

Machine Learning - Based q-RASPR Modeling of Power Conversion Efficiency of Organic Dyes in Dye-Sensitized Solar Cells

Souvik Pore, Arkaprava Banerjee, Kunal Roy*

Drug Theoretics and Chemoinformatics Laboratory,
Department of Pharmaceutical Technology, Jadavpur University,
188 Raja S C Mullick Road, 700032, Kolkata, India

Supplementary Materials SI-1

*Corresponding author

Prof. Kunal Roy, Phone: +91 9831594140; Fax: +91-33-2837-1078;

Email: kunalroy_in@yahoo.com; kunal.roy@jadavpuruniversity.in

Machine learning methods

- a) Ridge regression: It is a popular method, which is used to address the multicollinearity problem of MLR models without removing any independent variables. In this method, a small amount of bias (penalty) is added to get better predictions. It is an important regularization technique that helps to reduce the complexity of a model, and it is known as Tikhonov regularization. The generalized equation of ridge regression is:

$$L(x, y) = \text{Min} \left(\sum_{i=1}^n (y_i - w_i x_i)^2 + \lambda \sum_{i=1}^n (w_i)^2 \right) \quad \text{where } w_i \text{ is weightage of individual feature and } \lambda \text{ is penalty term.}^{42,43}$$

- b) Linear Support Vector Machine (LSVM): LSVM is a form of SVM algorithm, where the data domain can be classified linearly without any kind of transformation. The main steps in LSVM are mapping of data domain into response data, and then division of the data domain. The generalized equation for LSVM is: $\hat{y} = w^T X + b$.⁴⁴

- c) Support Vector Machine (SVM): It is a classification machine-learning algorithm, but it can also be used for regression problems, as known as Support Vector Regression (SVR). The main idea behind SVM is to draw a decision boundary between observations to perform its predictions. For nonlinear SVM, we have to transform the data into a feature space (using a kernel function) before mapping with the response, and the generalized equation for SVM (non-linear) is represented as follows: $\hat{y} = w^T \phi(X) + b$, where \hat{y} is predictions, w is the vector of weights, X is a vector of input features, ϕ is a kernel function and b is bias. This decision boundary is known as a plane for a three-dimensional space and known as a hyperplane for higher order space where a large number of features are present. The SVM method considers both margins which are formed by the area between the decision boundary and the closest training compound and the hyperplane for predictions. The margin is

$$\text{margin} = \frac{1}{w^T w}$$

represented by the following equation: SVM tries to maximize the distance between the two closest training compounds on either side of decision boundary.⁴⁵

- d) Random forest (RF): The RF algorithm builds multiple decision tree models and combines their outcome for more accurate and stable prediction. This method helps to overcome overfitting problem of a decision tree models. The RF algorithm is based on an ensemble learning method known as Bagging (bootstrap aggregating) which is a resampling technique applied to a dataset. In bootstrapping, observations are selected by random sampling with replacement, and random feature subsets are selected. In bagging, a large number of datasets

are created by bootstrapping the original dataset, multiple decision tree models are formed using these datasets, and finally average of predictions are taken.⁴⁵

- e) Gradient boosting (GB): Boosting is also an ensemble method that helps to form a strong learner by combining many weak learners. It is also a tree-based method in which decision trees are generated sequentially, and every tree tries to correct its predecessor. Gradient boosting (GB) is one of the boosting methods in which it tries to fit the current predictor with residual error made by the previous predictor.⁴²
- f) XGBoost: It stands for Extreme Gradient Boosting and was first implemented by the researchers of the University of Washington. This method was built based on the same algorithm of GB, but the main drawback of GB is that it searches for minimizing the loss function across all possible splits to create a new branch of a decision tree. Thus, GB method becomes time-consuming when thousands of features are present because there are thousands of possibilities to split the node of a decision tree. XGBoost handles these drawbacks by taking information of feature distribution across all data points in a single leaf node and by this way it reduces the search space. This method cannot generate multiple decision trees in parallel but can generate multiple branches of a decision tree in parallel.⁴⁶

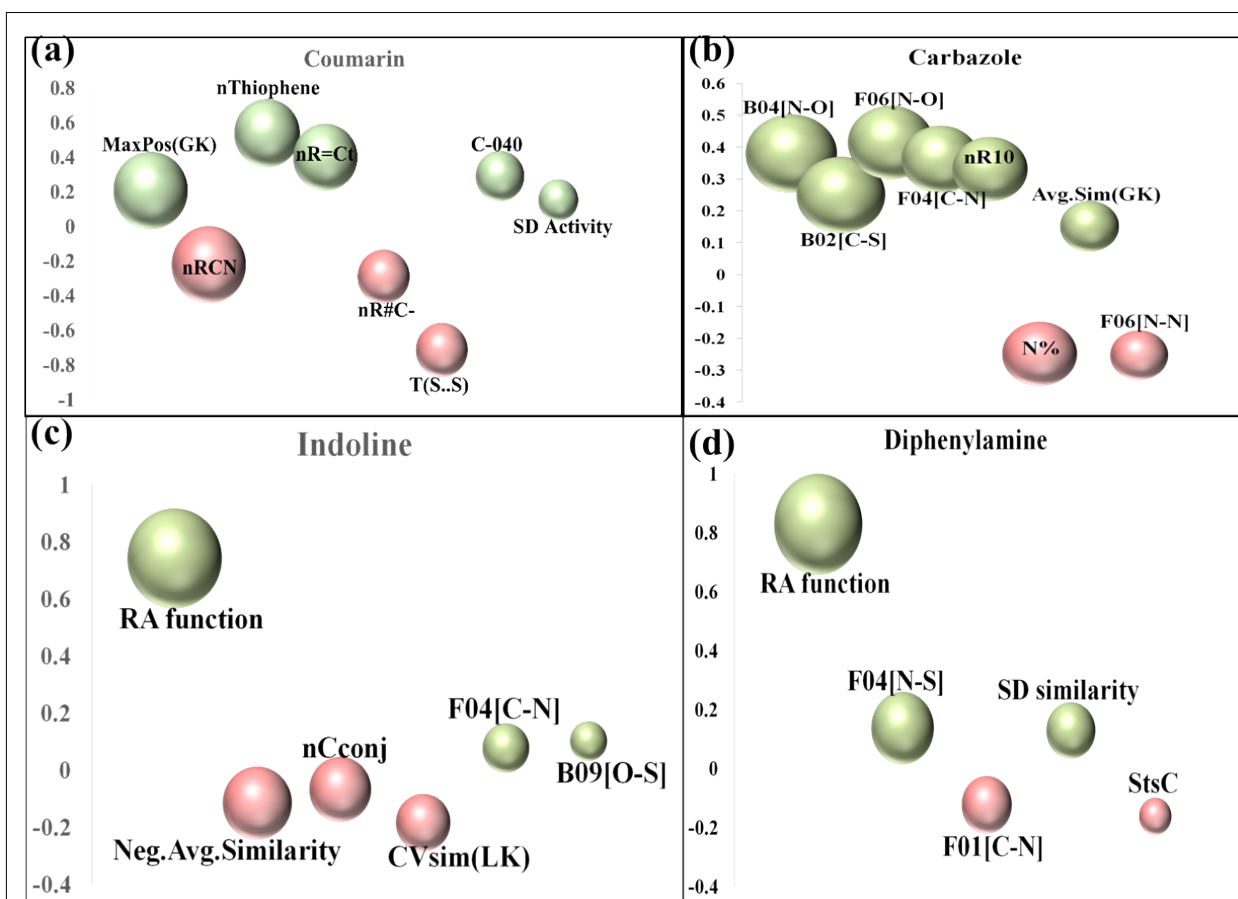


Figure S1. Bubble plot showing regression coefficients for individual descriptors in the PLS model (size of the bubble shows variable importance score)

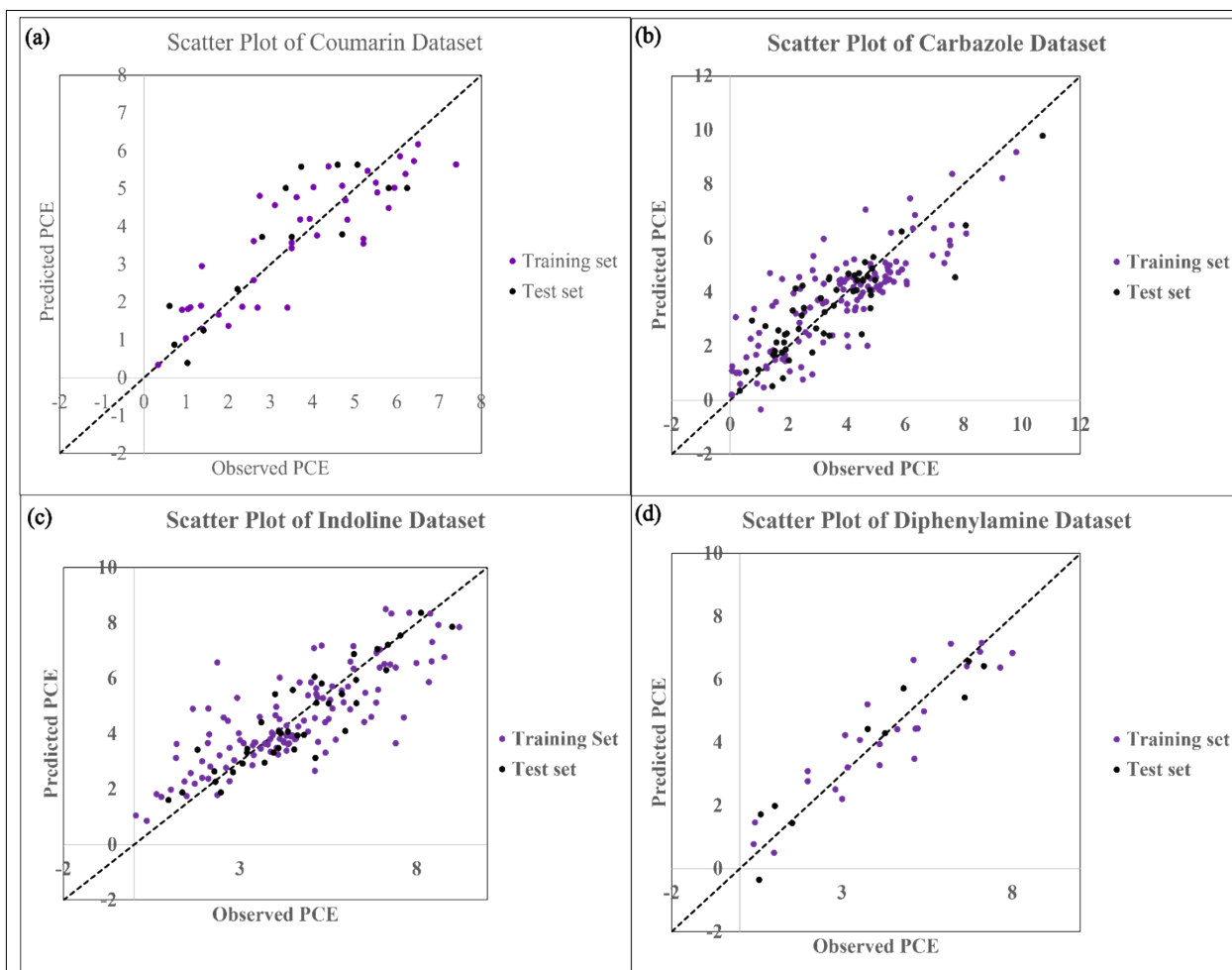


Figure S2. Scatter Plots indicating the prediction quality of the developed q-RASPR PLS models

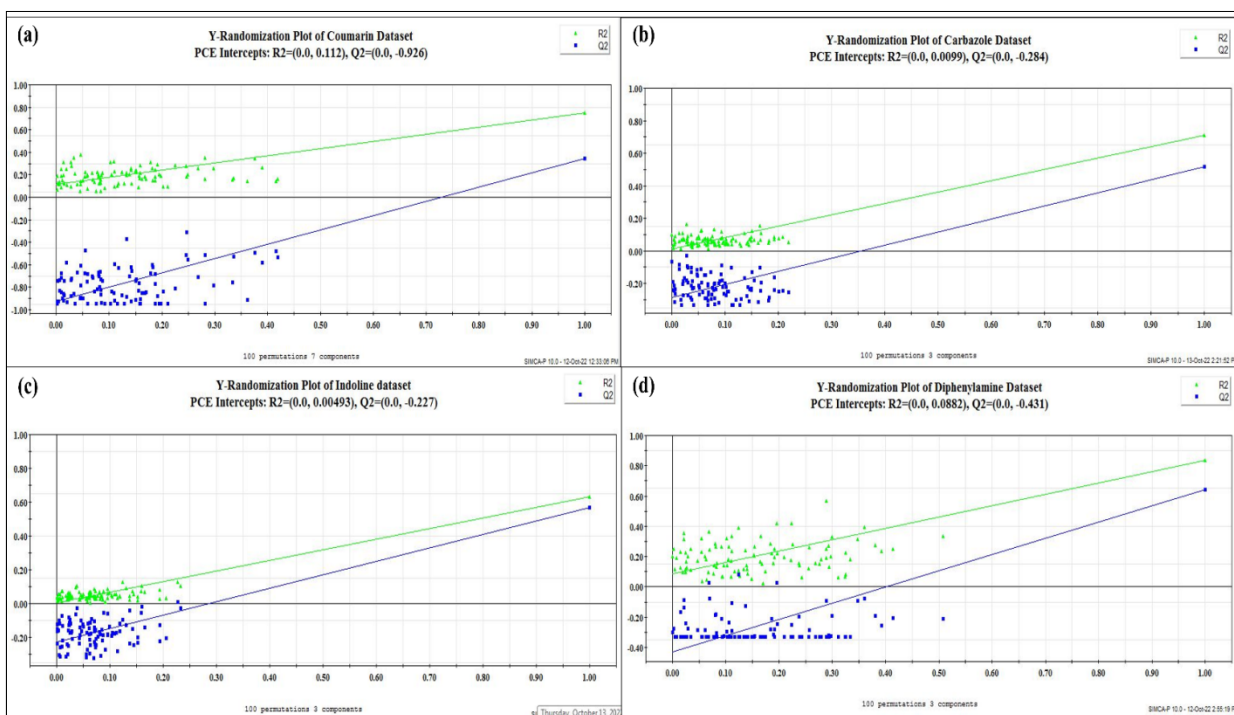


Figure S3. PLS randomization plots

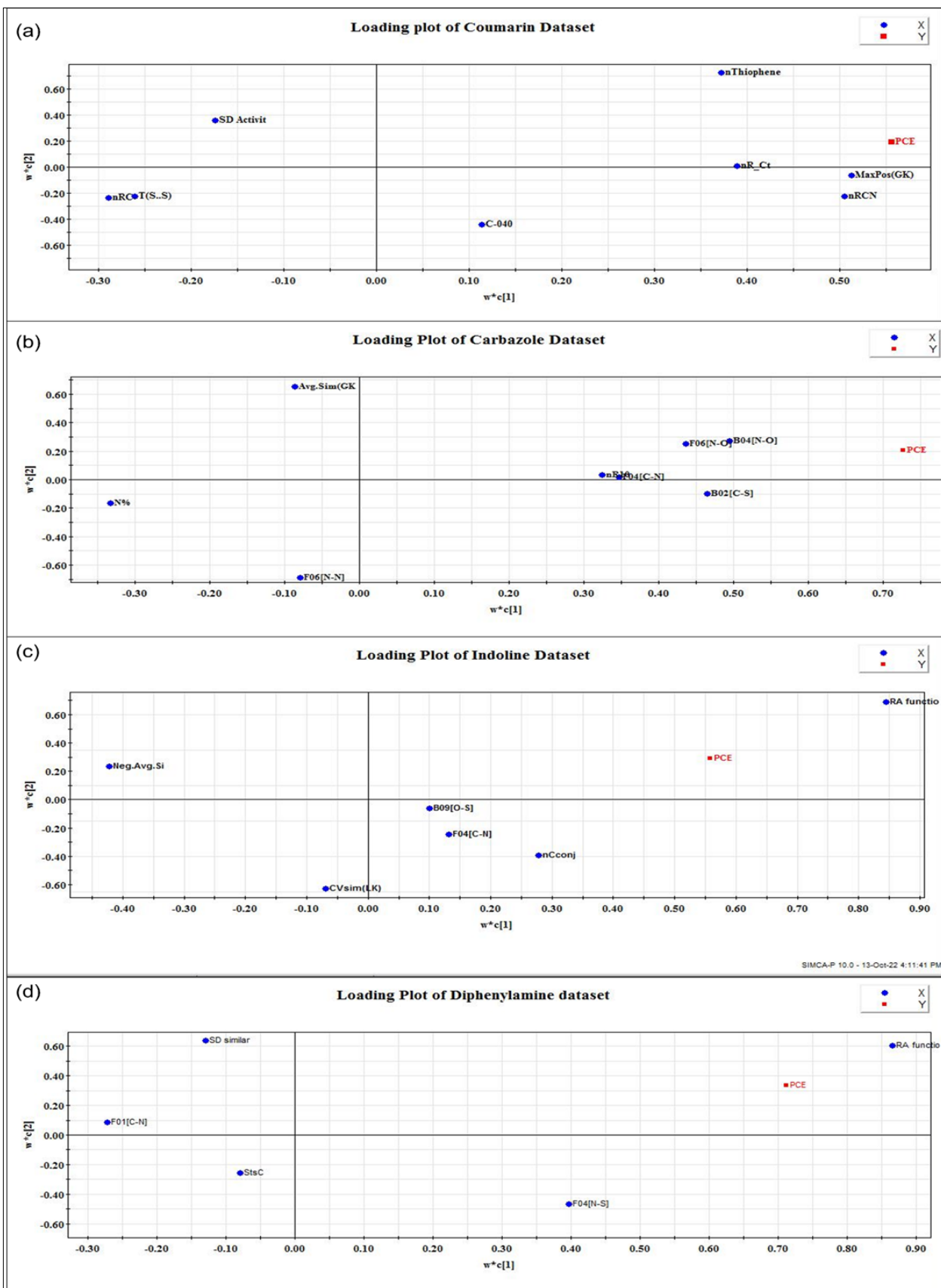


Figure S4. Loading Plots showing significance of the descriptors to PCE

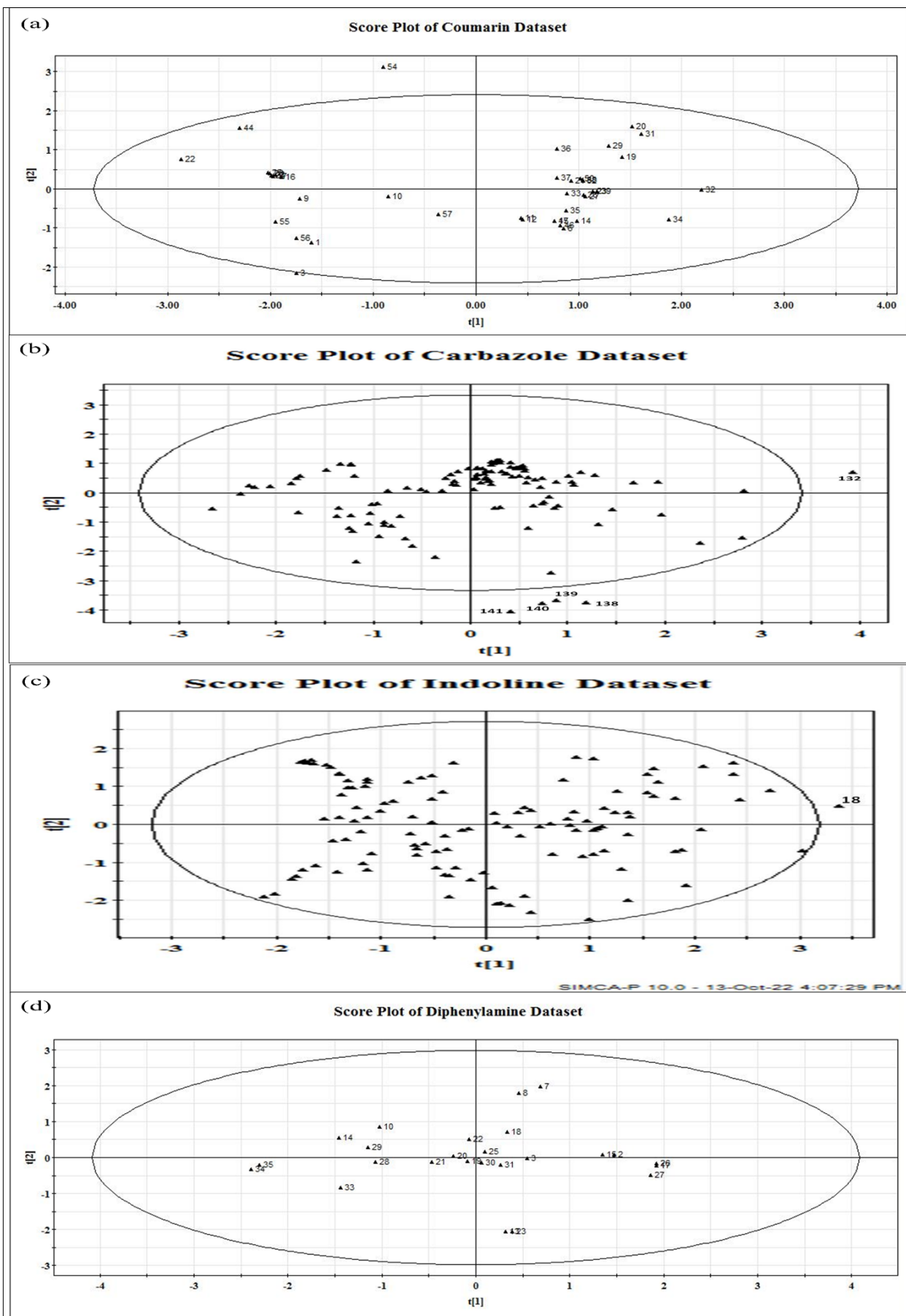


Figure S5. Score Plots indicating the applicability domain of PLS models

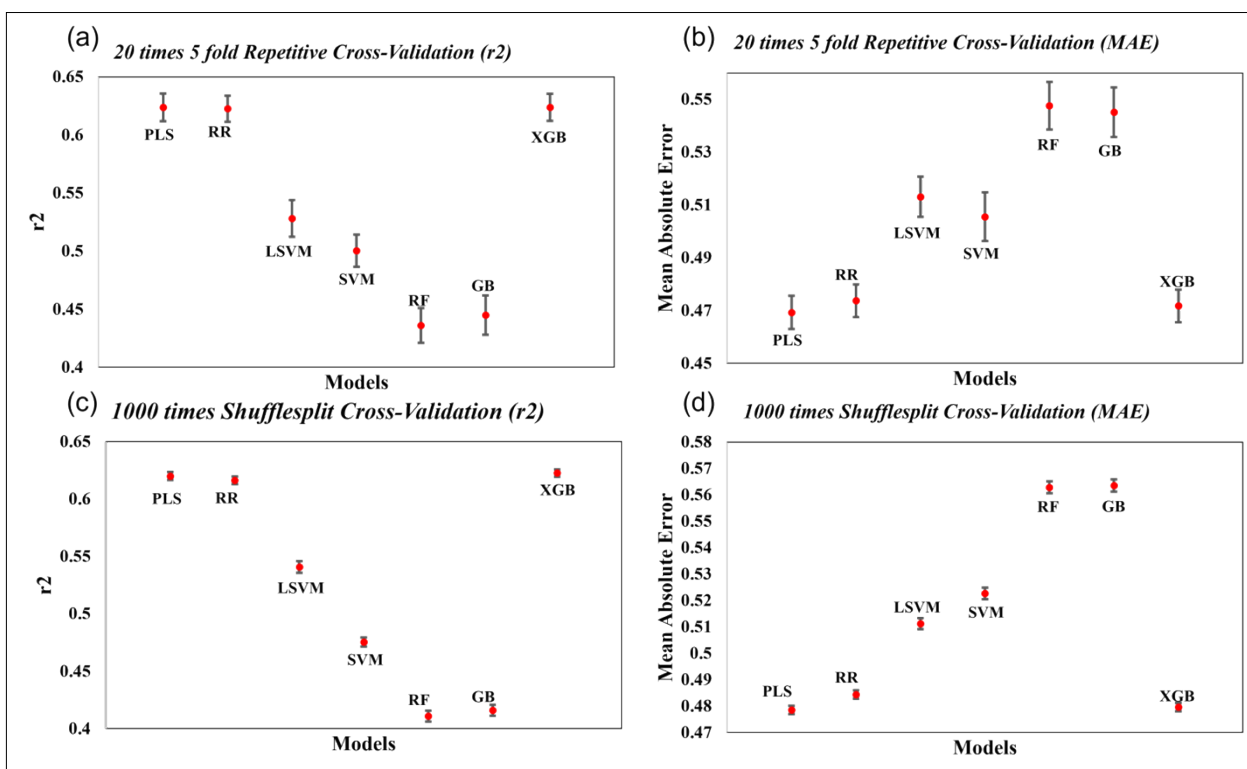


Figure S6. Cross-validation statistics based on 20 times 5-fold repetitive CV and 1000 shuffle split CV method (Mean \pm SEM) for the carbazole dataset

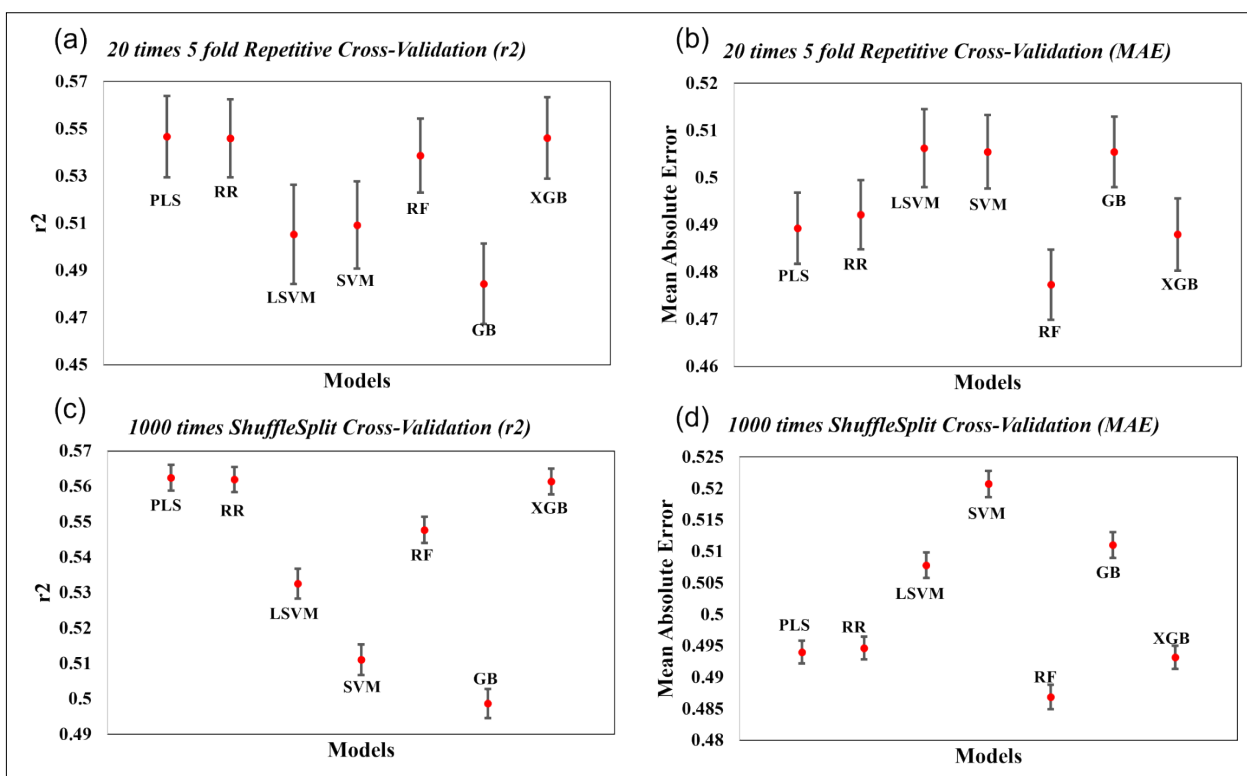


Figure S7. Cross-validation statistics based on 20 times 5-fold repetitive CV and 1000 shuffle split CV method (Mean \pm SEM) for the indoline dataset

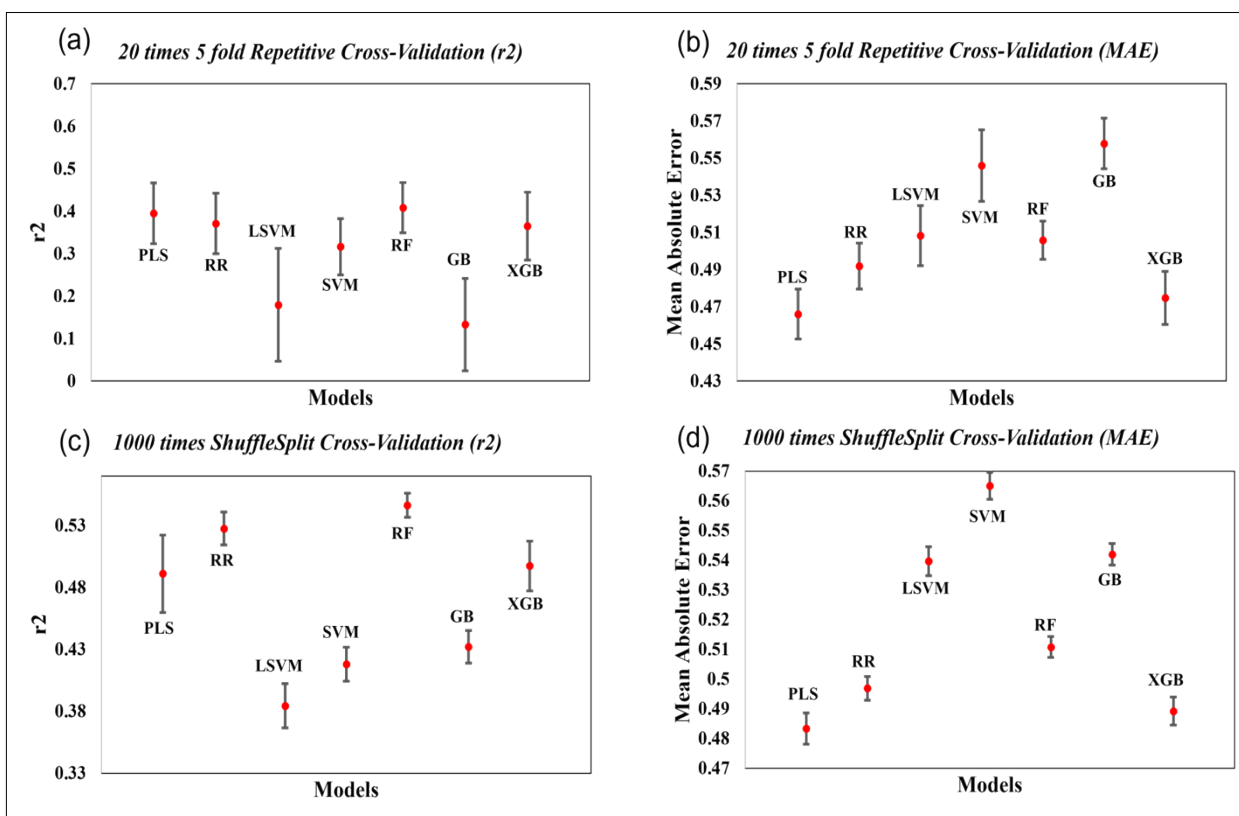


Figure S8. Cross-validation statistics based on 20 times 5-fold repetitive CV and 1000 shuffle split CV method (Mean \pm SEM) for the diphenylamine dataset

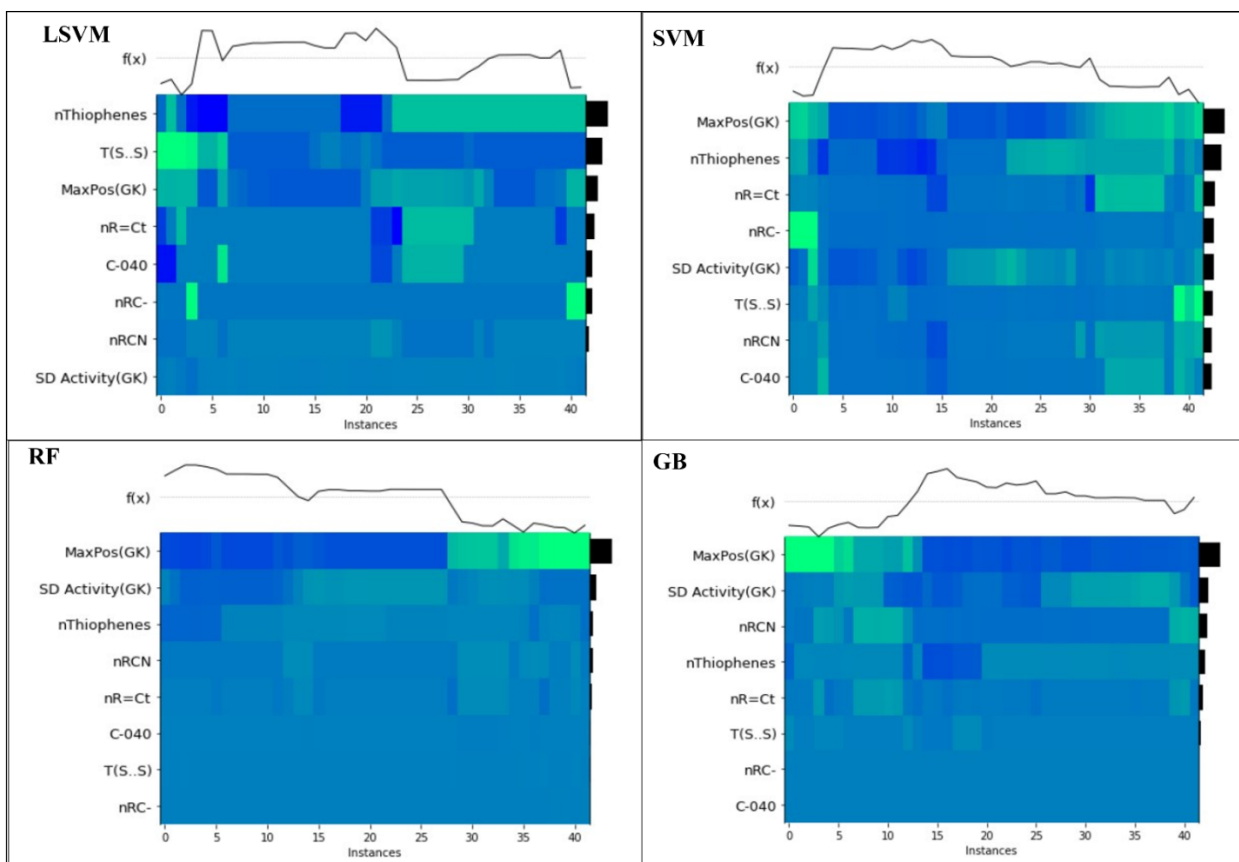


Figure S9. Heatmap plots of LSMV, SVM, RF and GB models for the coumarin dataset, indicating the relative importance of descriptors

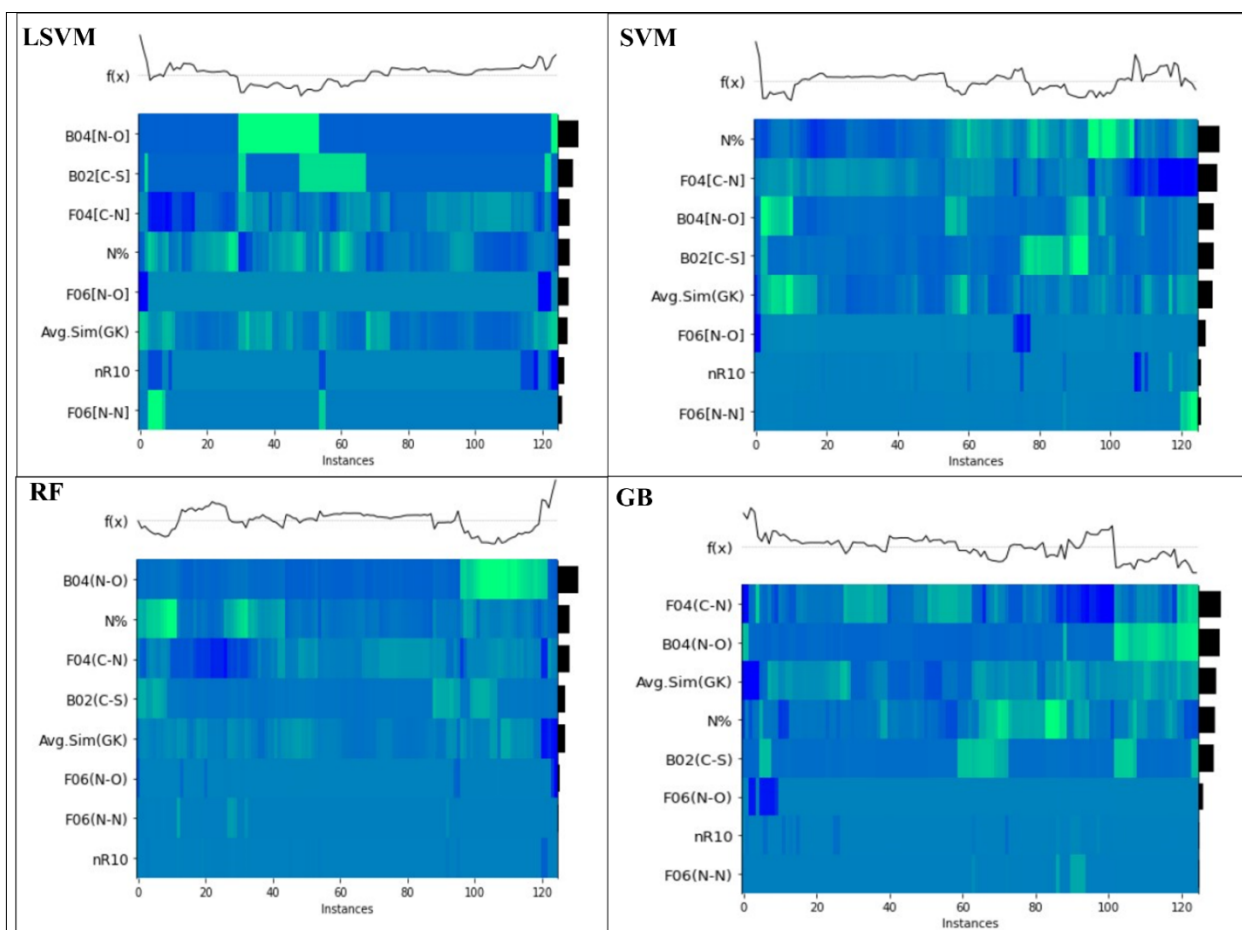


Figure S10. Heatmap plots of L1-SVM, SVM, RF and GB models for the carbazole dataset, indicating the relative importance of descriptors

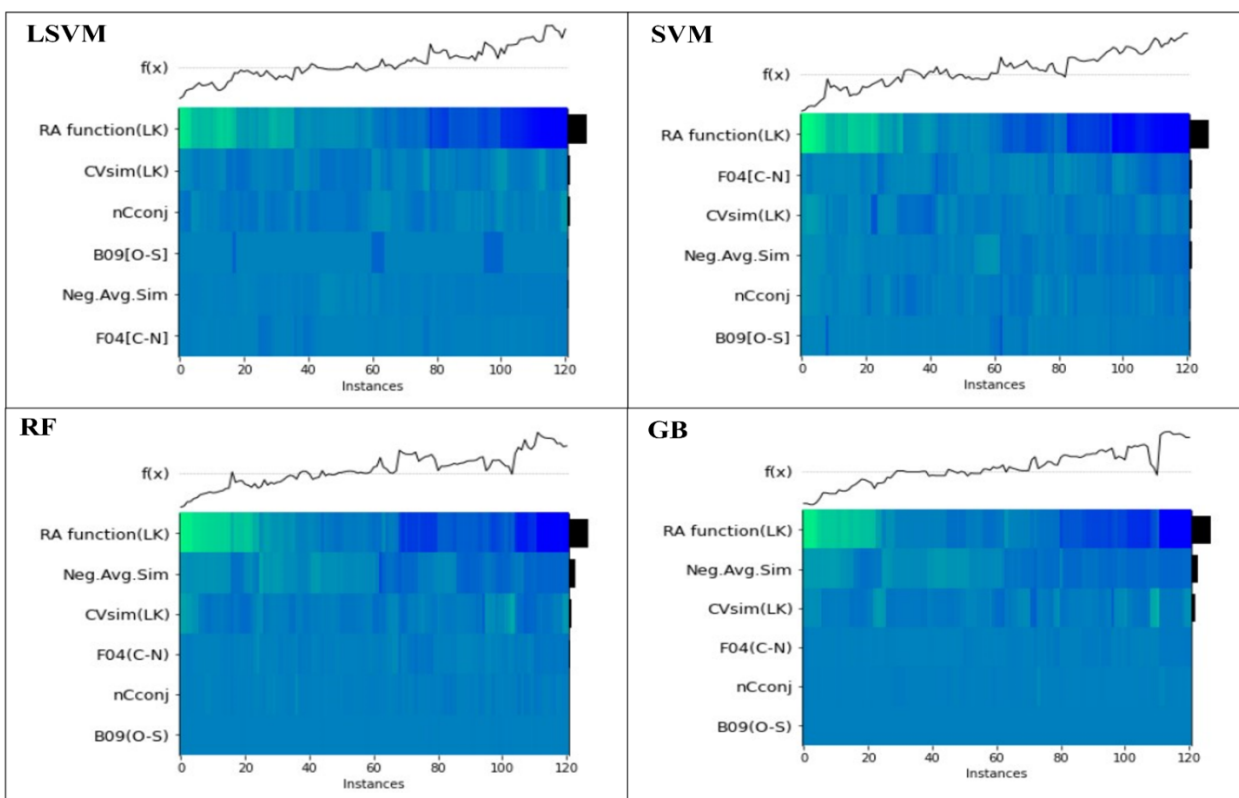


Figure S11. Heatmap plots of LSVM, SVM, RF and GB models for the indoline dataset, indicating the relative importance of descriptors

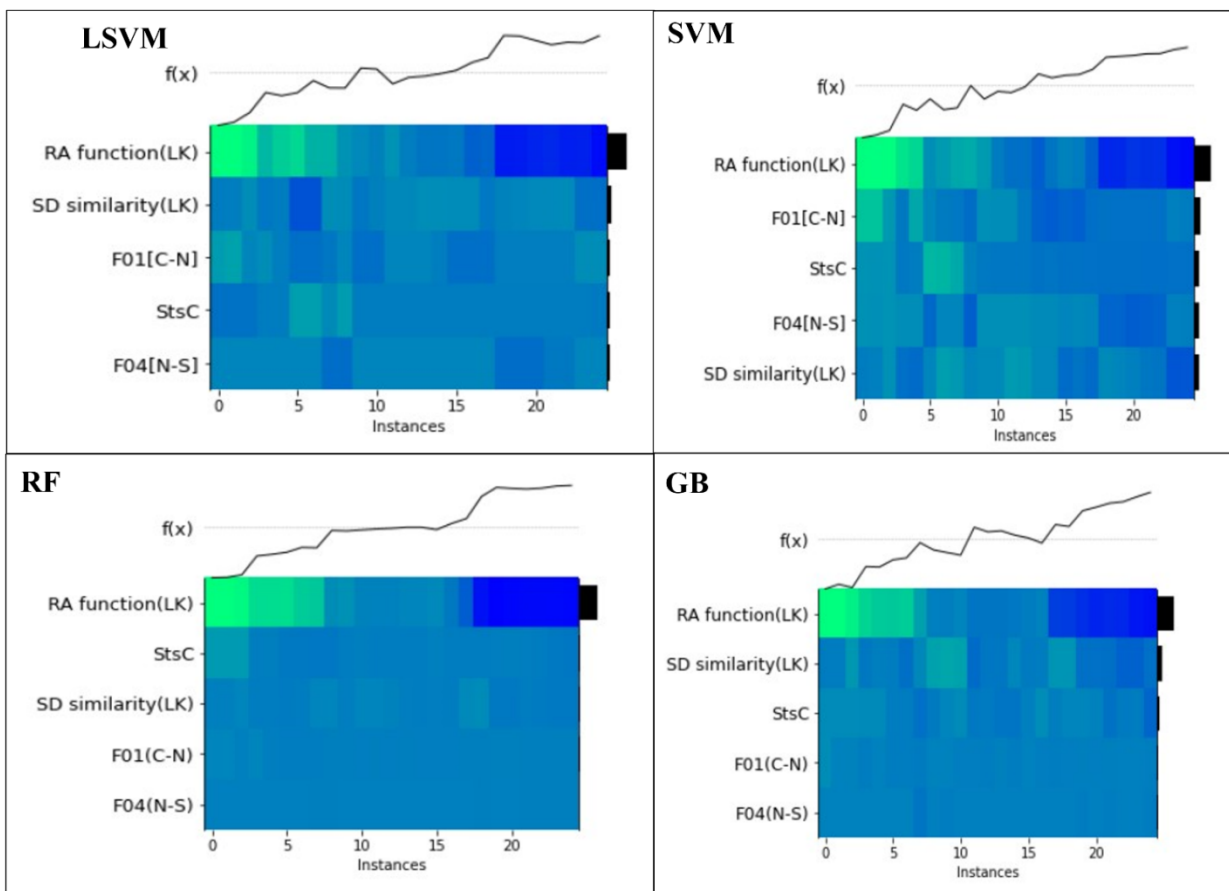
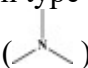
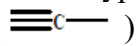


Figure S12. Heatmap plots of LSVM, SVM, RF and GB models for the diphenylamine dataset, indicating the relative importance of descriptors

Table S1. Initial pool of descriptors used for the read-across analysis

Dataset	Descriptors
Coumarin	<p>nRCN : number of aliphatic nitriles in the structure</p> <p>F08[N-S] : Frequency of N and S atoms at topological distance 8</p> <p>B09[S-S] : Presence or absence of S and S atoms at topological distance 9</p> <p>nCconj : number of non-aromatic conjugated C(sp²) atoms in the structure</p> <p>nArNR2 : number of aromatic tertiary amines in the structure</p> <p>B08[N-S] : Presence or absence of N and S atoms at topological distance 8</p> <p>C-034 : it is an atom centered fragment based descriptor which indicate the fragment - R-CR..X (R=any group linked through carbon atom, X=electronegative atoms like N, S, P, O, Halogens)</p> <p>nThiophenes : indicate the number of Thiophene rings in the structure</p> <p>nR#C- : number of a non-terminal carbon atom with the 'sp' hybridization</p> <p>nR=Ct : number of an aliphatic tertiary carbon atom with the 'sp²' hybridization</p> <p>T(S..S) : a 2D atom pair descriptor that indicates the sum of the topological distance between two sulfur atoms where they are part of two thiophene rings</p> <p>C-040 : is an atom-centered fragments descriptor that represents fragments like R-C(=X)-X/R-C#X/X = C = X (R: any group linked through carbon; X: any electronegative atom like N, S, P, O, halogen; #: triple bond)</p>
Carbazole	<p>F08[O-O] : Frequency of O and O atoms at a topological distance of 8</p> <p>NaasC : representing the Number of atoms of aasC (-C(-)-)</p> <p>F06[N-N] : representing the Frequency of N and N atoms at topological distance 6</p> <p>F06[C-C] : representing the Frequency of C and C atoms at topological distance 6</p> <p>nR10 : representing the number of 10-membered rings in the structure</p> <p>F04[C-N] : representing the Frequency of C and N atoms at topological distance 4</p> <p>B08[O-S] : representing the Presence or absence of O and S atoms at a topological distance 8</p> <p>B04[N-O] : representing the Presence or absence of N and O atoms at topological distance 4</p> <p>N% : total percentage of N atoms in the structure</p> <p>F06[N-O] : representing the Frequency of N and O atoms at topological distance 6</p> <p>B02[C-S] : representing the Presence or absence of C and S atoms at topological distance 2</p> <p>B10[C-S] : representing the Presence or absence of C and S atoms at topological distance 10</p> <p>B06[N-S] : representing the Presence or absence of N and S atoms at topological distance 6</p> <p>F06[O-S] : representing the Frequency of O and S atoms at topological</p>

	<p>distance 6 B04[O-S] : representing the Presence or absence of O and S atoms at topological distance 4</p>
Indoline	<p>SaaaC : it is an atom type E-state indices indicate the Sum of aaaC E-states (aaCa where a is aromatic bond) B07[N-N] : representing the Presence or absence of N and N atoms at topological distance 7 B06[N-N] : representing the Presence or absence of N and N atoms at topological distance 6 B04[S-S] : representing the Presence or absence of S and S atoms at topological distance 4 F04[C-N] : representing the Frequency of C and N atoms at topological distance 4 F07[N-S] : representing the Frequency of N and S atoms at topological distance 7 nCrq : number of ring quaternary C(sp³) atoms in the structure F10[C-N] : representing the Frequency of C and N atoms at topological distance 10 F07[N-O] : representing the Frequency of N and O atoms at topological distance 7 NsssN : it is an atom type E-state descriptors indicate the Number of atoms of type sssN () B05[O-S] : representing the Presence or absence of O and S atoms at topological distance 5 B09[O-S] : representing the Presence or absence of O and S atoms at topological distance 9 B05[S-S] : representing the Presence or absence of S and S atoms at topological distance 5 F04[S-S] : representing the Frequency of S and S atoms at topological distance 4 F07[N-N] : representing the Frequency of N and N atoms at topological distance 7 nCconj : representing the number of non-aromatic conjugated C (sp²) atoms in the structure F10[O-S] : representing the Frequency of O and S atoms at topological distance 10 B02[N-O] : representing the Presence or absence of N and O atoms at topological distance 2</p>
Diphenylamine	<p>F01[C-N] : representing the Frequency of C and N atoms at topological distance 1 F08[C-N] : representing the Frequency of C and N atoms at topological distance 8 StsC : It is an atom type E-state descriptor that indicates the sum of tsC E-states () nPyrimidines : It is functional group count descriptor indicate the number of Pyrimidines in the structure nCsp : It is a constitutional descriptor indicate the number of sp hybridized carbon atoms in the structure B08[N-N] : representing the Presence or absence of N and N atoms at topological distance 8 C-041 : It is an atom centered fragment corresponds to X-C(=X)-X,</p>

	<p>where, X can be any electro negative atom O, N, S, P, Se and halogens connected with the carbon atom</p> <p>nHAcc : It is a functional group count descriptor indicate the number of acceptor atoms for H-bonds (N,O,F)</p> <p>nR#C- : nR#C- : number of a non-terminal carbon atom with the 'sp' hybridization</p> <p>F04[N-S] : representing the Frequency of N and S atoms at topological distance 4</p> <p>ETA_dBeta : It is an extended topochemical atom descriptor, measuring the relative unsaturation content ($\Delta\beta$)</p>
--	--

Table S2. List of read-across derived descriptors and their definition

Descriptors	Description	Mathematical equation
RA function	It is a read-across weighted average prediction score for a target compound which is obtained based on the similarity between selected close source compound and target (or query) compound.	$RA\ function = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$ $w_i = \frac{S_i}{\sum_{i=1}^n S_i}$ <p>w_i = weightage given to individual selected close source compounds S_i = Similarity between individual selected close source compounds and target compounds x_i = observed response value of the selected close source compounds</p>
SD Activity ($S_{weighted}$)	It is the standard deviation of the observed response value of selected close source compounds for each query compounds.	$S_{weighted} = \sqrt{\frac{\sum_{i=1}^n (x_i - x_{wtd})^2}{\sum_{i=1}^n w_i}} \times \frac{n}{n-1}$ $x_{wtd} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$ <p>n = effective number of selected close source compound w_i = weightage given to individual selected close source compounds</p>

SE	It is the standard error of the observed response values of selected close source compounds for each query compounds.	$SE = \frac{SD \text{ Activity}}{\sqrt{n}}$
CVact (CV_{activity})	It is the coefficient of variance of the observed response (or activity) values of the selected close source compounds for each query compounds.	$CV_{activity} = \frac{S_{weighted}}{x_{wtd}^-}$
CVsim (CV_{similarity})	It is the coefficient of variance of the similarity values of the selected close source compounds for each query compounds.	$CV_{similarity} = \frac{SD \text{ Similarity}}{\bar{f}}$ \bar{f} = average similarity value of the selected close source compounds
MaxPos	Maximum similarity level to Positive close source compounds based on mean value of training set observed response.	
MaxNeg	Maximum similarity level to Negative close source compounds based on mean value of training set observed response.	
Abs MaxPos- MaxNeg (Abs.Diff.)	Absolute difference between MaxPos and MaxNeg.	$AbsDiff = MaxPos - MaxNeg $
Avg.Sim. (Average Similarity)	It is the mean of similarity values of the selected close source compounds for each query compounds.	$Avg.Sim. (\bar{f}) = \frac{\sum_{i=1}^n f_i}{n}$ f_i = similarity values of the selected close source compounds n=number of selected close source compounds
SD Similarity	It is the standard deviation of the similarity values of selected close source compounds for each query compounds.	$SD \text{ Similarity} = \sqrt{\frac{\sum_{i=1}^n (f_i - \bar{f})^2}{n - 1}}$ \bar{f} = average similarity value of the selected close source compounds n=number of selected close source compounds
g_m [Banerjee-Roy Coefficient]	A novel concordance measure	$g_m = (-1)^n Posfrac - 0.5 $ n=1 when MaxPos < MaxNeg n=2 when MaxPos >= MaxNeg Posfrac = Fraction of the close source compounds having a response value greater than the training set mean

		response
<i>gm*Avg.Sim</i>	Product of the values of g_m and <i>Average Similarity</i>	
<i>gm*SD Similarity</i>	Product of the values of g_m and <i>SD Similarity</i>	
<i>Pos.Avg.Sim</i>	Average similarity value of the selected positive close source compounds based on the mean observed response value of training set	
<i>Neg.Avg.Sim</i>	Average similarity value of the selected negative close source compounds based on the mean observed response value of training set	