of the Paclitaxel derivatives used to build QSAR methods with B3lyp/6-31g in gas phase.

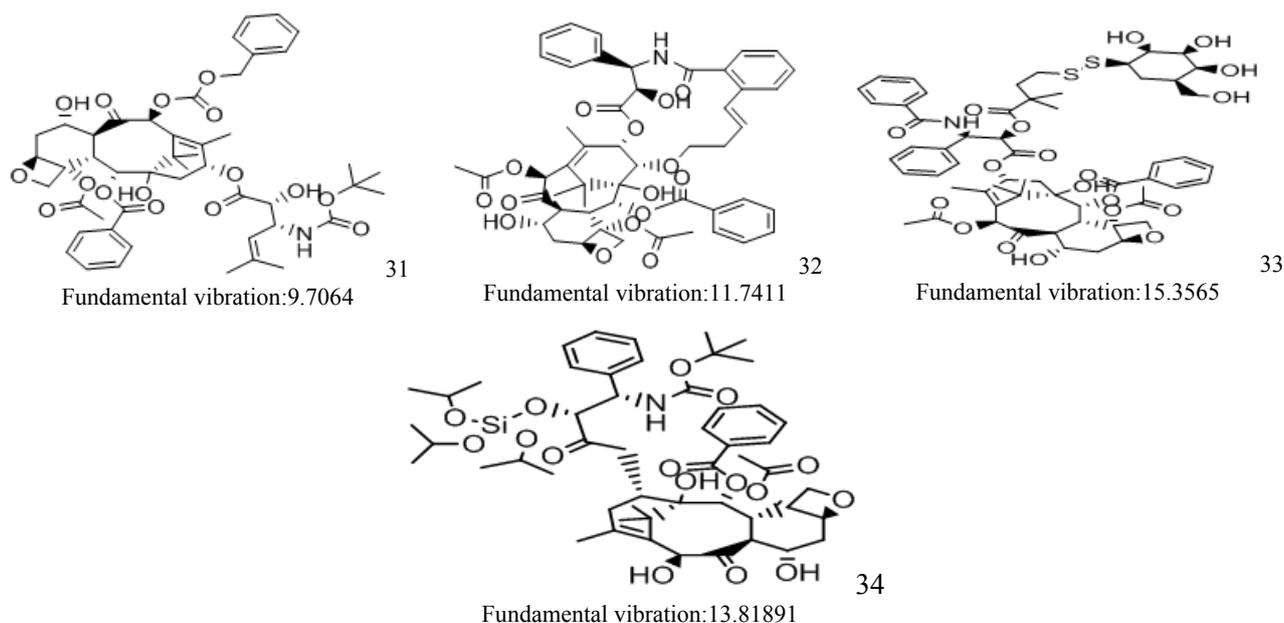


Fig. 2 – Optimized structure of the Paclitaxel derivatives used to build QSAR methods with B3lyp/6-31g in gas phase.

Table 1

Definition of the selected descriptors using MLR-MLR Method in gas and solvent phase

Descriptor	Definition	Type
IVDE	Mean in formation content on the vertex degree equality	Information indices
EEig04d	Eigenvalue 04 from edge adj.matrix weighted by dipole moments	Edge adjacency indices
C-026	R—CX—R	Atom-centred fragments
F08[O-O]	Frequency of O-O at topological distance 08	2D frequency fingerprints
MATS6m	Moran autocorrelation-log 6/ weighted by atomic masses	2D-auto correlations
Mor32e	3D-MoRSE signal 32/ weighted by atomic Sanderson electronegativities	3D-MoRSE descriptors
SIC3	Structural in formation content(neighborhood symmetry of 3-order)	Information indices
E2e	2nd component accessibility directional WHIM index/weighted by atomic Sanderson electronegativities	WHIM
Mor07u (Solvent)	3D-MoRSE signal 07/ unweighted	3D-MoRSE descriptors
Mor28m(Solvent)	3D-MoRSE signal 28/ weighted by atomic masses	3D-MoRSE descriptors
Mor26p(Solvent)	3D-MoRSE signal 28/ weighted by atomic masses	3D-MoRSE
L3m(Solvent)	3rd component size directional WHIM in dex/weighted by atomic masses	WHIM

Table 2

Definition of the selected descriptors using SA-ANN Method in gas and solvent phase

Descriptor	Definition	Type
H6u	H autocorrelation of log 6/unweighted	GETAWAY
HATS2e	Leverage-weighted autocorrelation of log 2/ weighted by atomic Sanderson electronegativities	GETAWAY
L3p	3 rd component size directional WHIM index/weighted by atomic polarizabilities	WHIM descriptors
piPC03	Molecular multiple path count of order 03	Walk and path counts
ATS6m	Broto –Moreau autocorrelation of a topological structure-log 6/weighted by atomic masses	2D autocorrelations
ATS1p	Broto –Moreau autocorrelation of a topological structure-log 1/weighted by atomic polarizabilities	2D autocorrelations
RDF085p	Radial distribution function-8.5/weighted by atomic polarizabilities	RDF descriptors
ESpm12u	Spectral moment 12 from edge adj.matrix	Edge adjacency indices

Table 2 (continued)

Mor26p(Solvent)	3D-MoRSE signal 26/ weighted by atomic polarizabilities	3D-MoRSE descriptors
ESpm11u(Solvent)	Spectral moment 11 from edge adj.matrix	Edge adjacency indices
BEHm1(Solvent)	highest eigenvalue n.1 of Burden matrix/weighted by atomic masses	Burden eigenvalues
BELP4(Solvent)	lowest eigenvalue n.4 of Burden matrix/weighted by atomic polarizabilities	Burden eigenvalues

Table 3

Definition of the selected descriptors using MLR-GA Method in gas and solvent phase

Descriptor	Definition	Type
SP19	Shape profile no.19	Randic molecular profiles
RDF050v	Radial distribution function-5.0/weighted by atomic van der Waals volumes	RDF descriptors
ESpm02u	Spectral moment 02 from edge adj.matrix	Edge adjacency indices
ESpm02x	Spectral moment 02 from edge adj.matrix weaghted by Edge degrees	Edge adjacency indices
DP16	Molecular profile no.16	Randic molecular profiles
Mor02v	3D-MoRSE signal 02/ weaghted by atomic van der waals volumes	3D-MoRSE descriptors
EEig04d	Eigenvalue 04 from edge adj.matrix weighted by dipole moments	Edge adjacency indices
ESpm06x	Spectral moment 06 from edge adj.matrix weighted by Edge degrees	Edge adjacency indices
GGL5(Solvent)	Topological charge index of order 5	Topological charge indices
MATS7m(Solvent)	Moran autocorrelation-log 7/ weighted by atomic masses	2D autocorrelations
Mor28m(Solvent)	3D-MoRSE signal 28/ weighted by atomic masses	3D-MoRSE descriptors
SP15(Solvent)	Shape profile no.15	Randic molecular profiles

Table 4

Definition of the selected descriptors using GA-ANN Method in gas and solvent phase

Descriptor	Definition	Type
E2u	2 nd component accessibility directional WHIM index/unweighted	WHIM descriptors
BELP1	lowest eigenvalue n.1 of Burden matrix/weighted by atomic polarizabilities	Burden eigenvalues
HATS6P	Leverage- weighted autocorrelation of log 6/ weighted by atomic polarizabilities	GETAWAY descriptors
piPC05	Molecular multiple path count of order 05	Walk and path counts
Mor14u	3D-MoRSE signal 14/ unweighted	3D-MoRSE descriptors
BELv8	lowest eigenvalue n.8 of Burden matrix/weighted by atomic vander Waals volumes	Burden eigenvalues
RDF120m	Radial distribution function-12.0/weighted by atomic masses	RDF descriptors
RDF025p	Radial distribution function-2.5/weighted by atomic polarizabilities	RDF descriptors
BEHe8(Solvent)	highest eigenvalue n.8 of Burden matrix/weighted by atomic Sanderson electronegativities	Burden eigenvalues
Mor07u(Solvent)	3D-MoRSE signal 07/un weighted	3D-MoRSE descriptors
H5u(Solvent)	H autocorrelation of log 5/unweighted	GETAWAY descriptors
EEig11r(Solvent)	Eigenvalue 11 from edge adj.matrix weighted by resonance integrals	Edge adjacency indices

Table 5

Statistical parameters of different linear QSAR models in gas and solvent phase

QSAR Models	Predict Set	
	R ²	RMSE
MLR-MLR	0.910	0.421
MLR-PCR	0.793	0.640
MLR-PLS1	0.895	0.456
MLR-PLS1(solvent)	0.6957	0.7773
MLR-PCR(solvent)	0.359	1.127
MLR-MLR(solvent)	0.6958	0.7771

In MLR-PCR, MLR-PLS1 and MLR-MLR models, the best descriptors were selected using MLR procedure of SPSS software in three steps described in theory and computational methods section. Thereafter, the selected descriptors were employed as input in unscramble software and

statistical parameters were calculated using PCR, PLS1 and MLR models (Table 5).

IVDE and SIC3 (Table 1) descriptors are information indices. The total information content (I) is obtained by multiplying the mean information content by the number of elements.¹²

EEig04d (Table 1), ESPm12u (Table 2), ESPm11u (Table 2), EEig11r (Table 4), ESpm02u, ESpm02x, EEig04d, ESpm06x (Table 3) descriptors are adjacency indices. The Edge adjacency relationships in molecular graphs have been utilized to define a new topographic index. Molecules as weighted graphs were employed for the calculation of the novel index, in which the elements of edges set were substituted by the bond orders between connected atoms in the molecule.²³

C-026 and F08 [O-O] (Table 1) descriptors are Atom-centered fragment and 2D-Frequency fingerprint, respectively. Fragment descriptors are representations of local atomic environments.¹³

MATs6m (Table 1), ATS6m, ATS1p (Table 2), MATS7m (Table 3) are 2D-autocorrelation descriptors that represent the topological structure of the compounds in nature are more complex than the classical topological descriptors.²⁴ Mor32e (Table 1), Mor07u, Mor28m (Table 1), Mor26p (Tables 1, 2), Mor 02v (Table 3), Mor28m (Table 3), Mor14u, Mor07u (Table 4) are 3D-MoRSE descriptors. The presence of a MoRSE descriptor indicates that the size of the inhibitor compound has certain effect on the extent of the interaction between the enzyme and compound.²³

E2e (Table 1), L3m (Table 2), E2u (Table 3) are WHIM descriptors that were built in such a way to capture the relevant molecular 3D information regarding molecular size, shape, and symmetry and atom distribution with respect to invariant reference frames.¹³

piPC03 (Table 2), piPC05 (Table 4) are Walk and path descriptors. The molecular multiple path counts (piPck) are defined as path counts weighted by the bond order.¹³

H6u, HATS2e (Table 2), HATS6P (Table 4), H5u (Table 4) are GETAWAY descriptors. GETAWAY (Geometry, Topology, and Atom-Weights Assembly) descriptor representations encode the geometrical information obtained from the molecular matrix, the topological information obtained from the molecular graph and the information obtained from atomic weights, which are specially designed with the aim of matching the 3D-molecular geometry.¹³

SP19, DP16, SP15 (Table 3) are Randic molecular profile descriptors. The Randic molecular profile DP_k , is derived from the distance distribution moments of the geometric matrix G as the average row sum of its entries raised to the k^{th} power and normalized by the factor $k!$.¹³

GGL5 (Table 3) descriptor is Topological charge indices that was proposed to evaluate the charge

transfer between pairs of atoms and, therefore, the global charge transfer in the molecule.¹³

RDF085p (Table 2), RDF050v (Table 3), RDF120m (Table 4), RDF025p (Table 4) are RDF descriptors. RDF descriptors are independent of the number of atoms, *i.e.*, the size of a molecule, it is unique regarding the three-dimensional arrangement of the atoms, and it is invariant against translational and rotational entire molecule.¹³

BEHm1, BLP4 (Table 2), BELP1, BELv8 (Table 4), BEHe8 (Table 4) are Burden eigenvalues descriptors. The B matrix has been defined as the number of atoms and bond order between two atoms or the electronegativity of the atoms.¹³

When the MLR-MLR model was utilized, the RMSE of the predicted activity was found to be 0.421 in the gas phase and 0.7771 in the solvent phase. In addition, the correlation coefficient (R^2) calculated for the PSET was 0.910 in the gas phase and 0.6958 in the solvent phase. It was demonstrated that MLR-MLR method is better than other linear methods (MLR-PLS1 and MLR-PCR, Table 5).

To establish the SA-ANN, MLR-GA and GA-ANN methods, the 1047 and 1110 descriptors in the gas and solvent phases were fed to the neural network to select the best descriptors.

In GA-ANN, SA-ANN, MLR-GA models, 80, 10 and 10% of data sets were randomly chosen as training, validation and test sets, respectively. Table 6 shows test and valid series in QSAR models in Paclitaxel compounds.

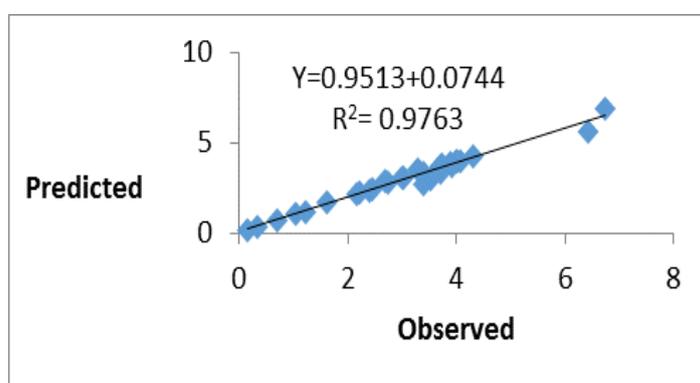
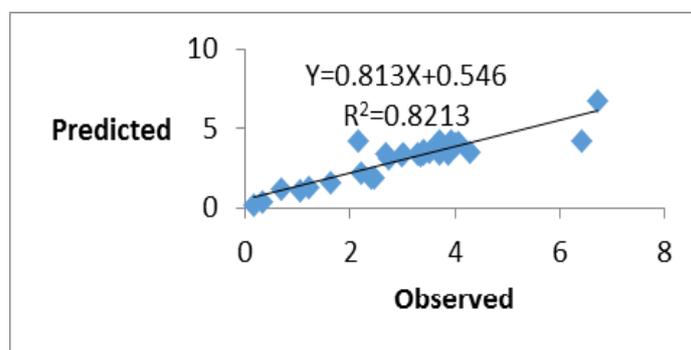
In the GA-ANN method, the RMSE and R-square were calculated as 0.0534 and 0.976 in the gas phase and 0.355 and 0.821 in the solvent phase for the predicted activity of Paclitaxel derivatives, respectively. Therefore, GA-ANN method was better than the other models and as such, only the descriptors utilized in this method were evaluated in this study as shown in Table 4. The plot showing the variation of observed *versus* predicted $-\log IC_{50}$ (empirical negative logarithm of half maximal inhibitory concentration) values are illustrated in Figs. 3 and 4.

High RMSE errors resulting from the method were due to possible errors in experimental data employed in this study and because RMSE is highly dependent on the range of the dependent variable.²⁴ The values of $-\log IC_{50}$ (experimental negative logarithm half maximal inhibitory concentration) in our dataset were in the range of 0.15 to 6.73. The RMSE of the predicted set in GA-ANN model were 0.0534 and 0.355 in the gas and solvent phases, respectively, which are acceptable in comparison with previous works.²³⁻²⁷

Table 6

Statistical parameters of different non-linear QSAR models in gas and solvent phase

QSAR models	Predicted		Train	
	R ²	RMSE	R ²	RMSE
GA-ANN(Gas)	0.976	0.0534	0.975	0.0574
GA-MLR(Gas)	0.922	0.160	0.940	0.159
SA-ANN(Gas)	0.965	0.0712	0.964	0.0527
GA-ANN(Solvent)	0.821	0.355	0.821	0.411
GA-MLR(Solvent)	0.773	0.462	0.793	0.478
SA-ANN(Solvent)	0.756	0.498	0.801	0.408

Fig. 3 – Plot between observed vs. predicted –log(IC₅₀) by using GA-ANN descriptors in gas phase.Fig. 4 – Plot between observed vs. predicted –log(IC₅₀) by using GA-ANN descriptors in solvent phase.

The graph of E2u, BELP1, HATS6p, PiPC05, Mor14u, BELv8, RDF120m, RDF025p descriptors in the gas phase *versus* empirical negative logarithm half maximal inhibitory concentration (–logIC₅₀) were plotted using Matlab program (Fig. 5).

The charts of the gas phase demonstrated that with increase in E2u up to 0.35, HATS6p, PiPC05,

BELv8, RDF120m descriptors, the response (–logIC₅₀) was reduced.

As the Mor14u (Factor 5) descriptor increased from 7.5 to 8, no change in response was observed. Thus during this period, a bar was seen in the response. In addition, as E2u of 0.35 BELp1, RDF025p (Factors 1, 2, 8) descriptors increased, the response was increased as well.

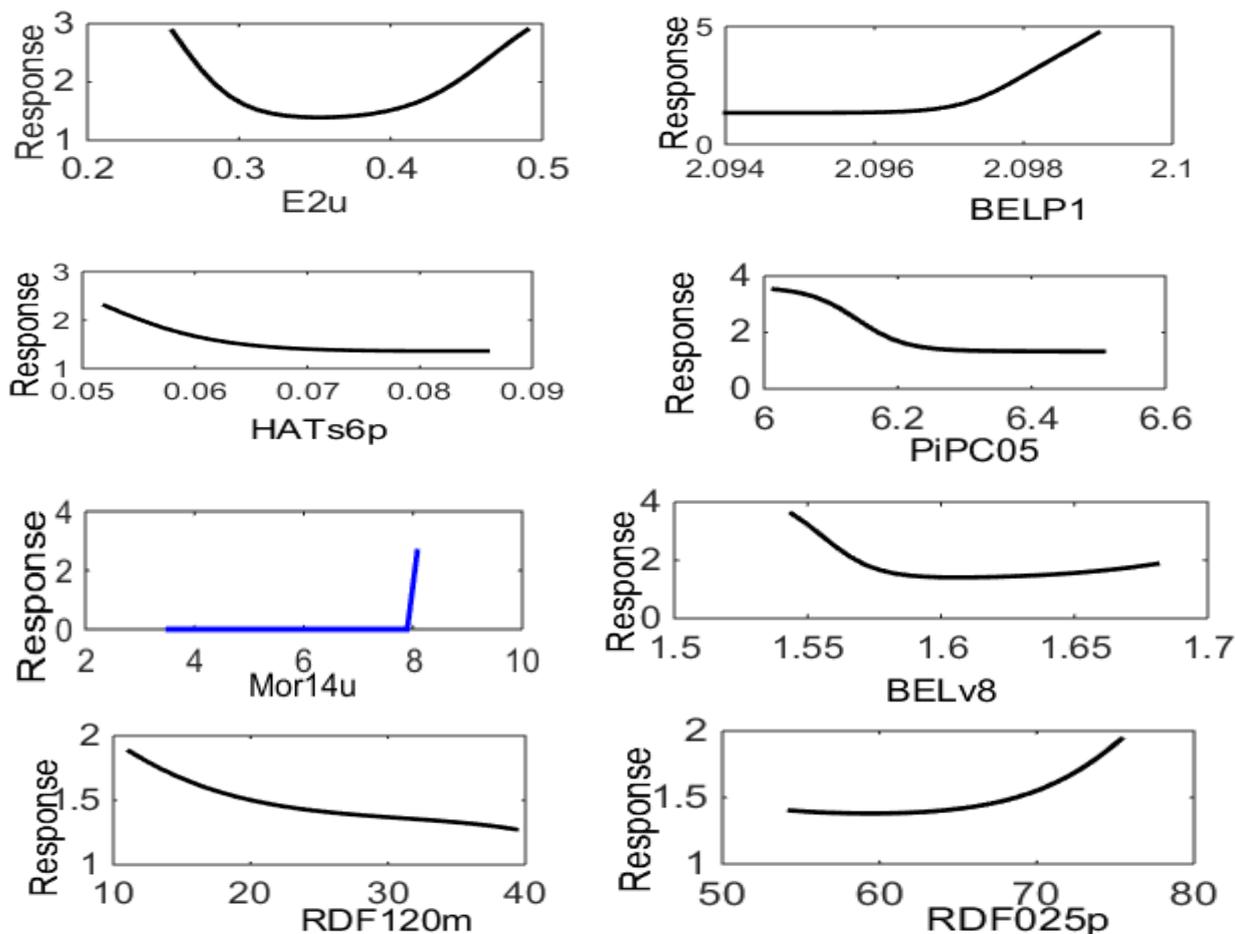


Fig. 5 – Plot between $-\log IC_{50}$ experimental (Response) versus of the E2u, BELP1, HATS6p, PiPC05, Mor14u, BELv8, RDF120m, RDF025p descriptors in gas phase.

CONCLUSION

Among the QSAR approaches employed in this study, the nonlinear feature selection methods were shown to be better than their linear methods, and the results of GA-ANN model were better than the other non-linear models used. These results also proved that E2u, BELP1, HATS6p, piPC05, Mor14u, BELv8, RDF120m, RDF025p descriptors in the gas phase and BEHe8, Mor07u, H5u, Eig11r descriptors in the solvent phase were more significant than other descriptors in building this QSAR model and predicting the biological activity of Paclitaxel substitution patterns.

Acknowledgement. The support provided by Rasht branch, Islamic Azad University is gratefully acknowledged.

REFERENCES

- M. W. Saville, J. Lietzau, J. M. Pluda, W. H. Wilson, R. W. Humphrey, E. Feigel, S. M. Steinberg and S. Broder, *The Lancet*, **1995**, 346, 26.
- B. Rajnish and Y. Hongtao, *Oncogene*, **2004**, 23, 2016.
- D. A. Brito, Z. Yang and C. L. Rieder, *J. Cell Biology*, **2008**, 182, 623.
- D. Hadjipavlou-Litina, *Med. Res. Rev.*, **1998**, 18, 91.
- P. Gramatica and E. Papa, *QSAR. Comb. Sci.*, **2003**, 22, 374.
- C. Hansch, A. Kurup, R. Garg and H. Gao, *Chem. Rev.*, **2001**, 101, 619.
- D. Horvath and B. Mao, *QSAR Comb. Sci.*, **2003**, 22, 49.
- S. Putta, J. Eksterowicz, C. Lemmen and R. Stanton, *J. Chem. Inf. Comput. Sci.*, **2000**, 43, 1623.
- S. Gupta, M. Singh and A. K. Madan, *J. Chem. Inf. Comput. Sci.*, **1999**, 39, 272.
- V. Consonni, R. Todeschini and M. Pavan, *J. Chem. Inf. Comput. Sci.*, **2000**, 42, 693.
- E. Borges de Melo and M. M. C. Ferreira, *Eur. J. Med. Chem.*, **2009**, 44, 3577.
- J. H. Schuur, P. Selzer and J. Gasteiger, *J. Chem. Inf. Comput. Sci.*, **1996**, 36, 334.
- R. Todeschini and V. Consonni, "Hand Book of Molecular Descriptors", Wiley-VCH Verlag GmbH, Weinheim, Germany, 2008.
- M. C. Hemmer, V. Steinhauer and J. Gasteiger, *Vibr. Spectrosc.*, **1999**, 19, 151.
- V. Consonni, R. Todeschini and M. Pavan, *J. Inf. Comput. Sci.*, **2002**, 42, 682.

16. V. Consonni, R. Todeschini and M. Pavan, *J. Inf. Comput. Sci.*, **2002**, *42*, 693.
17. P. Gramatica, V. Consonni and R. Todeschini, *Chemosphere*, **1999**, *38*, 371.
18. P. Gramatica, V. Consonni, R. Todeschini, *Chemosphere*, **2000**, *4*, 763.
19. M. H. Fatemi and S. Gharaghani, *Med. Chem.*, **2007**, *15*, 7746.
20. R. Todeschini, *Milano Chemometrics and QSAR Research Group. Molecular descriptors. An introduction.* <http://www.disat.unimib.it/chem> (accessed **2000**).
21. M. Jalali-Heravi and F. Parastar, *J. Inf. Comput. Sci.*, **200**, *40*, 147.
22. K. Levenberg, *Quarterly of Applied Mathematics*, **1994**, *2*, 164.
23. S. H. Sadat Hayatshahi, P. Abdolmaleki, M. Ghiasi and S. Safarian, *FEBS Lett.*, **2007**, *581*, 506.
24. M. Nirouei, G. Ghasemi, P. Abdolmaleki, A. Tavakoli and S. Shariati, *Indian J. Biochem. Biophys.*, **2012**, *49*, 202.
25. D. Sisodiya and K. Dashora, *Int. J. Phytopharmacy.*, **2015**, *4*, 153.
26. G. Melagraki, A. Afantitis, H. Sarimveis, O. Igglessi-Markopoulou, *Bioorg. Med. Chem.*, **2006**, *14*, 1108.
27. D. Jaiswal, C. Karthikeyan, S. K. Shirastava and P. Trivedi, *Internet Electronic J. Molecul. Design*, **2006**, *5*, 345.