

1 **The Effect of Leverage and Influential on Structure-Activity Relationships**

2

3 Sorana D. BOLBOACĂ¹ and Lorentz JÄNTSCHI^{2,3,*}

4

5 ¹"Iuliu Hațieganu" University of Medicine and Pharmacy Cluj-Napoca, Department of Medical
6 Informatics and Biostatistics, 6 Louis Pasteur, 400349 Cluj-Napoca, Cluj, Romania.

7 ²Technical University of Cluj-Napoca, Department of Chemistry, 103-105 Muncii Bdv., 400641
8 Cluj-Napoca, Romania.

9 ³University of Agricultural Science and Veterinary Medicine Cluj-Napoca, 3-5 Calea Mănăştur,
10 400372 Cluj-Napoca, Romania.

11 E-mails: sbolboaca@umfcluj.ro; lorentz.jantschi@gmail.com

12

13 *Author to whom correspondence should be addressed; E-Mail: lorentz.jantschi@gmail.com;

14 Tel.: +4-0264-401-775; Fax: +4-0264-401-768.

15

16

17 Running title: Leverage and Influential on QSARs

18

19 **The Effect of Leverage and Influential on Structure-Activity Relationships**

20

21

22 **Abstract**

23 Quantitative Structure-Activity Relationship approaches have established as the main computational
24 molecular modeling method. In spirit of reporting valid and reliable models the aim of our research was to
25 assess how the analysis of leverage with Hat matrix (h_i) and of the influential using Cook's distance (D_i) of
26 QSAR models reflects in the model reliability and its characteristics. The datasets included in this research
27 was collected from previously published manuscripts. Seven datasets accomplished the imposed inclusion
28 criteria and were analyzed. Three models were obtained for each dataset (full-model, h_i -model and D_i -model)
29 and several validation criteria (statistical criteria) were defined to assess and to compare the model. The
30 analysis of the obtained results revealed that in 5 out of 7 sets the correlation coefficient increase when both
31 compounds with h_i and respectively D_i higher than thresholds were removed. The number of withdrawn
32 compounds varied from 2 to 4 for h_i -model and from 1 to 13 for D_i -model. The analysis of validation
33 statistics showed that D_i -models obtained systematically better results compared to both full-models and h_i -
34 models. Identification of influential compound in data set could significantly improve the model and should
35 be conducted any time when a regression analysis is desired. Cook's distance approach is recommended to
36 be used to identify influential compounds in dataset whenever the linear regression analysis for QSAR
37 models is applied.

38

39 **Keywords:** model sensitivity; quantitative structure-activity relationship (QSAR); leverage (h_i);
40 Cook's distance (D_i); model validation.

41

42

43 **Introduction**

44 Translation of structural features of chemical compounds in the activity by incorporation of
45 physico-chemical mechanisms into statistical models led to development of QSAR/QSPR
46 (Quantitative Structure-Activity/Property Relationship) computational molecular modeling
47 methodologies. In view of the fact that the capabilities of collecting and storing (such as PubChem)
48 from one hand and analyzing of data from other hand due to rapid development of information and
49 communication technologies have significant increased, QSAR modeling could be seen as an
50 approach of statistical analyses as well as application of data-mining.

51 Guidance regarding the correct procedures in QSAR development has been published in scientific
52 literature [1-3]. The detailed description of QSAR modeling techniques, methodologies and trends is
53 beyond the aim of the present manuscript. It is well known that the main characteristic of a QSAR
54 model is its predictivity, translated in how well the model is able to predict the activity on
55 compounds not used to develop the model. Guidelines for validation of QSAR models have been
56 developed by experts [4-6]. Beside good practice principles, other QSARs problems were addressed
57 by researchers. Mekenyan and Veith [7] pointed out two general problems of QSAR: various
58 environments used to study the property/activity and proliferation of molecular descriptors. Dearden
59 and co-authors identified 21 types of errors in QSAR modeling, errors classified according to
60 OECD principles [8]. From statistical point of view, the identified errors were as follow [8]:

- 61 ▪ Collinearity of molecular descriptors which is mainly reflected in the instability of the
62 regression coefficients [1,2].
- 63 ▪ Outlier detection and removal. Removal of a significant outlier led to a more significant model
64 [9].
- 65 ▪ Lack of/inadequate statistics. In most of published QSAR models, neither considerations of
66 linear regression assumptions nor considerations of distribution of residuals are addressed
67 [10,11]. Recommended statistics are as follow: n (sample size), r^2 or R^2 (determination
68 coefficient), q^2 or Q^2 (determination coefficient in leave-one-out analysis); R^2_{adj} (adjusted
69 determination coefficient), s (standard error of estimate - measure of error) and F-statistics
70 (including p-value) [8]. Moreover, other methods of error are recommended: standard deviation,
71 root mean square error and mean absolute error (ignore the sign of an error – provide
72 information about random error [12]), mean error (consider the sign of an error – very low value
73 indicates the absence of systematic error [12]; similar mean error and absolute mean error
74 indicate the presence of systematic errors).
- 75 ▪ Misuse/Misinterpretation of statistics. The application of linear regression technique without
76 investigation of its assumption is one of most frequent misuse of statistics [13]. The inclusion in
77 the model of additional independent variable(s) is another example [14].

78 Staying in the field of statistics for QSAR/QSPR models the following was the hypothesis of the
79 present research: Model sensitivity analysis translated through influential point(s) could identify a
80 stable and reliable QSAR/QSPR. Our aim was to assess how the analysis of leverage and influential
81 using Cook's distance of QSAR models reflects in the model reliability and its characteristics.

82

83

84 **Materials and Methods**

85 ***Data sets***

86 Several datasets previously published in International Journal of Molecular Science (MDPI
87 Publishing, Basel, Switzerland) were included in our analysis. The search was conducted on April
88 2012 using the following search strategy:

Where? (Field)	What?
Title/Keyword	QSAR OR Quantitative Structure-Activity Relationship
Journal	IJMS
Article Type	Article OR Review
Time period	2000 to date

89

90 There were included in the study the dataset available in the previously published manuscripts
91 that respected the following inclusion criteria: ① quantitative continuous dependent variable AND ②
92 values of descriptors provided in manuscript or supplementary material(s) AND ③ sample size > 20
93 AND ④ simple/multiple linear regression model with determination coefficient higher than or equal
94 to 0.6.

95

96 ***Analysis of Influential***

97 Model sensitivity in linear regression analysis refers to how estimates are affected by subgroups of
98 the data. Three main issues could be used to assess the model sensitivity: residuals (large value
99 identify the outliers), leverage (large value identify the point significantly far from the center point
100 of the predictor space) and influential (large effect on an estimate) but just two of them are
101 addressed in the present research.

102 The following steps were applied to accomplish the aim of the research:

- 103 ▪ **Step 1:** Test the normality of observed/measured activity using Kolmogorov-Smirnov [15]
104 and/or Chi-Square goodness-of-fit [16] → If data normal construct the SLM (Simple Linear
105 Regression) / MLR (Multiple linear regression)
- 106 ▪ **Step 2:** Identify the best SLM / MLR model → If $R^2 < 0.5$ STOP analysis. The dataset is
107 removed from further analysis.

- 108 ▪ **Step 3:** Identify the influential using:
- 109 a. Hat matrix - leverage (h_i). Leverage are “a measure of the geometric distance of the i^{th}
- 110 predictor point ($X_{i1}, X_{i2}, \dots, X_{ik}$) from the center point of the predictor space” [17]. The
- 111 formula applied to identify the leverage was: $h_i = 1/n + (x_i - x_m)^2 / \sum[(x_i - x_m)^2]$, where $h_i =$
- 112 leverage of the i^{th} compound, $n =$ sample size, $x_i =$ the value of predictor variable for the
- 113 i^{th} compound, $x_m =$ the average mean for predictor x . The leverage indicates those
- 114 compounds that may have potential influence in the model being used also as
- 115 applicability domain of the QSAR models [18,19]. The leverage threshold (h_t) was set to
- 116 $2*(k+1)/n$ for regression models with intercept and $2*k/n$ for models without intercept
- 117 (where $k =$ number of descriptors in the model; $n =$ sample size) [17]. → If $h_i > h_t$
- 118 withdrawn the influential till no leverage exceed the threshold value or no improvement
- 119 in the determination coefficient is observed.
- 120 b. Cook’s distance (D_i). Cook's distance combines residual and leverage in one indicator to
- 121 identify influential in regression models [20,21]. Any compound was considered as
- 122 influential if $D_i > 4/n$ (where $n =$ sample size) [22]. → If $D_i > 4/n$ withdrawn the
- 123 influential till no exceed of the threshold value is observed or no increase in the
- 124 determination coefficient is observed.
- 125 ▪ **Step 4:** Construct and evaluate the final SLM / MLR. The criteria used for assessment and
- 126 validation of QSAR models are presented in Table 1. The correlated correlation analysis was
- 127 apply to test if correlation coefficients obtained by full-model, h_i -model and D_i -model are
- 128 statistically significant different at a significance level of 5% [23].
- 129 ▪ **Step 5:** Take two sets of compounds and split the dataset in training (~2/3 compounds) and test
- 130 set using a simple random approach [24] (leave-many-out analysis) in order to assess the
- 131 behavior of the full-model and respectively model with higher correlation coefficient and
- 132 smaller standard error.

133 To test the overall performances of leverage and influential withdrawn on QSAR models compared

134 to full-model the Fisher's Chi-Squared (abbreviated as F-C-S) was applied [32]. The F-C-S- test

135 was applied to test the following null hypothesis "The correlation coefficient on a specific model

136 (such as h_i -model or D_i -model) is statistically higher compared to another model (full-model or h_i -

137 model when D_i -model was compared to h_i -model)".

138

Table 1. Criteria for validation of regression models.

Criterion	Interpretation/Remark
Goodness-of-fit	
R^2 = determination coefficient	A descriptive measure. It does not measure the quality of the regression model. The higher the better
R^2_{adj} = adjusted determination coefficient	Its value decrease if an added predictor does not reduce the unexplained variance Used as a measure of usefulness of introducing a new variable in the model Closeness to the R^2 the better
R^2_{loo} = determination coefficient in leave-one-out analysis [25]	Internal validity of the model Underestimates the true predictive error when small samples are used to develop the model [26] Closeness to the R^2 the better
s = standard error of estimate	Measure of the dispersion around the regression line of observed values Smaller the better
s_{loo} = standard error of predicted	
F-value (p-value)	Ration between explained and unexplained variance of a given number of df – degrees of freedom p-value associated to F-value as significance of the level of correlation [27] The higher the better
F_{loo} (p-value)	
t-value (p-value)	Significance of the coefficients in the regression model t-value - the higher the better vs. p-value - the lower the better
Validation statistics	
RMS = residual mean square	Error variance The lower the better
APV = average prediction variance [28]	The lower the better
TSE = total squared error [29]	The lower the better
APMSE = Average Prediction Mean Squared Error [30]	The lower the better
%PredErr = percentage prediction error [31]	Defined as prediction error (module of the difference between observed and estimated) divided by the highest activity
Predictive Power – Fisher's approach ($t_{pp} - p_{pp}$)	Evaluate if the mean of residuals is statistically different by the expected mean (where expected mean = 0); p_{pp} : the lower the better
RMSE = root-mean-square error	Measures the average magnitude of the error The lower the better
MAE = mean absolute error	Measures the average magnitude of the errors Could be also used to compare two models - The lower the better
MAPE = mean absolute percentage error	The lower the better
SEP = standard error of prediction	The lower the better
REP% = relative error of prediction	The lower the better

140

141 **Results**

142 Sixty-four manuscripts were identified using the applied search strategy. Fifteen manuscripts
 143 provided the experimental/observed values as well as values of molecular descriptors. Seven
 144 manuscripts accomplished all inclusion criteria and their sets of compounds were included in the
 145 analysis. The main characteristics of the previously published models (not necessary linear models)
 146 are presented in Table 2.

147 The identified sets of compounds were investigated in order to assess how the influential affect
 148 the model validity and characteristics. The best performing regression models for each set on the
 149 whole data set, on the sample after removal of compounds with leverage higher than threshold and

150 on the sample after removal of compounds with Cook's distance higher than threshold are presented
 151 in Table 3.

152

153 **Table 2.** Datasets included in analysis and basic summary of previously reported models.

Set [Ref]	Model characteristics
Set1 [33]	$R^2=0.9992$; $s=0.929$; $F=3534$; $n=60$; $k=5$
Set2 [34]	$R^2=0.7779$; $F=133$; $R^2_{loo}=0.774$; $n=79$; $k=2$
Set3 [35]	$R^2=0.820$; $R^2_{loo}=0.716$, $s=0.440$, $F=22.805$; $n=31$, $k=5$ (outliers: 5 & 15)
Set4 [36]	$R^2=0.9571$; $R^2_{cv}=0.8521$; $s=0.2825$; $F=28.8207$; $n=29$; $k=5$
Set5 [37]	$R^2=0.840$; $R^2_{cv}=0.777$; $F=31.54$; $s=0.034$; $n=36$; $k=5$
Set6 [38]	n.a.
Set7 [39]	$R_t=0.870$; $s=0.206$; $R_{test}=0.835$, $s_{test}=0.232$; $R_{loo}=0.925$; $s_{loo}=0.198$; $n=46$; $k=5$

R=correlation coefficient; R^2 =determination coefficient; loo=leave-one-out analysis; s=standard error of estimate; F=Fisher's statistics; n=sample size; k=number of independent variables used by the reported model; tr=training set; test=test set; n.a. = not available

154

155 **Table 3.** Regression characteristics: full-model (whole dataset), h_i -model (withdrawn of compounds
 156 with $h_i > h_t$) and D_i -model (withdrawn of compounds with $D_i > 4/n$, where n = sample size).

Set1: $\hat{Y}_{HF} = a + b_1 \times \chi^2 + b_2 \times H^* + b_3 \times J^*$ where \hat{Y} = estimated heat of formation; χ^2 = generalized connectivity index; H^* = Harary index; J^* = Balaban index; HF = heats of formation; a = intercept; b_i = regression coefficients		
n=60	whole dataset	$R^2=0.985$; $R^2_{adj}=0.985$; $s=3.46$; $F=1256$ ($p=2.63 \cdot 10^{-55}$); $R^2_{loo}=0.983$; $s_{loo}=3.76$, $F_{loo}=1061$ ($p=2.91 \cdot 10^{-51}$); $RMS=11.774$; $APV=0.003$; $TSE=4$; $APMSE=0.210$; %PredErr= 6.190; $t_{pp}=5.23 \cdot 10^{-14}$ $(p_{pp}=1)$; $RMSE=3.462$; $MAE=2.881$; $MAPE=0.396$; $SEP=3.191$; $REP(\%)=26.581$
n=56	$h_i > 2 \cdot (k+1)/n$ withdrawn (1, 38, 39, 40)	$R^2=0.987$; $R^2_{adj}=0.986$; $s=3.35$; $F=1318$ ($p=5.08 \cdot 10^{-49}$); $R^2_{loo}=0.986$; $s_{loo}=3.54$, $F_{loo}=882$ ($p=3.67 \cdot 10^{-49}$); $RMS=10.980$; $APV=0.003$; $TSE=4$; $APMSE=0.211$; %PredErr= 5.529; $t_{pp}=2.71 \cdot 10^{-13}$ $(p_{pp}=1)$; $RMSE=3.345$; $MAE=2.758$; $MAPE=0.354$; $SEP=3.253$; $REP(\%)=27.928$
n=54	$D_i > 4/n$ withdrawn (1, 2, 3, 16, 20, 23)	$R^2=0.989$; $R^2_{adj}=0.988$; $s=3.04$; $F=1441$ ($p=1.57 \cdot 10^{-48}$); $R^2_{loo}=0.987$; $s_{loo}=3.24$; $F_{loo}=1268$ ($p=2.67 \cdot 10^{-49}$); $RMS=9.059$; $APV=0.003$; $TSE=4$; $APMSE=0.181$; %PredErr= 4.866; $t_{pp}=1.25 \cdot 10^{-13}$ $(p_{pp}=1)$; $RMSE=3.040$; $MAE=2.517$; $MAPE=0.290$; $SEP=2.749$; $REP(\%)=29.9702$
Set2: $\hat{Y}(\log(1/EC_{50})) = a + b_1 \times \log P + b_2 \times MTD^*$ where $\hat{Y}(\log(1/EC_{50}))$ = estimated $\log(1/EC_{50})$ - EC_{50} = level that produces a 50% protection of MT-4 cells against HIV-1 cytopathic effect; $\log P$ = hydrophobicities; MTD^* = minimal topological difference descriptor [34]; a = intercept; b_i = regression coefficients		
n=79	whole dataset	$R^2=0.754$; $R^2_{adj}=0.747$; $s=0.68$; $F=116$ ($p=7.59 \cdot 10^{-24}$); $R^2_{loo}=0.733$; $s_{loo}=0.70$; $F_{loo}=104$ ($p=9.75 \cdot 10^{-23}$); $RMS=0.4516$; $APV=0.4630$; $TSE=3$; $APMSE=0.0059$; %PredErr= 4.636; $t_{pp}=1.74 \cdot 10^{-14}$ $(p_{pp}=1)$; $RMSE=0.676$; $MAE=0.503$; $MAPE=0.083$; $SEP=0.668$; $REP(\%)=10.546$
n=77	$h_i > 2 \cdot (k+1)/n$ withdrawn (57, 61)	$R^2=0.761$; $R^2_{adj}=0.754$; $s=0.66$; $F=118$ ($p=1.01 \cdot 10^{-23}$); $R^2_{loo}=0.714$; $s_{loo}=0.68$, $F_{loo}=106$ ($p=1.06 \cdot 10^{-22}$); $RMS=0.4275$; $APV=0.4386$; $TSE=3$; $APMSE=0.0058$; %PredErr=4.636; $t_{pp}=1.40 \cdot 10^{-14}$ $(p_{pp}=1)$; $RMSE=0.658$; $MAE=0.482$; $MAPE=0.080$; $SEP=0.650$; $REP(\%)=10.336$
n=66	$D_i > 4/n$ withdrawn (14,34,50,51,57-62,64,71,75)	$R^2=0.899$; $R^2_{adj}=0.895$; $s=0.41$; $F=279$ ($p=4.83 \cdot 10^{-32}$); $R^2_{loo}=0.891$; $s_{loo}=0.43$, $F_{loo}=256$ ($p=1.45 \cdot 10^{-31}$); $RMS=0.1964$; $APV=0.2024$; $TSE=3$; $APMSE=0.0031$; %PredErr=2.719; $t_{pp}=2.66 \cdot 10^{-1}$ $(p_{pp}=0.7910)$; $RMSE=0.412$; $MAE=0.353$; $MAPE=0.060$; $SEP=0.440$; $REP(\%)=7.059$

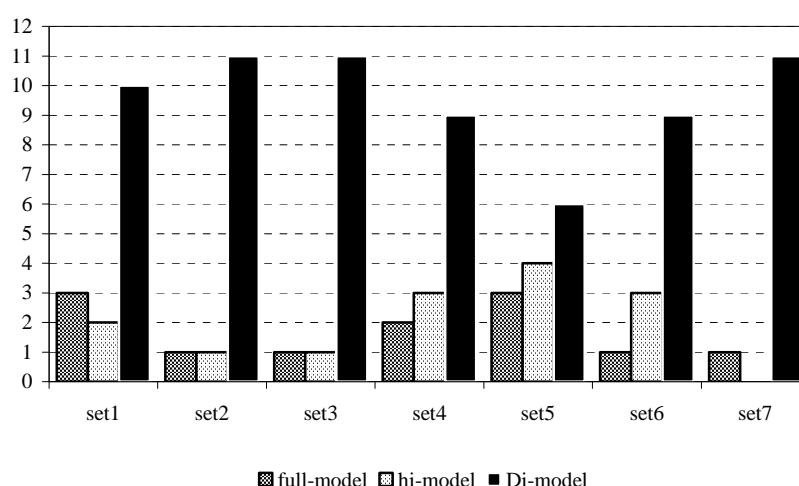
Set3: $\hat{Y}(\log K_i) = a + b_1 \times L + b_2 \times B_1 + b_3 \times B_3 + b_4 \times \text{FPSA}_3 + b_5 \times \rho$ where \hat{Y} = estimated activity; K_i = binding affinity; L = sterimol parameter; B_1, B_3 = sterimol width parameters; FPSA_3 = fractional charged partial surface area; ρ = density; a = intercept; b_i = regression coefficients		
n=33	whole dataset	$R^2=0.524$; $R^2_{\text{adj}}=0.436$; $s=0.69$; $F=6$ ($p=0.001$); $R^2_{100}=0.287$; $s_{100}=0.88$; $F_{100}=1.52$ ($p=0.2155$); $\text{RMS}=0.4588$; $\text{APV}=0.5283$; $\text{TSE}=6$; $\text{APMSE}=0.0170$; $\% \text{PredErr}=37.0833$; $t_{\text{pp}}=-1.39 \cdot 10^{-14}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.690$; $\text{MAE}=0.494$; $\text{MAPE}=0.679$; $\text{SEP}=0.634$; $\text{REP}(\%)=39.601$
n=31	$h_i > 2 \cdot (k+1)/n$ withdrawn (1, 8)	$R^2=0.555$; $R^2_{\text{adj}}=0.466$; $s=0.65$; $F=6$ ($p=0.001$); $R^2_{100}=0.254$; $s_{100}=0.96$; $F_{100}=0.81$ ($p=0.5518$); $\text{RMS}=0.4993$; $\text{APV}=0.3267$; $\text{TSE}=5$; $\text{APMSE}=0.0192$; $\% \text{PredErr}=34.5228$; $t_{\text{pp}}=-1.35 \cdot 10^{-14}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.698$; $\text{MAE}=0.490$; $\text{MAPE}=0.680$; $\text{SEP}=0.660$; $\text{REP}(\%)=41.023$
n=26	$D_i > 4/n$ withdrawn (1, 2, 5, 13, 15, 21, 30)	$R^2=0.858$; $R^2_{\text{adj}}=0.821$; $s=0.41$; $F=23$ ($p=1.87 \cdot 10^{-7}$); $R^2_{100}=0.767$; $s_{100}=0.52$; $F=13$ ($p=1.05 \cdot 10^{-5}$); $\text{RMS}=0.3267$; $\text{APV}=0.3831$; $\text{TSE}=3$; $\text{APMSE}=0.0192$; $\% \text{PredErr}=17.0907$; $t_{\text{pp}}=-2.48 \cdot 10^{-14}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.427$; $\text{MAE}=0.289$; $\text{MAPE}=0.294$; $\text{SEP}=0.529$; $\text{REP}(\%)=32.175$
Set4: $\hat{Y}(\text{MPmg}) = a + b_1 \times \text{RPCG} + b_2 \times \text{Q10} + b_3 \times \text{F}_{\text{H}_2\text{O}}$ where $\hat{Y}(\text{MPmg})$ = estimated mutagenic potencies for <i>M. gilvum</i> ; RPCG = (charge of the most positively charged atom) / (sum of total positive charge); Q10 = charges on position 10; $\text{F}_{\text{H}_2\text{O}}$ = desolvation free energy for waterA; a = intercept; b_i = regression coefficients		
n=29	whole dataset	$R^2=0.652$; $R^2_{\text{adj}}=0.610$; $s=0.41$; $F=16$ ($p=6.38 \cdot 10^{-6}$); $R^2_{100}=0.477$; $s_{100}=0.51$; $F=7$ ($p=0.0013$) $\text{RMS}=0.1610$; $\text{APV}=0.1776$; $\text{TSE}=4$; $\text{APMSE}=0.0064$; $\% \text{PredErr}=3.834$; $t_{\text{pp}}=8.80 \cdot 10^{-16}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.4091$; $\text{MAE}=0.1443$; $\text{MAPE}=3.2906$; $\text{SEP}=0.3866$; $\text{REP}(\%)=116.6703$
n=27	$h_i > 2 \cdot (k+1)/n$ withdrawn (10, 26)	$R^2=0.643$; $R^2_{\text{adj}}=0.596$; $s=0.64$; $F=14$ ($p=2.34 \cdot 10^{-5}$); $R^2_{100}=0.495$; $s_{100}=0.46$; $F=7$ ($p=0.0016$); $\text{RMS}=0.1401$; $\text{APV}=0.1557$; $\text{TSE}=5$; $\text{APMSE}=0.0061$; $\% \text{PredErr}=3.2149$; $t_{\text{pp}}=3.25 \cdot 10^{-14}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.399$; $\text{MAE}=0.125$; $\text{MAPE}=1.228$; $\text{SEP}=0.360$; $\text{REP}(\%)=89.281$
n=23	$D_i > 4/n$ withdrawn (10,13,16,26)	$R^2=0.568$; $R^2_{\text{adj}}=0.506$; $s=0.34$; $F=9$ ($p=4.38 \cdot 10^{-4}$); $R^2_{100}=0.407$; $s_{100}=0.41$; $F=4$ ($p=0.0145$); $\text{RMS}=0.1104$; $\text{APV}=0.1236$; $\text{TSE}=4$; $\text{APMSE}=0.0053$; $\% \text{PredErr}=2.6091$; $t_{\text{pp}}=-2.62 \cdot 10^{-15}$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.340$; $\text{MAE}=0.097$; $\text{MAPE}=2.390$; $\text{SEP}=0.318$; $\text{REP}(\%)=102.864$
Set5: $\hat{Y}(\text{pKI}) = b_1 \times ^2\text{AIC} + b_2 \times \text{NBR} + b_3 \times \text{NCA}$ where $\hat{Y}(\text{pKI})$ = estimated inhibitory activity against CA II isozyme; ^2AIC = average information content (order 2); NBR = number of benzene rings; NCA = number of C atoms; b_i = regression coefficients		
n=38	whole dataset	$R^2=0.586$; $R^2_{\text{adj}}=0.533$; $s=0.29$; $F=16$ ($p=8.79 \cdot 10^{-7}$); $R^2_{100}=0.532$; $s_{100}=0.31$; $F=13$ ($p=8.99 \cdot 10^{-6}$); $\text{RMS}=0.0816$; $\text{APV}=0.0880$; $\text{TSE}=5$; $\text{APMSE}=0.0024$; $\% \text{PredErr}=3.4285$; $t_{\text{pp}}=-0.0453$ ($p_{\text{pp}}=0.9641$); $\text{RMSE}=0.2856$; $\text{MAE}=0.0751$; $\text{MAPE}=0.1334$; $\text{SEP}=0.2778$; $\text{REP}(\%)=14.7242$
n=34	$h_i > 2 \cdot k/n$ withdrawn (C23, C24, C25, C32) $b_2 - p=0.1093$	$R^2=0.448$; $R^2_{\text{adj}}=0.380$; $s=0.29$; $F=8$ ($p=3.42 \cdot 10^{-4}$); $R^2_{100}=0.360$; $s_{100}=0.32$; $F=5$ ($p=4.12 \cdot 10^{-3}$); $\text{RMS}=0.0863$; $\text{APV}=0.0939$; $\text{TSE}=5$; $\text{APMSE}=0.0029$; $\% \text{PredErr}=3.0648$; $t_{\text{pp}}=0$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.2938$; $\text{MAE}=0.0787$; $\text{MAPE}=0.1307$; $\text{SEP}=0.2847$; $\text{REP}(\%)=14.5028$
n=37	$D_i > 4/n$ withdrawn (C8)	$R^2=0.597$; $R^2_{\text{adj}}=0.544$; $s=0.28$; $F=17$ ($p=8.60 \cdot 10^{-7}$); $R^2_{100}=0.541$; $s_{100}=0.31$; $F=13$ ($p=9.48 \cdot 10^{-6}$); $\text{RMS}=0.0810$; $\text{APV}=0.0875$; $\text{TSE}=5$; $\text{APMSE}=0.0025$; $\% \text{PredErr}=3.3326$; $t_{\text{pp}}=0$ ($p_{\text{pp}}=1$); $\text{RMSE}=0.2845$; $\text{MAE}=0.0744$; $\text{MAPE}=0.1321$; $\text{SEP}=0.2765$; $\text{REP}(\%)=14.5992$
Set6: $\hat{Y}(\text{HE-Mlog}(1/\text{MRC}_{50})) = b_1 \times \log P + b_2 \times E_{\text{tot}}$ where $\hat{Y}(\text{HE-Mlog}(1/\text{MRC}_{50}))$ = estimated toxicity on <i>Hydractinia echinata</i> ; $\log P$ = hydrophobicity; E_{tot} = total optimized energy; b_i = regression coefficients		
n=28	whole dataset	$R^2=0.631$; $R^2_{\text{adj}}=0.579$; $s=1.25$; $F=22$ ($p=2.81 \cdot 10^{-6}$); $R^2_{100}=0.550$; $s_{100}=1.42$; $F_{100}=15$ ($p=5.32 \cdot 10^{-5}$); $\text{RMS}=1.5644$; $\text{APV}=1.6761$; $\text{TSE}=4$; $\text{APMSE}=0.0626$; $\% \text{PredErr}=3.4705$; $t_{\text{pp}}=-0.0574$ ($p_{\text{pp}}=0.9546$); $\text{RMSE}=1.2507$; $\text{MAE}=1.4526$; $\text{MAPE}=2.4174$; $\text{SEP}=1.2274$; $\text{REP}(\%)=35.2585$
n=26	$h_i > 2 \cdot k/n$ withdrawn	$R^2=0.692$; $R^2_{\text{adj}}=0.638$; $s=1.19$; $F=27$ ($p=9.26 \cdot 10^{-7}$);

	(C8, C25) $b_1 - p > 0.05$	$R^2_{loo}=0.649$; $s_{100}=1.30$; $F=21$ ($p=6.28 \cdot 10^{-6}$); RMS=1.4097; APV=1.5182; TSE=4; APMSE=0.0613; %PredErr=3.0244; $t_{pp}=0.7156$ ($p_{pp}=0.4804$); RMSE=1.1873; MAE=1.3013; MAPE=2.3819; SEP=1.1633; REP(%)=32.9849
n=23	$D_i > 4/n$ withdrawn (C5, C8, C21, C25, C27) $b_1 - p > 0.05$	$R^2=0.674$; $R^2_{adj}=0.611$; $s=1.13$; $F=22$ ($p=9.72 \cdot 10^{-6}$); $R^2_{loo}=0.627$; $s_{100}=1.24$; $F_{loo}=17$ ($p=5.51 \cdot 10^{-5}$); RMS=1.2801; APV=1.3914; TSE=4; APMSE=0.0640; %PredErr=2.4588; $t_{pp}=1.5758$ ($p_{pp}=0.1267$); RMSE=1.1314; MAE=1.1688; MAPE=3.1657; SEP=1.1054; REP(%)=31.6591
Set7: $\hat{Y}(\log ED_{50}) = b \times DCW^3$ where \hat{Y} = estimated antiepileptic activities (dose at which 50% of individuals reach the desired effect); DCW^3 = descriptor calculated with Monte Carlo simulation [39]; b = regression coefficient		
n=51	whole dataset	$R^2=0.737$; $R^2_{adj}=0.717$; $s=0.21$; $F=140$ ($p=5.65 \cdot 10^{-16}$); $R^2_{loo}=0.729$; $s_{100}=0.21$; $F_{loo}=131$ ($p=1.89 \cdot 10^{-15}$); RMS=0.0427; APV=0.0435; TSE=3; APMSE=0.0009; %PredErr=3.2254; $t_{pp}=-0.1155$ ($p_{pp}=0.9085$); RMSE=0.2066; MAE=0.0418; MAPE=0.1116; SEP=0.2066; REP(%)=12.9209
n=	$h_i > 2 \cdot k/n$ withdrawn (none)	no h_i value higher than threshold was identified
n=48	$D_i > 4/n$ withdrawn (C2, C19, C26, C36, C46, C51)	$R^2=0.838$; $R^2_{adj}=0.816$; $s=0.15$; $F=228$ ($p=8.22 \cdot 10^{-19}$); $R^2_{loo}=0.835$; $s_{100}=0.16$; $F_{loo}=213$ ($p=1.75 \cdot 10^{-18}$); RMS=0.0230; APV=0.0235; TSE=3; APMSE=0.0005; %PredErr=2.2033; $t_{pp}=-0.1733$ ($p_{pp}=0.8632$); RMSE=0.1516; MAE=0.0225; MAPE=0.0892; SEP=0.1516; REP(%)=9.6954

R^2 = determination coefficient; R^2_{adj} = adjusted correlation coefficient; s = standard error of estimate; F =F-value (p = p-value); R^2_{loo} = determination coefficient in leave-one-out analysis; s_{100} = standard error of predicted; F_{loo} = Fisher's value and associated significance in leave-one-out analysis; RMS= residual mean square; APV= average prediction variance; TSE= total squared error; APMSE= average prediction mean squared error; %PredErr= prediction error; t_{pp} , p_{pp} = t-statistics for intercept and regression coefficients; RMSE= root-mean-square error; ME= mean error; MAE= mean absolute error; MAPE= mean absolute percentage error; SEP= standard error of prediction; REP(%)=relative error of prediction

157

158 Classification of QSAR models (full-model, model obtained after withdrawn of compound(s) with
159 h_i - h_i -model and respectively with D_i higher than thresholds - D_i -model) according to applied
160 validation statistics is presented in Figure 1.



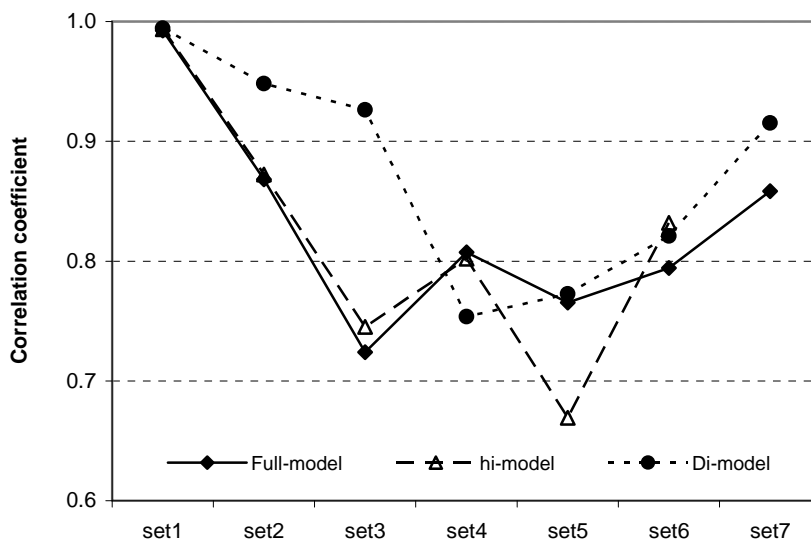
161

162 **Figure 1.** Full model & h_i -model & D_i -model: classification according to validation criteria.

163

164 The highest correlation coefficient was obtained in 5 cases out of 7 by the model after removal the
165 compounds with the Cook's distance higher than threshold. The full model obtained the higher

166 correlation coefficient in the fourth set, while the model obtained after removal of the compounds
 167 with leverage higher than threshold obtained the higher correlation coefficient in the sixth set. The
 168 evolution of correlation coefficients is presented in Figure 2.



169
 170 **Figure 2.** Full model - h_i model - D_i model: evolution of correlation coefficient
 171

172 Statistical significant increases in correlation coefficient have been identified in the second and
 173 third sets when both the full-model and the h_i -model were compared to D_i -model (Table 4). The
 174 Fisher's Chi-Square statistic (F-C-S) was applied to test if overall one model is better than other and
 175 the results are presented in Table 4.
 176

177 **Table 4.** Steiger's Z test for correlation coefficients comparisons and overall significance: results

Set	Full-model vs. h_i -model Z (p-value)	Full-model vs. D_i -model Z (p-value)	h_i -model vs. D_i -model Z (p-value)
set1	0.3760 (0.3535)	0.855 (0.1963)	0.4820 (0.3149)
set2	0.1040 (0.4586)	2.861 (0.0021)	2.7450 (0.0030)
set3	0.1740 (0.4309)	2.583 (0.0049)	2.3810 (0.0086)
set4	0.0560 (0.4777)	0.465 (0.3210)	0.4040 (0.3431)
set5	0.8100 (0.2090)	0.073 (0.4709)	0.8760 (0.1905)
set6	0.3840 (0.3505)	0.255 (0.3994)	0.1130 (0.4550)
set7	n.a.	1.312 (0.0948)	n.a.
F-C-S (p-value)	6.0139 (0.4216)	18.2757 (0.0108)	15.2359 (0.0185)

F-C-S = Fisher's Chi-Square statistic; p-value = probability

178
 179 The leave-many-out analyses were conducted on set1 and set2 to assess the usefulness of influential
 180 identification and withdrawn on the QSARs abilities. Characteristics of the obtained models are
 181 presented in Table 5.
 182

Table 5. Leave-many-out analysis: results.

Set	Split	n	R ²	F	p _F	full-model		D _i -model	
						n	R ²	F	p _F
set1	Training	40	0.9875	950	2.59·10 ⁻³⁴	38	0.9890	1020	2.32·10 ⁻³³
	Test	20	0.9802	223	2.89·10 ⁻¹³	16	0.9869	300	1.50·10 ⁻¹¹
set2	Training	53	0.7539	77	5.55·10 ⁻¹⁶	45	0.9097	211	1.18·10 ⁻²²
	Test	26	0.7609	33	1.58·10 ⁻⁷	21	0.8810	67	4.77·10 ⁻⁹

n = sample size; R² = determination coefficient;
 F = Fisher's statistics; p_F = significance of F statistics;

185 The plot of full-model versus D_i-model for set1 and set2 are presented in Figures 3 and 4.

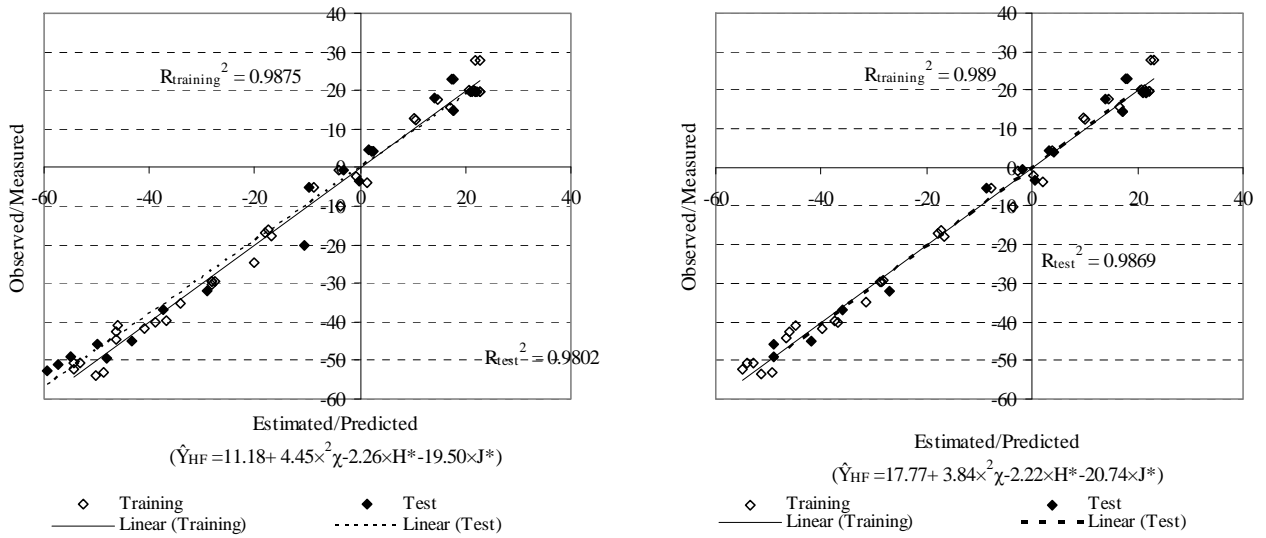


Figure 3. Set1 full-model (left-hand) vs. D_i-model (right-hand): observed/Measured vs.

estimated/predicted

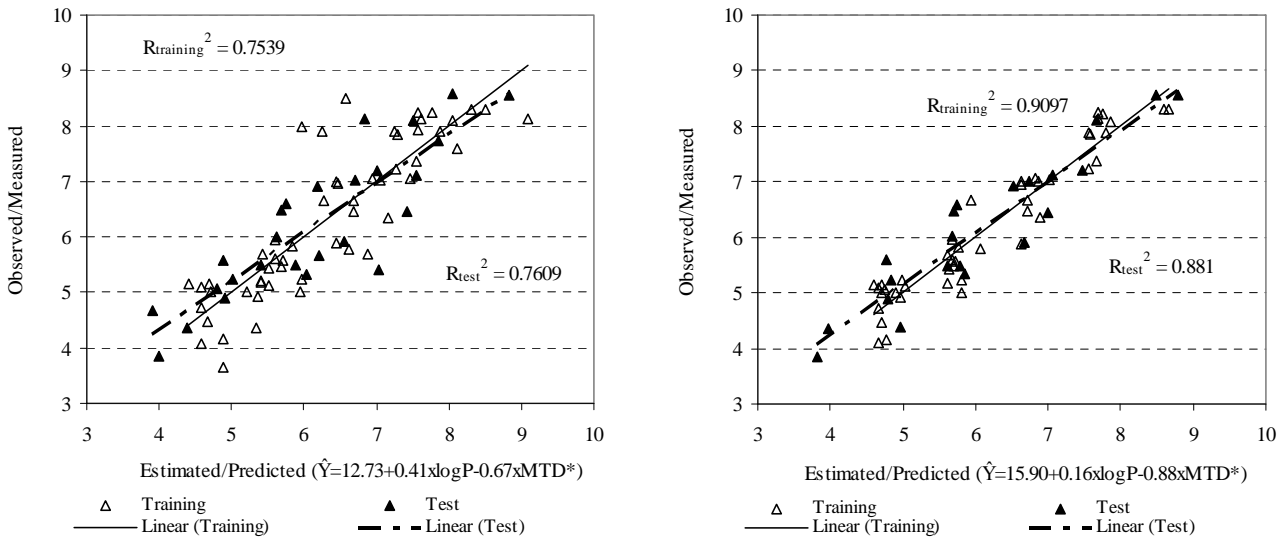


Figure 4. Set2 full-model (left-hand) vs. D_i-model (right-hand): observed/Measured vs.

estimated/predicted

192

193 The leave-one-out cross-validation determination coefficient for training set1 was of 0.9846 while
194 for training set set2 was of 0.7208 when full-models were investigated. The leave-one-out cross-
195 validation determination coefficient for training set1 was of 0.9870 while for training set2 was of
196 0.8975 when the D_i -models were investigated. A statistically significant increase of correlation
197 coefficient has been identified for the training set of the set2 in D_i -model compared to full-model (Z
198 = 2.609, p-value = 0.0045).

199

200 **Discussion**

201 The assessment of influential withdrawn using leverage and Cook's distance has successfully
202 accomplished. Seven data sets with sample sizes range from 28 (set6) to 79 (set2) were analyzed.
203 Three linear regression models were investigated for each set included in analyzes whenever
204 appropriate (full-model, h_i -model and D_i -model). The present study tried to answer to the following
205 research question: "Hat-matrix approach is more appropriate than Cook's distance approach to
206 identify influential in regression analysis?".

207 The analysis of the obtained results revealed that in 5 out of 7 sets the correlation coefficient
208 increase when both compounds with h_i and respectively D_i higher than thresholds were removed
209 (Table 3). The number of withdrawn compounds varied from 2 to 4 for h_i -model and from 1 to 13
210 for D_i -model (Table 3). In just few cases the same compound was identified as influential by both
211 leverage and Cook's methods: 1 compound (in set1, set2, and set3) and 2 compounds (in set set4
212 and set6).

213 Some independent variable proved not to have a statistically contribution to the model (see Table
214 3): h_i -model set5 (translated also to a lower determination coefficient compared to full-model) and
215 set6 and D_i -model set6 (the determination coefficient had a higher value for h_i -model compared to
216 D_i -model for set6). In these cases, it is correct to construct the models without those descriptors
217 identified with no statistically contribution to the model. With one exception represented by set4,
218 determination coefficients for D_i -models were higher than determination coefficients obtained in
219 full-models (Table 3 and Figure 2). The highest increase of determination coefficient was observed
220 in D_i -model of set3. The difference between determination coefficient and adjusted determination
221 coefficient varied from 0 to 0.088 (for full-model – set3), 0.089 (for h_i -model – set3) and 0.063 (for
222 D_i -model – set6). The difference between determination coefficient and its corresponding value in
223 leave-one-out analysis varied from 0.002 to 0.237 (full-model), 0.001 to 0.148 (h_i -model), and from
224 0.002 to 0.161 (D_i -model).

225 The analysis of validation statistics showed that D_i -models obtained systematically better results
226 compared to full-models (Table 3 and Figure 1). Furthermore, even if goodness-of-fit is not a good

227 statistics for model predictivity [40,41], no statistically significant differences between correlation
228 coefficients obtained in full-model compared to those obtained in h_i -models were identified (Table
229 4). However, the correlation coefficients obtained by D_i -models proved statistically significant
230 higher compared to those obtained in both full-model and h_i -model for set2 and set3 (Table 4).
231 Furthermore, the F-C-S statistic showed that overall, the D_i -model was better than both full-model
232 and h_i -model ($p < 0.05$, Table 4). The above-presented facts let to the conclusion that analysis of
233 influential should be conducted by applying the Cook's distance approach.

234 The external validation of the Cook's distance approach was furthermore assessed in leave-many-
235 out analysis on two datasets (set1 and set2), one with statistically increase of correlation coefficient
236 (set2) and one without statistically increase of correlation coefficient (set1). Similar results are
237 obtained when training and test sets are compared (Table 5). The significant increase of
238 determination coefficient in both training and test sets is transmitted also in leave-many-out
239 analyzes for the second dataset (set2), the increase being of 0.156 for training set and 0.120 for test
240 set. The spread of point in the plots of full-model and D_i -model is similar for set1 (Figure 3) but the
241 difference are obvious when set2 is investigated (Figure 4). A reliable and valid regression model
242 must look as set2 D_i -model not as set2 full-model (Figure 4).

243 Scientifically literature recommend not to trust a QSAR model when correlation coefficient is lower
244 than 0.6, which known to be is an insufficient condition for assessment of predictive power of a
245 model [42]. This analysis show that a determination coefficient < 0.6 could be significantly
246 improved with analyses and withdrawn of influential in order to obtain a model with good
247 performance in prediction (see Table 3, set3). In our opinion, the predictivity power of a model
248 stands in correct application of statistical methods to identify the QSAR models. Identification of
249 influential compound in data set could significantly improve the model and should be conducted
250 any time when a regression analysis is desired. Fit the model with and without the influential
251 compound(s) and look to the effect on regression characteristics (R^2 , R^2_{adj} , F-value (p -value), s ,
252 regression coefficients and their significance, validation criteria presented in Table 1) as well as on
253 the plot of the models. It is the task of a statistician to examine the influential compounds and to
254 identify important cases before presentation of results but this task could be done by any researcher
255 with experience in statistics. Based on the presented results, it is showed that Cook's distance
256 approach is more suitable to proper identification of influential in dataset and we recommend its
257 application in linear regression analysis for QSAR models. The leverage approach could be used on
258 the D_i -model to analyze the membership of compounds in the model to the structural model domain
259 [43].

260 Based on the results obtained in this study we recommend that either to accept (if leave-one-out,
261 leave-many-out analyses and external validation sustain the model) or to reject the QSAR model

262 obtained after removal of influential(s) and never accept a model that contains influential
263 compounds (their presence lead to instability of the QSAR model).

264
265

266 **Conclusion**

267 The use of leverage methodology led to improvement of QSAR models characteristic and
268 performances. Better QSAR models were obtained when Cook's distance approach was used
269 compared to both full-model and h_i -model. Cook's distance approach is recommended to be used to
270 identify influential compounds in dataset whenever the linear regression analysis for QSAR models
271 is applied.

272

273 **Conflict of Interest**

274 The authors declare that there is no conflict of interest.

275

276 **References**

- 277 [1] Eriksson, L.; Jaworska, J.; Worth, A.P.; Cronin, M.T.D.; McDowell, R.M.; Gramatica, P.
278 Methods for reliability and uncertainty assessment for applicability evaluations of classification-
279 and regression-based QSARs. *Environ. Health Persp.*, **2003**, *111*, 1361-1375.
- 280 [2] Walker, J.D.; Jaworska, J.; Comber, M.H.I.; Schultz, T.W.; Dearden, J.C. Guidelines for
281 developing and using quantitative structure-activity relationships. *Environ. Toxicol. Chem.*, **2003**,
282 *22*, 1653-1665.
- 283 [3] Tropsha, A.; Golbraikh, A. Predictive QSAR modeling workflow, model applicability
284 domains, and virtual screening. *Curr. Pharm. Des.*, **2007**, *13*, 3494-3504.
- 285 [4] Jaworska, J.S.; Comber, M.; Auer, C.; Van Leeuwen, C.J. Summary of a workshop on
286 regulatory acceptance of (Q)SARs for human health and environmental endpoints. *Environ. Health*
287 *Persp.*, **2003**, *111*, 1358-1360.
- 288 [5] Guidance Document on the Validation of (Quantitative)Structure-Activity Relationships
289 [(Q)SAR] Models. [online] [Accessed 31 May 2012]. ENV/JM/MONO (OECD Environment
290 Health and Safety Publications) 2007;2. Available from: URL:
291 [http://www.oecd.org/officialdocuments/displaydocumentpdf/?cote=env/jm/mono\(2007\)2&doclanguage=en](http://www.oecd.org/officialdocuments/displaydocumentpdf/?cote=env/jm/mono(2007)2&doclanguage=en)
292 age=en
- 293 [6] Tropsha, A. Best Practices for QSAR Model Development, Validation, and Exploitation.
294 *Mol. Inform.*, **2010**; *29(6-7)*, 476-488.
- 295 [7] Mekenyan, O.G.; Veith, G.D. Synergism between QSAR Modeling and Physico-Chemical
296 Principles. *SAR QSAR Environ. Res.*, **1995**, *4*, 155-165.

- 297 [8] Dearden, J.C.; Cronin, M.T.D.; Kaiser, K.L.E. How not to develop a quantitative structure-
298 activity or structure-property relationship (QSAR/QSPR). *SAR QSAR Environ. Res.*, **2009**, *20*(3-4),
299 241-266.
- 300 [9] Cronin, M.T.D.; Schultz, T.W. Pitfalls in QSAR. *J. Theoret. Chem. (Theochem)*, **2003**, *622*,
301 39-51.
- 302 [10] Bolboacă, S.D.; Jäntschi, L. Distribution Fitting 3. Analysis under Normality Assumptions.
303 *Bulletin of University of Agricultural Sciences and Veterinary Medicine Cluj-Napoca. Horticulture*,
304 **2009**, *66*(2), 698-705.
- 305 [11] Bolboacă, S.D.; Jäntschi, L. Modelling the property of compounds from structure: statistical
306 methods for models validation. *Environ. Chem. Lett.*, **2008**, *6*, 175-181.
- 307 [12] Abraham, M.H.; Sanchez-Moreno, R.; Cometto-Muniz, J.E. A quantitative structure-activity
308 analysis on the relative sensitivity of the olfactory and the nasal trigeminal chemosensory systems.
309 *Chem. Senses*, **2007**, *32*, 711-719.
- 310 [13] Yaffe, D.; Cohen, Y.; Espinosa, G.; Arenas, A.; Giralt, F. A fuzzy ARTMAP based on
311 quantitative structure-property relationships (QSPRs) for predicting aqueous solubility of organic
312 compounds. *J. Chem. Inf. Comput. Sci.*, **2001**, *41*, 1177-1207.
- 313 [14] Topliss, J.G.; Costello, R.J. Chance correlations in structure-activity studies using multiple
314 regression analysis. *J. Med. Chem.*, **1972**, *15*, 1066-1068.
- 315 [15] Kolmogorov, A. Confidence Limits for an Unknown Distribution Function. *Ann. Math. Stat.*,
316 **1941**, *12*(4), 461-463.
- 317 [16] Pearson, K. On the criterion that a given system of deviations from the probable in the case of
318 a correlated system of variables is such that it can be reasonably supposed to have arisen from
319 random sampling. *Philosophical Magazine*, **1900**, *50*, 157-175.
- 320 [17] Kleinbom, D.G.; Kupper, L.L.; Muler, K.E.; Nizam, A. *Applied Regression Analysis and*
321 *other Multivariate Methods*. Duxbury Press: Pacific Grove, **1998**.
- 322 [18] Tropsha, A.; Gramatica, P.; Gombar, V.K. The Importance of Being Earnest: Validation is
323 the Absolute Essential for Successful Application and Interpretation of QSPR Models. *QSAR*
324 *Comb. Sci.*, **2003**, *22*(1), 69-77.
- 325 [19] Gramatica, P. Principles of QSAR models validation: internal and external. *QSAR Comb.*
326 *Sci.*, **2007**, *26*(5), 694-701.
- 327 [20] Cook, R.D. Detection of Influential Observations in Linear Regression. *Technometrics*
328 *(American Statistical Association)*, **1977**, *19*(1), 15-18.
- 329 [21] Cook, R.D. Influential Observations in Linear Regression. *J. Amer. Statistical Assoc.*, **1979**,
330 *74*(365), 169-174.

- 331 [22] Bollen, K.A.; Jackman, R. *Regression diagnostics: An expository treatment of outliers and*
332 *influential cases*. In: *Modern Methods of Data Analysis*, Fox, J.; Scott, J. Long, Eds. Sage: Newbury
333 Park, **1990**, pp. 257-91.
- 334 [23] Steiger, J.H. Tests for comparing elements of a correlation matrix. *Psychol. Bull.*, **1980**, *87*,
335 245-251.
- 336 [24] Bolboacă, S.D. Assessment of Random Assignment in Training and Test Sets using
337 Generalized Cluster Analysis Technique. *Appl. Med. Inform.*, **2010**, *28(2)*, 9-14.
- 338 [25] Kohavi, R. *A study of cross-validation and bootstrap for accuracy estimation and model*
339 *selection*. In: Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence
340 **1995**, *2(12)*, pp. 1137-1143.
- 341 [26] Martens, H.A.; Dardenne, P. Validation and verification of regression in small data sets.
342 *Chemometr. Intell. Lab.*, **1998**, *44*, 99-121.
- 343 [27] Besalú, E.; de Julián-Ortiz, J.V. Equivalence of the Pecka-Ponec Correlation Probability and
344 the Statistical F Significance for MLR Models. *J. Math. Chem.*, **2004**, *36(4)*, 361-363.
- 345 [28] Mallow, C.L. Some comments on Cp. *Technometrics*, **1973**, *15*, 661-675.
- 346 [29] Gorman, J.W.; Toman, R.J. Selection of variables for fitting equations to data.
347 *Technometrics*, **1966**, *8*, 27-51.
- 348 [30] Tukey, J.W. Discussion. *J. R. Statisti. Soc.*, **1967**, *29*, 47-48.
- 349 [31] Milac, A.-L.; Avram, S.; Petrescu, A.-J. Evaluation of a neural networks QSAR method
350 based on ligand representation using substituent descriptors: Application to HIV-1 protease
351 inhibitors. *J. Mol. Graph. Model.*, **2006**, *25(1)*, 37-45.
- 352 [32] Fisher, R.A. Combining independent tests of significance. *Amer. Statistician*, **1948**, *2*, 30.
- 353 [33] Mercader, A.; Castro, E.A.; Toropov, A.A. Maximum Topological Distances Based Indices
354 as Molecular Descriptors for QSPR. 4. Modeling the Enthalpy of Formation of Hydrocarbons from
355 Elements. *Int. J. Mol. Sci.*, **2001**, *2*, 121-132.
- 356 [34] Duda-Seiman, C.; Duda-Seiman, D.; Dragos, D.; Medeleanu, M.; Careja, V.; Putz, M.V.;
357 Lacrama, A.-M.; Chiriac, A.; Nutiu, R.; Ciubotariu, D. Design of Anti-HIV Ligands by Means of
358 Minimal Topological Difference (MTD) Method. *Int. J. Mol. Sci.*, **2006**, *7*, 537-555.
- 359 [35] Kim, D.; Hong, S.-I.; Lee, D.-S. Triazoloquinazolines as Human A3 Adenosine Receptor
360 Antagonists: A QSAR Study. *Int. J. Mol. Sci.*, **2006**, *7*, 485-496.
- 361 [36] Kim, D.; Hong, S.-I.; Lee, D.-S. The Quantitative Structure-Mutagenicity Relationship of
362 Polycyclic Aromatic Hydrocarbon Metabolites. *Int. J. Mol. Sci.*, **2006**, *7*, 556-570.
- 363 [37] Eroglu, E. Some QSAR Studies for a Group of Sulfonamide Schiff Base as Carbonic
364 Anhydrase CA II Inhibitors. *Int J Mol Sci.*, **2008**, *9(2)*, 181-197.

- 365 [38] Chicu, S.A.; Putz, M.V. Köln-Timișoara Molecular Activity Combined Models toward
366 Interspecies Toxicity Assessment. *Int. J. Mol. Sci.*, **2009**, *10*, 4474-4497.
- 367 [39] Garro Martinez, J.C.; Duchowicz, P.R.; Estrada, M.R.; Zamarbide, G.N.; Castro, E.A. QSAR
368 Study and Molecular Design of Open-Chain Enaminones as Anticonvulsant Agents. *Int. J. Mol.*
369 *Sci.*, **2011**, *12*, 9354-9368.
- 370 [40] Novellino, E.; Fattorusso, C.; Greco, G. Use of comparative molecular field analysis and
371 cluster analysis in series design Original. *Pharm. Acta Helv.*, **1995**, *70*(2), 149-154.
- 372 [41] Norinder, U. Single and domain mode variable selection in 3D QSAR applications. *J*
373 *Chemomet.*, **1996**, *10*(2), 95-105.
- 374 [42] Golbraikh, A.; Tropsha, A. Beware of q²! *J. Mol. Graphics Mod.*, **2002**, *20*(4), 269-276.
- 375 [43] Atkinson, A.C. *Plots, Transformations and Regression*. Oxford: Clarendon Press, **1985**.