Mini Review

Open Access    CrossMark

# Review on: quantitative structure activity relationship (QSAR) modeling

## Abstract

Quantitative Structure Activity Relationship (QSAR) are mathematical models that seek to predict complicated physicochemical /biological properties of chemicals from their simpler experimental or calculated properties .QSAR enables the investigator to establishes a reliable quantitative relationship between structure and activity which will be used to derive an insilico model to predict the activity of novel molecules prior to their synthesis. The past few decades have witnessed much advances in the development of computational models for the prediction of a wide span of biological and chemical activities that are beneficial for screening promising compounds with robust properties. This review covers the concept, history of QSAR and also the components involved in the development of QSAR models.

**Keywords:** QSAR, model development, applicability domain, molecular descriptor, virtual screening

## Umma Muhammad,[1] Adamu Uzairu,[2] David Ebuka Arthur[2]

[1]Department of Pre-nd Sci & Tech. School of General Studies, Kano State Polytechnic, Nigeria
[2]Department of Chemistry, Ahmadu Bello University Zaria, Nigeria

**Correspondence:** Umma Muhammad, Department of Pre-nd Sci & Tech. School of General Studies, Kano State Polytechnic, Nigeria, Email umjidda58@gmail.com

## Introduction

Quantitative structure – activity relationship (QSAR) modeling pertains to the construction of predictive models of biological activities as a function of structural and molecular information of a compound library. The concept of QSAR has typically been used for drug discovery and development and has gained wide application for correlating molecular information with not only biological activities but also with other physicochemical properties, which has therefore been termed quantitative structure – property relationship (QSPR). QSAR is widely accepted predictive and diagnostic process used for finding associations between chemical structures and biological activity. QSAR has emerged and has evolved trying to fulfill the medicinal chemist's need and desire to predict biological response.[1] It found its way into the practice of agro chemistry, pharmaceutical chemistry, and eventually most facets of chemistry.[2]

QSAR is the final result of computational processes that start with a suitable description of molecular structure and ends with some inference, hypothesis, and predictions on the behavior of molecules in environmental, physicochemical and biological system under analysis.[3] The final outputs of QSAR computations are set of mathematical equations relating chemical structure to biological activity.[4–6] Multivariate QSAR analysis employs all the molecular descriptors from various representations of a molecule (1D, 2D and 3D representation) to compute a model, in a search for the best descriptors valid for the property in analysis. This review covers the concepts, history and the steps involved in the development of QSAR models.

## History of QSAR

Cros[2] proposed a relationship which existed between the toxicity of primary aliphatic alcohols with their water solubility.[2] In 1868 Crum-Brown and Fraser published an equation which is considerable to be the first generation formulation of a quantitative structure-activity relationship, in their investigations of different alkaloids.[7] Systematic QSAR began with the work of[8] on the narcotic activity of various drugs.[9] Hammett[10] introduced a method to account for

substituent effects on reaction mechanism.[10] Taking Hammetts model into account Taft proposed in 1956 an approach for separating polar, steric, and resonance effects of substituents in aliphatic compounds.[11] Classical approach to QSAR/QSPR was led by the pioneering works of Hansch et al.[12] in the development of linear Hansch equation.[12]

QSAR/QSPR received a big boost with the development of newer, more complex descriptors, soft ware's and computers. This has been instrumental in the application of the prediction techniques that were either not feasible or were previously too time consuming.

## QSAR methodology

QSAR methodologies have the potential of decreasing substantially the time and effort required for the discovery of new medicines.[13] A major step in constructing the QSAR models is to find a set of molecular descriptors that represents variations of the structural properties of the molecule.[14] The QSAR analysis employs statistical methods to derive quantitative mathematical relationship between chemical structure and biological activity.[15] The process of QSAR modelling can be divided into three stages: development, model validation and application.

## Development

For the development of the model the compounds gathered from literature source could be divided into training and test set. The training sets are used in model construction while the test set for external validation.

The structures of the complexes under study could be drawn in 2D ChemDraw. These could be converted into 3D objects using the default conversion procedure implemented in the CS Chem 3D ultra. The generated 3D structures of the complex were then subjected to energy minimization and geometry optimization using Spartan.[16] Molecular descriptors could be calculated using chemical software's such as Dragon,[17] Gaussian,[18] PADEL,[19] etc. Molecular descriptors can be defined as the essential information of a molecule in terms of its physicochemical properties such as constitutional, electronic, geometrical, hydrophobic, lipophilicity, solubility, steric, quantum

chemical and topological descriptors.[20] Multivariate analysis such as multi linear regression, Partial least Square etc could be carried out for correlating molecular descriptors with observed activity.

**Internal model validation**

The developed models were validated internally by leave- one- out (LOO) cross- validation technique. In this technique, one compound is eliminated from the data set at random in each cycle and the model is built using the rest of the compounds. The model thus formed is used for predicting the activity of the eliminated compound. The process is repeated until all the compounds are eliminated once. The Cross-validated squared correlation coefficient, R2cv (Q2) was calculated using the expression:

$$Q^2 = 1 - \frac{\sum\left(Y_{Obs} - Y_{Pred}\right)^2}{\sum\left(Y_{Obs} - \overline{Y}\right)^2}$$

Where $Y_{OBS}$ represents the observed activity of the training set compounds, $Y_{pred}$ is the predicted activity of the training set compounds and $Y$ corresponds to the mean observed activity of the training set compounds. Also calculated was the adjusted $R^2$ ($_{adj}R^2$) which is a modification of $R^2$ that adjust the number of explanatory terms in a model. Unlike $R^2$ in which addition of descriptors to the developed QSAR model increases its value, the value of $_{adj}R^2$ increases only if the new term improves the model more than what would be expected by chance.[21] Hence $_{adj}R^2$ overcomes the draw backs associated with the value of $R^2$ and was calculated using the expression:

$$adj\, R^2 = \frac{(n-1)R^2 - p}{n - p - 1}$$

Where p is the number of predictor variables used in the model development. In other to judge the overall significance of the regression coefficients, the variance ratio, F value (the ratio of regression mean square to deviations mean square), was also calculated using the relation:

$$F = \frac{\left(\dfrac{\sum\left(Y_{cal} - \overline{Y}\right)^2}{p}\right)}{\left(\dfrac{\sum\left(Y_{obs} - Y_{cal}\right)^2}{N - P - 1}\right)}$$

**External model validation**

External validation was employed in order to determine the predictive capacity of the developed model as judged by its application for the prediction of test set activity values and calculation of predictive $R^2$($R^2$pred) value as given by the expression:

$$R^2_{pred} = 1 - \frac{\sum\left(Y_{pred\left(Test\right)} - Y_{\left(Test\right)}\right)^2}{\sum\left(Y_{\left(Test\right)} - \overline{Y}_{\left(Training\right)}\right)^2}$$

Where $Y_{pred\left(Test\right)}$ and $Y_{\left(Test\right)}$ indicate predicted and observed activity values respectively, of the test compounds. $\overline{Y}_{\left(Training\right)}$ indicates mean activity value of the training set. $R^2_{pred}$ is the predicted correlation coefficient calculated from the predicted activity of all

the test set compounds. It has been observed that $R^2$ pred may not be sufficient to indicate the external predictability of a model since its value is controlled by $\sum\left(Y_{\left(Test\right)} - \overline{Y}_{\left(Training\right)}\right)^2$. Thus $R^2_{pred}$ depends on the training set mean and may not truly reflect the predictive capability of the developed model with regards to a new data set.[22] this may result in considerable numerical difference between the observed and predicted values in spite of maintaining a good overall intercorrelation.

**Randomization test**

The Robustness of the developed QSAR model was checked using Y-randomization technique in which model randomization was employed. In Y-randomization, validation was performed by permutating the response values, Activity (Y) with respect to the descriptor (X) matrix which was unaltered. The deviation in the values of the squared mean correlation coefficient of the randomized model ($Rr^2$) from the squared correlation coefficient of the non-random model ($R^2$) is reflected in the value of $R^2_p$ parameter computed from the expression.[23]

$$R^2_p = R^2 \times \sqrt{\left(R^2 - R^2_r\right)}$$

In an ideal case, it is observed that the average value of R2 (Rr2) for randomized models should be zero. This implies that the value of Rp2 should be equal to the value of R2 for the developed QSAR model. This led Todeschini[25] to suggest a correction for Rp2 which is defined as:

$$cR^2_p = R \times \sqrt{R^2 - R^2_r}$$

In other to penalize the developed models for the difference between the squared correlation coefficients of the randomized and the non-randomized models, the values $cR_p^2$ was calculated for each model. This procedure ensures that the model is not due to a chance. The Y-randomization results were generated using the program "MLR Y-Randomization Test 1.2".[24]

**Application**

The application of QSAR models depends on statistical significance and predictive ability of the models. The prediction of a modeled response using QSAR is valid only if the compound being predicted is within the applicability domain of the model. The applicability domain is a theoretical region of the chemical space, defined by the model descriptors and modeled response and thus by the nature of the training set molecules.[25] It is possible to check whether a new chemical lies within applicability domain using the leverage approach. A compound will be considered outside the applicability domain when the leverage values is higher than the critical value of 3p/n, where p is the number of model variables plus 1 and n is the number of objects used to develop the model. Other approach includes training set interpolation by Jawors.[26] Cluster – based approach by Stan forth et al.[27]

## Conclusion

The QSAR models are useful for various purposes including the prediction of activities of untested chemicals. It helps in the rational design of drugs by computer aided tools via molecular modeling, simulation and virtual screening of promising candidates prior to synthesis. In this review article the concept, brief history and components involved in modeling were discussed.

## Acknowledgements

None.

## Conflict of interest

The author declares that there is no conflict of interest.

## References

1. Hansch CLA. *Substituent constants for correlation analysis in chemistry and biology*. New York: John Wiley and Sons; 1979.

2. Cros A. *Action de l'alcool amylique sur l'organisme*. 1863.

3. Eriksson L, Jaworska J, Worth AP, et al. Methods for reliability and uncertainty assessment and for applicability evaluations of classification–and regression–based QSARs. *Environ Health Perspect*. 2003;111(10):1361–1375.

4. Golbraikh A, Shen M, Xiao Z, et al. Rational selection of training and test sets for the development of validated QSAR models. J *Comput Aided Mol Des*. 2003;17(2–4):241–253.

5. Hansch C, Sinclair JF, Sinclair PR. Induction of Cytochrome P450 by Barbiturates in Chick Embryo Hepatocytes: A Quantitative Structure Activity Analysis. *Quantitative Structure Activity Relationships*. 1990;9(3):223–226.

6. Wedebye EB, Dybdahl M, Nikolov NG, et al. QSAR screening of 70,983 REACH substances for genotoxic carcinogenicity, mutagenicity and developmental toxicity in the ChemScreen project. *Reprod Toxicol*. 2015;55:64–72.

7. Crum–Brown AFT. *On the connection between chemical constitution and physiological action. Pt 1. On the physiological action of the salts of the ammonium bases, derived from Strychnia, Brucia. Thebia, Codeia, Morphia, and Nicotia.* T Roy Soc Edin. 1868;25:151–203.

8. Cantor RS. Breaking the Meyer–Overton rule: predicted effects of varying stiffness and interfacial activity on the intrinsic potency of anesthetics. *Biophys J*. 2001;80(5):2284–2297.

9. Pohorille A, Wilson MA, New MH, et al. Concentrations of anesthetics across the water–membrane interface; the Meyer–Overton hypothesis revisited. *Toxicol Lett*. 1998;100:421–430.

10. Hammett LP. Some Relations between Reaction Rates and Equilibrium Constants. *Chemical Reviews*. 1935;17(1):125–136.

11. Taft R. *Separation of Polar, Steric and Resonance Effects in Reactivity in Steric Effects in Organic Chemistry*. New York; John Wiley and Sons: 1956.

12. Fujita T, Iwasa J, Hansch C. A new substituent constant, $\pi$, derived from partition coefficients. *Journal of the American Chemical Society*. 1964;86(23):5175–5180.

13. Gramatica P, Giani E, Papa E. Statistical external validation and consensus modeling: A QSPR case study for K oc prediction. *Journal of Molecular Graphics and Modelling*. 2007;25(6):755–766.

14. Gramatica P. Principles of QSAR models validation: internal and external. *QSAR and Combinatorial Science*. 2007;26(5):694.

15. Ghafourian T, Cronin MT. The impact of variable selection on the modelling of oestrogenicity. *SAR QSAR Environ Res*. 2005;16(1–2):171–190.

16. Hehre WJ, Huang WW. *Chemistry with Computation: An introduction to SPARTAN: Wavefunction*. 1995.

17. Mauri A, Consonni V, Pavan M, et al. Dragon software: An easy approach to molecular descriptor calculations. *Match*. 2006;56(2):237–248.

18. Salahub DR, Fournier R, Młynarski P, et al. Gaussian–based density functional methodology, software, and applications. *Density functional methods in chemistry*.1991;77–100.

19. Yap CW. PaDEL descriptor: An open source software to calculate molecular descriptors and fingerprints. *J Comput Chem*. 2011;32(7):1466–1474.

20. Todeschini R, Consonni V. *Molecular Descriptors for Chemoinformatics*. USA; John Wiley & Sons. 2010;41(2).

21. K Roy, I Mitra, S Kar, et al. Comparative Studies on some metrics for external validation of QSAR model. *J Chem Inf Mdel*. 2012;52(2):396–408.

22. Kar S, Roy K. QSAR of phytochemicals for the design of better drugs. *Expert Opin Drug Discov*. 2012;7(10):877–902.

23. Roy K, Kar S, Das RN. Background of QSAR and Historical Developments. In: KRKN Das, editor. *Understanding the Basics of QSAR for Applications in Pharmaceutical Sciences and Risk Assessment*. Boston: Academic Press; 2015;1–46.

24. Roy K, Kar S, Ambure P. On a simple approach for determining applicability domain of QSAR models. *Chemometrics and Intelligent Laboratory Systems*.145, 22–29.

25. Todeschini R, Consonni V, Pavan M. *Milano chemometrics and QSAR research group*. KOALA–Software for Kohonen Artificial Neural Networks. 2007.

26. Jaworska JS, Comber M, Auer C, et al. Summary of a workshop on regulatory acceptance of (Q) SARs for human health and environmental endpoints. *Environ Health Perspect*. 2003;111(10):1358–1360.

27. Stanforth RW, Kolossov E, Mirkin B. A measure of domain of applicability for QSAR modelling based on intelligent K–means clustering. *QSAR and Combinatorial Science*. 2007;26(7):837–844.