

1 **Estimation of Influential Points in any Data Set from**
2 **Coefficient of Determination and its Leave-one-out**
3 **Cross-validated Counterpart**

4
5 Gergely Tóth¹, Zsolt Bodai¹, Károly Héberger^{2,*}

6 ¹ Institute of Chemistry, Loránd Eötvös University,
7 Pázmány sétány 1/a, Budapest, H-1117 Hungary

8 ² Institute of Materials and Environmental Chemistry, Research
9 Centre of Natural Sciences, Hungarian Academy of Sciences,
10 Pusztaszeri út 59-67, Budapest, H-1025 Hungary

11
12
13
14
15
16
17
18
19
20
21 *corresponding author, E-mail: heberger.karoly@ttk.mta.hu,

22 phone:+36 1 438 11 03

23 fax:+36 1 438 11 43

1 ABSTRACT

2 Coefficient of determination (R^2) and its leave-one-out cross-validated analogue
3 (denoted by Q^2 or R^2_{cv}) are the most frequently published values to characterize the predictive
4 performance of models. In this article we use R^2 and Q^2 in a reversed aspect to determine
5 uncommon points, i.e. influential points in any data sets. The term $(1-Q^2)/(1-R^2)$ corresponds
6 to the ratio of predictive residual sum of squares (*PRESS*) and the residual sum of squares
7 (*RSS*). The ratio correlates to the number of influential points in experimental and random test
8 data sets. We propose an (approximate) *F*-test on $(1-Q^2)/(1-R^2)$ term to quickly pre-estimate
9 the presence of influential points in training sets of models. The test is founded upon the
10 routinely calculated Q^2 and R^2 values and warns the model builders to verify the training set,
11 to perform influence analysis or even to change to robust modeling.

12

13 KEYWORDS

14 Coefficient of determination,
15 leave-one-out cross-validation,
16 influence analysis,
17 quantitative structure activity relationships,
18 prediction,
19 training set

20

1. INTRODUCTION

Model validation and evaluation of predictive ability are basic steps in chemometrics, bioinformatics, quantitative structure activity relationship (QSAR) and quantitative structure retention relationship (QSRR). The coefficient of determination (R^2) and the leave-one-out cross-validated R^2 (Q^2 or R^2_{cv}) e.g. in ref. [1] are performance parameters calculated in most studies. In the last decades there is a plenty of discussion on the qualitative and the quantitative meaning of these parameters in the validation and prediction processes alike. There are other ways of calculations for performance parameters, e.g. they can be calculated on the training set and on the test set of the data [2-7]. The former is called internal validation, the latter is called external one. We can use the mean of the test set or of the training set in external Q^2 calculations [6,7]. Further functions can be defined, if we take into account the degrees of freedom of the sums of squares in the calculations [1]. In the case of Q^2 , most of the calculations are performed with the leave-one-out cross-validation method, but there are many examples for different number of data to leave out [3,8].

Though the interpretation of R^2 is usually straightforward, in the case of Q^2 , the interpretation is not unified or even dubious. Some authors only take into account the value of Q^2 to R^2 . If Q^2 is only “slightly” less than the corresponding R^2 , the model is considered to be validated [2,6]. However, the measure for “slight” difference cannot be given, especially not without the degree of freedom. Other users concentrate on the numerical value of Q^2 without the degree of freedom and without any comparison to R^2 . If it is larger than e.g. 0.5, the model is thought to be validated [2,9-11]. It is not necessary that Q^2 calculated on the training set correlates to the external predictive ability as it is stated in the article entitled “Beware of Q^2 ” written by Golbraikh and Tropsha on QSAR in 2002 [2]. Doweiko repeated this observation in his paper entitled “QSAR: dead or alive?” [12].

1 The literature on Q^2 is connected mostly on model validation and predictive ability.
2 Leave-one-out Q^2 on the training set is a measure of internal predictive power and it is not the
3 standalone best choice to quantify predictive performance in general, e.g. [4,13-15].

4 In this article we focus on a different aspect of R^2 and Q^2 . Originally, we tried to
5 develop a statistical test to be used in model validation, where the input data are R^2 and Q^2
6 calculated on the training set. We tried with different formulas, but none of them indicated
7 reliable correlation to the expected validity of the models. Looking through the calculation
8 details of R^2 and Q^2 , we realized that our methods were not connected to the validity or the
9 predictive ability of the models, but they were connected to a different feature of the training
10 set. Here, we suggest using a statistical test to pre-estimate the presence of influential points
11 in the training set. In our study we focus on the average model builders of QSAR or QSRR
12 ones, where ordinary or partial least square regressions are applied, and R^2 and Q^2 are
13 routinely calculated. Influence analysis, identification of x and y outliers, and comparison to
14 robust regression are usually outside of scope in the average QSAR/QSRR publication.
15 Therefore, an introduction of a method that alerts model builders is a valuable aim. There are
16 two general ways to investigate the data set and the model building for uncommon points (e.g.
17 ref. [16]). The first one is the regression diagnostics pioneered by Cook [17,18]. Here, the
18 model is fitted to the whole data set first, and thereafter the influential points, x and y outliers
19 are detected *via* different criteria [1,18]. The other way is to use robust regression, where the
20 model is built on a subset or on a weighted set of data. Here, the uncommon feature of the
21 points is taken into account in the model building. The outliers are quantified with large
22 robust residuals in the y direction and robust distances in the predictor space. The latter points
23 are usually termed as leverages. Since the aim of our study is to introduce a quick test to alert
24 uncommon points in the data set of QSAR studies, where the model building has been already
25 performed with ordinary or partial least square methods, we have limited ourselves to the first
26 type of regression diagnostics. It does not mean, that we question or neglect the results

1 obtained in the last decades with robust regression, simple the aim and the corresponding
 2 preconditions do not allow its use. The suggested test uses *PRESS* (predictive residual sum of
 3 squares) and its meaning is at least questionable in the combination of the leave-one-out
 4 method and robust regression.

5

6 2. THEORY

7 2.1 Calculation of sum of squares

8 We denote with *TSS*, *RSS* and *MSS* the total, residual and model sum of squares of *n*
 9 data, y_i . The average of the experimental data is denoted by \bar{y} and \hat{y}_i -s are the data
 10 calculated by a model. The number of the parameters (including intercept) in the model is *p*.

$$11 \quad TSS = \sum_{i=1}^n (y_i - \bar{y})^2, \quad RSS = \sum_{i=1}^n (\hat{y}_i - y_i)^2, \quad MSS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (\text{Eq. 1})$$

12 The coefficient of determination is defined as

$$13 \quad R^2 = \frac{MSS}{TSS} = 1 - \frac{RSS}{TSS} \quad (\text{Eq. 2})$$

14 because $TSS = MSS + RSS$.

15 In the case of internal cross-validation with leave-one-out method, we can calculate the
 16 predictive residual sum of square as:

$$17 \quad PRESS = \sum_{i=1}^n (\hat{y}_{i/i} - y_i)^2 \quad (\text{Eq. 3})$$

18 where $\hat{y}_{i/i}$ denotes the value calculated for the *i*-th experiment leaving out the *i*-th experiment
 19 in the parameterization of the model. The cross-validated correlation coefficient is defined in
 20 Eq.4.

$$21 \quad Q^2 = R_{cv}^2 = 1 - \frac{PRESS}{TSS} \quad (\text{Eq. 4})$$

1 $R^2 \in [0;1]$, but Q^2 can be negative, if the model performs weakly (worse than modeling with a
2 simple average), therefore $Q^2 \in (-\infty,1]$.

3 The basic assumption in our test is that the ratio of two variances sampling from the
4 same normal distribution follows F -distribution with the corresponding degree of freedom.
5 Both RSS and $PRESS$ are sum of squares with df_{RSS} and df_{PRESS} degrees of freedom. If our
6 data set (training set) is correctly chosen, we can reasonably expect that

7
$$\frac{PRESS / df_{PRESS}}{RSS / df_{RSS}} = \frac{(1 - Q^2) / df_{PRESS}}{(1 - R^2) / df_{RSS}} \approx F - distributed \quad (\text{Eq. 5})$$

8 Strictly speaking the F -distribution in (Eq. 5) is only valid when $PRESS$ and RSS are
9 independent. The $PRESS$ is higher or equal to RSS , i.e. they are not fully independent. Hence
10 we emphasize the approximate sign in Eq. (5).

11 Therefore, a traditional F -test (known also as variance ratio test) gives us information that the
12 models on the reduced data sets obtained by the leave-one-out way are significantly different
13 in the aspect of the variance from the one derived on whole data set. Of course it is not easy to
14 identify the direct link between the meaning of “difference in the aspect of the variance” and
15 “model validation”.

16 One of the reviewers suggested that the $PRESS/RSS$ test (i.e. the traditional parametric
17 F -test) might be substituted with a non-parametric alternative. However, we have not found
18 any reasonable algorithm (i.e. using bootstrap) for our case, where a given $PRESS/RSS$ ratio is
19 available from the literature. The bootstrap on the given data set provides very important (but
20 different) information. Namely, it provides the uncertainty (confidence interval, histogram) of
21 the $PRESS/RSS$ on the given data set. This issue is detailed in the result and discussion part
22 (3.3).

23

24 *2.2 Identification of influential points*

1 There is no unique mathematical definition of an influential observation in the
2 literature, therefore we used the following “... compared to other observations it has a
3 relatively large impact on the estimated quantities like response, regression coefficient,
4 standard error, etc.” [1]. One or more parameters are extremely sensitive to the influential
5 observation. If we omit the observation, there is a reasonable difference in the parameter set
6 causing different models. Outliers and influential points are similar but not identical concepts.
7 Many of the outliers are influential points as well, despite that they are outliers in the y
8 direction (termed often as outliers) or in the x direction (known as leverages). There are
9 mathematical definitions for outliers, there are methods to detect them despite the masking
10 effect, but it is a mismatch to use the outlier definitions for influential observations.

11 To identify the influential points in data sets, we selected some basic methods. A good
12 survey of the methods was published in 1986 [18]. A comparison of some new methods to
13 robust methods was performed recently, as well [16]. As we mentioned earlier, we did not use
14 all available methods to identify influential points, because we concentrated on the studies,
15 where Q^2 and R^2 are calculated and the model can be obtained with ordinary least squares and
16 partial least square regressions. Robust methods are very efficient to build models with
17 correct treatise of outliers and leverages, but an average QSAR or QSRR developer avoids
18 using robust methods. The aim of our study was to develop a quick pre-estimation tool on
19 uncommon data points for average model builders, who are not interested in model building
20 with special knowledge on robust statistics or influence analysis. Therefore, we deliberately
21 chose several non-robust methods being differently sensitive on the presence of influential
22 points in dissimilar data sets. Here, we outline the selected ones as described in the handbook
23 of Frank and Todeschini [1]. For details see the references therein and the surveys mentioned
24 in refs. [16,18]. We performed some calculations with robust methods as well (e.g. [16,19]),
25 but the results are questionable for the aim of the study and because of the combination of
26 ordinary least squares and robust regression.

1 HD: The i -th observation is called influential point, if the corresponding diagonal
 2 element of the hat matrix (h_{ii}) is

$$3 \quad h_{ii} > 2p/n \quad (\text{Eq. 6})$$

4 The hat matrix is calculated as $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ for ordinary least square regression, where \mathbf{X} is
 5 the predictor matrix. Since h_{ii} is proportional to the square of the Mahalanobis distance of the
 6 observation from the mean of observations, this definition often used to identify leverages, as
 7 well.

8 SR: t_i denotes the studentized residuals and calculated as

$$9 \quad t_i = \frac{r_i}{s_{/i}\sqrt{1-h_{ii}}} \quad (\text{Eq. 7})$$

10 where $r_i = \hat{y}_i - y_i$. $s_{/i} = \sqrt{\frac{(n-p)s^2 - r_i^2/(1-h_{ii})}{n-p-1}}$ means the standard error without the i -th
 11 observation. The i -th observation is influential, if $t_i > 2$.

12 COOK: Cook's method is used for regression, where r_{si} denotes the standardized
 13 residuals and s_r is the residual standard deviation. In this method a d_i value is defined and
 14 tested in an F -test with p and $n-p$ degrees of freedom.

$$15 \quad COOK_i = \frac{r_{si}^2 h_{ii}}{p(1-h_{ii})} \quad (\text{Eq. 8})$$

16 Practically, the F -test can be replaced by the comparison of $COOK_i$ to different limit values
 17 (constants). We defined influential points as COOK-1, if $COOK_i > 1$ and COOK-2, if
 18 $COOK_i > 4/n$ [17]. According to the classification of Chatterjee and Hadi [18], Eq. 8 belongs
 19 to the influence function type definitions.

20 COVRATIO: The covariance ratio method measures the influence of the i -th
 21 observation on the variance of the regression coefficients.

$$22 \quad COVRATIO_i = \left(\frac{s_{/i}}{s}\right)^{2p} \frac{1}{1-h_{ii}} = \frac{1}{(1-h_{ii})\left[(n-p-1)/(n-p) + t_i^2/(n-p)\right]^2} \quad (\text{Eq. 9})$$

1 We used the definitions of influential points with $|COVRATIO_i - 1| > 3p/n$ [20]. It differs
 2 from the (maybe mistyped) definition in the book of Frank and Todeschini [1]. Eq. 9 is
 3 related to the volume of confidence ellipsoids according to the classification of ref. [18].

4 DFBETAS is calculated using the b_j estimated regression coefficient, its $b_{j/i}$ estimation
 5 when the i -th experiment is omitted and c_{jj} is the diagonal of the $(\mathbf{X}^T \mathbf{X})^{-1}$ matrix.

$$6 \quad DFBETAS_{ij} = \frac{b_j - b_{j/i}}{s_{/i} \sqrt{c_{jj}}} \quad (\text{Eq. 10})$$

7 The i -th data is treated as influential observation, if $DFBETAS_{ij} \geq 2/\sqrt{n}$ [20]. The definition
 8 is related to partial influence [18].

9 DFFITS: The scaled variable is the difference between the predicted and the response of
 10 the i -th observation with and without using the observation in the model. It is scaled by the
 11 standard error of the observations.

$$12 \quad DFFITS_i = t_i \left(\frac{h_{ii}}{1 - h_{ii}} \right)^{1/2} \quad (\text{Eq. 11})$$

13 A point is influential, if $DFFITS_i > 2\sqrt{p/n}$ (case DFFITS-1) or $DFFITS_i > 2$ (case DFFITS-
 14 2). It belongs to the influence function type definitions [18].

15

16 3. RESULTS AND DISCUSSIONS

17 3.1 Simulated data sets

18 Data sets were generated using random numbers to mimic experimental data to be
 19 modeled with multivariate ordinary least squares regression. The superscripts denote the
 20 dimension of the variables. At first \mathbf{y}^p and $\mathbf{X}^{p \times (p-1)}$ variables were filled with uniform random
 21 numbers of $[0;1)$. $\mathbf{X}^{p \times (p-1)}$ was extended with a p -th column containing ones. The solution of
 22 the set of linear equations $\mathbf{X}^{p \times (p-1)} \mathbf{p}^p = \mathbf{y}^p$ provided a set of regression parameters \mathbf{p}^p , where p_p
 23 was the intercept in the regression. The number of the rows of \mathbf{X} was extended from p to n

1 and the new rows were filled with random numbers from uniform $[0;1)$ distribution. The last
2 column of \mathbf{X} contained only ones. The dimension of the column vector \mathbf{y} was extended from p
3 to n . The new elements were calculated using the previously obtained \mathbf{p}^p regression
4 parameters and the generated new rows of $\mathbf{X}^{n \times p}$ using the equation $\mathbf{y}^n = \mathbf{X}^{n \times p} \mathbf{p}^p$. Finally, a
5 white noise of $w \cdot \varepsilon$ was added to each elements of \mathbf{y} , where ε was a random number chosen
6 from standard normal distribution and w was a predefined factor. Nine-9 parameter sets were
7 used because of practical reasons: $n=10, p=5$; $n=20, p=5$; $n=20, p=10$; with combinations of
8 $w=0.05$, $w=0.10$ and $w=0.25$ weights of white noise; 10^5 random model calculations were
9 performed for each parameter sets resulted all together $9 \cdot 10^5$ datasets. The limit correlation
10 coefficients for chance correlation for $n=20$ (or $n=10$) is 0.444 (or 0.632) at the 5 % level
11 according to the Table C-3 of Bevington [21]; i.e. the medium range for correlation
12 coefficients were used, where the distortions can be effectively observed. Such a way the data
13 sets will contain outliers, influential points randomly.

14 If we used an F -like test for the ratio defined in Eq. 5, we need to know the degrees of
15 freedom both for RSS and $PRESS$. In the case of ordinary least square regression $df_{RSS}=n-p$.
16 We found in the literature that $df_{PRESS}=n-p$ is used without any proof or explanation. To test
17 this we determined df_{PRESS} numerically. $PRESS$ is a sum of squares of residual quantities. If
18 we accept the reasonable assumption of OLS regression that the $PRESS$ residuals are not
19 serially correlated and they are normally distributed, their sum of squares shows χ^2
20 distribution. The shape of the χ^2 distribution functions can be used to determine the degrees of
21 freedom [22] as it depends strongly on them. We calculated the histograms of our $PRESS$ -s
22 (10^5 $PRESS$ -s for each parameter set). We scaled the histograms with their standard
23 deviations. Thereafter we calculated the overlap integral of the scaled and normalized
24 histograms and theoretical χ^2 -distributions with different degrees of freedom. The maximal
25 overlap integrals (0.96-0.98) were obtained for the theoretical distributions with $n-p-1$ or $n-p$

1 degrees of freedom for all of the nine cases. The results encouraged us to use $df_{PRESS}=n-p$. It
2 has the advantage of simplifying Eq. 5, as well.

3 The number of influential points for all the $9 \cdot 10^5$ datasets were calculated with the
4 methods detailed above. Different number of influential points was provided according to the
5 different definitions. We found good correlation among the number of the influential points
6 and the *PRESS/RSS* ratios for the methods SR, COOK-1, COOK-2, DFBETAS, DFFIT-1 and
7 DFFIT-2. We did not detect reliable correlation for the method called HD, and we got
8 negative correlation for the COVRATIO one. This negative correlation is not surprising,
9 because a large SR value causes small COVRATIO, especially, if the HD method did not
10 seem to be decisive for our data sets. The lack of positive correlation of the HD and the
11 COVRATIO methods mean that leverage points are not necessarily influential observations,
12 because these quantities are suitable (only) to identify leverages and not influential
13 observations. HD is related to the Mahalanobis distance of the corresponding point to the
14 centre of the points, and COVRATIO is related to the volume of the confidence ellipsoids
15 [18]. We performed calculations, where *PRESS/RSS*-s were calculated by ordinary least
16 squares fit and the leverage points were detected by extreme Mahalanobis distances or by
17 extreme robust distances after different robust regression methods. In these cases we did not
18 observe correlation between the number of leverage points and *PRESS/RSS* values, similarly
19 to non-correlation with the HD and COVRATIO terms. We calculated also the correlation
20 between *PRESS/RSS* from ordinary least squares regression and the number of the omitted or
21 down-weighted observations in robust regressions, but we did not observe any significant
22 correlations. Without going into details and repeating all specific aspects of robust regression,
23 the lack of correlation in these incoherent comparisons can be caused by the differences in the
24 definitions of Euclidean and Mahalanobis distances, by the so-called masking effect and
25 differences in the breakdown points.

1 We calculated the relative frequencies of the number of data sets with different number
2 of influential points versus the *PRESS/RSS* of the data sets. In Figure 1 the relative
3 frequencies are shown for the COOK-2 method in the case of $n=10$, $p=5$ and $w=0.05$
4 parameter set. This method identifies for the most data sets 2-4 influential points. This range
5 is not surprising due to relatively small n/p ratio. It is also known that there is some
6 connection between the expected number of influential observations and the number of the
7 parameters in a model [22]. Data sets with larger number of influential points (defined as
8 COOK-2) had larger *PRESS/RSS* values. We plotted three percentiles of *F*-distributions for
9 85%, 90% and 95% with $\nu_1 = 5$ and $\nu_2 = 5$ degrees of freedom. Three or four influential points
10 were in the data sets, if the *PRESS/RSS* ratio was higher than the 90% percentile. Five
11 uncommon points were found, if *PRESS/RSS* was larger than the 95% percentile.

12

13 (Figure 1)

14

15 The results of an even more sensitive method can be seen in Figure 2. The DFBETAS
16 method identified 3-6 influential points for the most of the cases in the same set ($n=10$, $p=5$
17 and $w=0.05$) 7 and more influential points were found mostly with *PRESS/RSS* larger than the
18 90% percentiles. The lack or the small number of influential points depended on the
19 *PRESS/RSS* as well. Zero to two influential points were found mostly for *PRESS/RSS* smaller
20 than the 95% percentiles.

21

(Figure 2)

22 3.2 Experimental data

23 We tested the method on the results of Zhang *et al.* [23]. They performed a quantitative
24 structure retention relationship (QSRR) study on the gas chromatographic retention indices
25 using molecular descriptors. They built a multivariate regression model on the experimental
26 retention data of 161 hydrocarbons using a constant and two descriptors: the total number of

1 non-H bonds (constitutional descriptor), R autocorrelation of lag 3 weighted by atomic van
2 der Waals volumes (GETAWAY descriptor) [24].

3 In order to test the relation between *PRESS/RSS* and the number of influential points we
4 performed resampling on their data. We chose $n=10$ or $n=20$ molecules. We performed the
5 regression with $p=3$ parameters. We calculated *PRESS*, *RSS* and the number of the influential
6 points with the different methods. We repeated the random resampling for 10^5 cases both for
7 $n=10$ and $n=20$.

8 (Figure 3)

9 The relative frequencies of the number of data sets versus *PRESS/RSS* are shown for the
10 DFFITS-1 method ($n=10$, $p=3$) in Figure 3. There was zero or one influential point in the
11 majority of the resampled sets. Three influential points were seldom found and it coincided
12 mostly with *PRESS/RSS* values larger than the 85% percentage of the corresponding *F*-
13 distribution. The results are shown for the COOK-1 method on the same sets in Figure 4. This
14 method selected only few influential points, the most of the data sets were without any
15 influential points. Two influential points were found mostly for data sets with larger
16 *PRESS/RSS* than the 90% percentile of the corresponding *F*-distribution. This percentile
17 served also as an upper limit for the data sets with zero influential point.

18 (Figure 4)

19 We applied the f_i coefficient e.g. in ref. [1] to quantify the correlation among
20 *PRESS/RSS* and influential point methods. It is defined as:

$$21 \quad f_i = \frac{ad - bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}, \quad (\text{Eq. 12})$$

22 where the meaning of a , b , c and d is detailed in Table 1.

23 (Table 1)

24 In this calculation we distinguished according to the absence/presence of influential
25 points, but it was an approximation in the case of our simulated data, because methods

1 identifying the influential points observed many influential points in most of the cases. We
2 did not predefine an F_{crit} value, but we scanned the possible $PRESS/RSS$ range (x-axis of
3 Figures 1-4) to find an F_{crit} which maximizes the f_i coefficient. The ranges of the obtained
4 maximal f_i -s for the different parameter sets in the scanned F_{crit} range were as follows: SR
5 (0.1-0.2), COOK-1 (0.6-0.8), COOK-2 (0.5-0.6), DFBETAS (0.4-0.5), DFFITS-1 (0.2-0.3),
6 DFFITS-2 (0.3-0.4). We note again, that the absence/presence criterion fails due to the high
7 number of the influential points found by the methods for the most data sets.

8 In Figures 1-4 we showed that higher number of influential points (L) causes shift of the
9 relative frequency curves along the $PRESS/RSS$ axis and that the $PRESS/RSS$ ratio positively
10 correlates with the number of the influential points. If the correlation is strong, there is a
11 possibility to use $PRESS/RSS$ to detect, or at least to pre-estimate the presence of influential
12 points. An F -test on $PRESS$ and RSS may be used for this purpose, but we have to predefine a
13 significance level and a corresponding F_{crit} value. A reasonable significance level can be
14 identified, if we search the F_{crit} value, where the $PRESS/RSS$ and one of the identification
15 methods shows the maximal f_i . It means an F_{crit} value, where the separation of data sets with
16 and without influential points is maximal. We show the F_{crit} dependence of f_i for the gas
17 chromatographic retention data of Zhang *et al.* [23]. The COOK-1 method was chosen for the
18 detection of influential points. It can be seen in Figure 5, that a clear maximum is obtained at
19 $F_{crit}=2.4$ here. $F_{crit}=2.4$ corresponds to a percentile of 86%.

20 (Figure 5)

21 We collected Q^2 and R^2 values of 247 QSAR models from the literature. The sources of
22 the data were mostly collections, for QSAR details see references [6,10,25] and references
23 therein. We calculated the $PRESS/RSS$ ratios for these data and also the F -percentiles, because
24 n and p were accessible in the data collections. The histogram of the F -percentiles is shown in
25 Figure 6. Obviously, there are no data less than 0.5, because $1 \leq PRESS/RSS$ and $v_1=v_2$ cause
26 a minimal 0.5 percentile of the corresponding F -distribution. There were no influential points

1 in the training sets of the two thirds of the models according to our test, but one third of the
2 QSAR models would fail on an F -test of the $PRESS/RSS$ ratio. It means the training sets of
3 these models probably contained influential points. It can be interpreted that there were
4 problems already with the internal predictive character of the models. We note here that the
5 part of the models were taken in the collections of references [6,10] to show the existence of
6 better models. We note as well that the most of the models with F -percentiles larger than 0.95
7 (models with larger probability of influential observations) were 3D QSAR ones collected or
8 calculated by Cramer and Wendt [10].

9 (Figure 6)

10 The determination of the number of the degrees of freedom is not straightforward in the
11 case of partial least square regression (PLS). There are different assumptions and methods to
12 calculate so-called pseudo degrees of freedom for PLS regression [22,26,27]. Unfortunately,
13 we were not able to calculate pseudo degrees of freedom for these cases with PLS, because it
14 needs more details on the data sets and the models, than it was accessible in the used literature
15 sources of Q^2 and R^2 . Anyway, we plot a second histogram in Figure 6 (frequencies against
16 percentiles of the F distribution ($p=0.05$), where the pseudo degrees of freedom of the model
17 was defined as $4*p$ causing $df_{PRESS}=n-4p$ and $df_{RSS}=n-4p$. The factor 4 was chosen as an
18 extremely large difference between conventional degrees of freedom and pseudo degrees of
19 freedom. Figure 6 clearly shows that though the majority of the models are acceptable about
20 80 models are wrong (percentile is above 95%). The ambiguity problem of degrees of
21 freedom in case of PLS (or principal component regression) cannot cause a serious limitation.
22 The problem disappears asymptotically as ' $n-p$ ' approximates n , whereas there is some
23 uncertainty in the p -value only. The test detected the same models as wrong ones even if the
24 degree of freedom value was calculated by other multiplier than 4.

25 3.3 The uncertainty of $PRESS/RSS$ data

1 We used the bootstrap method of resampling residuals [28] to assess the uncertainty and
2 the confidence intervals of PRESS/RSS calculation. We generated 500 bootstrap samples for
3 each of 500 random datasets corresponding to our test sets with given n , m and w . The
4 bootstrap PRESS/RSS averages, the corresponding 2.5 and 97.5 percentiles are plotted versus
5 traditional PRESS/RSS values of the data sets (Figure 7).

6 (Figure 7)

7 In the case of low PRESS/RSS values, the bootstrap means are usually larger than the
8 standard one, while at medium and large PRESS/RSS they are smaller. The correlation of the
9 bootstrap averages and the standard ones are strong with less than unit slope. It means, the use
10 of bootstrap average PRESS/RSS enhances the conservative feature of our proposed test. The
11 lower and the upper confidence limits depend strongly on the datasets and they provide
12 rather large uncertainty.

13 4. CONCLUSIONS

14 The *PRESS/RSS* ratio calculated from leave-one-out Q^2 and R^2 correlates well with the
15 number of influential points in the training sets. Different identification methods on both
16 simulated and experimental data support the conclusion. The correlation is strong enough, so
17 we suggested a variance ratio test on the *PRESS/RSS* ratio to pre-estimate the presence of
18 influential points in the training set, if degrees of freedoms (df_{PRESS} and df_{RSS}) are known.
19 Some ambiguity in the degrees of freedom does not limit the applicability, because the test is
20 conservative in this sense: i.e. it will detect only the “largely” contaminated models as wrong
21 ones. However, any leave-one-out at a time diagnostic will fail, if influential points are shown
22 up in groups (e.g. in pairs).

23 There are two possible applications of our results. Q^2 , R^2 , n and p are usually calculated
24 and published in modeling, especially in QSAR studies. The rapid calculation of *PRESS/RSS*
25 and the F -test on it is a fast method to pre-estimate the presence of influential points or with
26 other words the internal predictive character of a model. If a model fails in this test, it is

1 worthwhile to consider changes in the training data. As many fortuitous QSAR models appear
2 in the literature, editors and reviewers can check the submitted models easily: if a model fails
3 the above variance ratio test the model has little generalization ability, if at all.

4 The other possibility is to apply the method for influential point detection, where not
5 specific data is declared as an influential one, but the whole set is marked as influential point
6 free or infected one. Of course, a hypothesis is necessary for the F -test. Our examples
7 suggested using 85-95 % percentiles as critical F values to make decisions between the H_0
8 hypothesis of influential observation free or H_a alternative hypothesis of presence of
9 influential points.

10

1 REFERENCES

- 2 [1] I.E. Frank, R. Todeschini, *The data analysis handbook*, first ed., Elsevier, Amsterdam,
3 1994.
- 4 [2] A. Golbraikh, A. Tropsha, *J. Mol. Graph. Model.* 20 (2002) 269–276.
- 5 [3] R. Todeschini, V. Consonni, A. Mauri, M. Pavan, *Anal. Chim. Acta* 515 (2004) 199–208.
- 6 [4] H. Kubinyi, in I. Yalcin, E. Aki Sener (Eds) *QSAR & Molecular Modelling in Rational*
7 *Design of Bioactive Molecules (Proceedings of the 15th European Symposium on*
8 *QSAR & Molecular Modelling, Istanbul, Turkey, 2004)*, CADD Society, Ankara,
9 2006, pp. 30–33.
- 10 [5] V. Consonni, D. Ballabio, R. Todeschini, *J. Chem. Inf. Model.* 49 (2009) 1669–1678.
- 11 [6] P.P. Roy, S. Paul, I. Mitra, K. Roy, *Molecules* 14 (2009) 1660–1701.
- 12 [7] V. Consonni, D. Ballabio, R. Todeschini, *J. Chemometr.* 24 (2010) 194–201.
- 13 [8] A.T Manvar, R.R.S. Pissurlenkar, V.R. Virsodia, K.D. Upadhyay, D.R. Manvar, A.K.
14 Mishra, H.D. Acharya, A.R. Parecha, C.D. Dholakia, A.K. Shah, E.C. Coutinhi, *Mol.*
15 *Divers.* 14 (2010) 285–305.
- 16 [9] A. Golbraikh, M. Shen, Z. Xiao, Y.D. Xiao, K.H. Lee, A. Tropsha, *J. Comput. Aided Mol.*
17 *Des.* 17 (2003) 241–253.
- 18 [10] R.D. Cramer, B. Wendt, *J. Comput. Aided Mol. Des.* 21 (2007) 23–32.
- 19 [11] E. Jiménez-Contreras, D. Torres-Salinas, R. Bailón-Moreno, R. Ruiz-Baños, E. Delgado-
20 López-Cózar, *Scientometrics* 79 (2008) 201–218.
- 21 [12] A. M. Doweiko, *J. Comput. Aided Mol. Des.* 22 (2008) 81–89.
- 22 [13] N. Chirico, P. Gramatica, *J. Chem. Inf. Model.* 51 (2011) 2320-2335.
- 23 [14] N. Chirico, P. Gramatica, *J. Chem. Inf. Model.* 52 (2012) 2044-2058.
- 24 [15] K. Roy, I. Mitra, P.K. Ojha, S. Kar, R.N. Das, H. Kabir, *Chemom. Intell. Lab. Syst.* 118
25 (2012) 200-210.
- 26 [16] A. Bagheri, H. Midi, M. Ganjali, S. Eftekhari, *Appl. Math. Sci.* 4 (2010) 1367-1386.

- 1 [17] D.R. Cook, S. Weisberg, Residuals and influence regression, Chapman & Hall, New
2 York, 1982.
- 3 [18] S. Chatterjee, A.S. Hadi, 1 (1986) 379-416.
- 4 [19] P. Rousseeuw, M. Hubert, in: Y. Dodge (Ed.), Lab statistical procedures and related
5 topics: Papers from the 3rd International Conference on Lab-Norm Related Methods
6 Neuchatel 1997, Ins. Math Stat., Hayward, 1997, pp. 201–214.
- 7 [20] D.A. Belsley, E. Kuh, R.E. Welsch, R. E., Regression Diagnostics: Identifying influential
8 data and sources of collinearity. John Wiley, New York, 1980.
- 9 [21] P.R. Bevington, Data Reduction and Error Analysis for the Physical Sciences, McGraw-
10 Hill Book Co., New York, 1969
- 11 [224] H. van der Voet, J. Chemometr. 13 (1999) 195–208.
- 12 [232] X. Zhang, L. Ding, Z. Sun, L. Song, T. Sun, Chromatographia 70 (2009) 511–518.
- 13 [243] I.V. Tetko, J. Gasteiger, R. Todeschini, A. Mauri, D. Livingstone, P. Ertl, V.A.
14 Palyulin, E.V. Radchenko, N.S. Zefirov, A.S. Makarenko, V.Y. Tanchuk, V.V.
15 Prokopenko, J. Comput. Aid. Mol. Des. 19 (2005) 453–63.
- 16 [254] J.C. Dearden, T.I. Netzeva, J Pharm. Pharmacol. 56 (2004) Suppl: S–82. 53, and
17 http://www.ukqsar.org/slides/John_Dearden2004.pdf.
- 18 [265] H.A. Seipel, J.H. Kalivas, J. Chemometr. 18 (2004) 306–311.
- 19 [276] L. Zhang, S. Garcia-Munoz, Chemometr. Intell. Lab. Syst. 97 (2009) 152–158.
- 20 [28] J. Fox, Applied Regression Analysis and Generalized Linear Models, second ed., SAGE
21 Publications, Thousand Oaks, 2008.
- 22

1

2

Table 1

3

Frequencies a , b , c and d denote the number of occurrences of the sub cases where L is the

4

number of the influential points in the data set

	$0 < L$	$L = 0$
$F_{\text{crit}} \leq \text{PRESS}/\text{RSS}$	a	b
$\text{PRESS}/\text{RSS} < F_{\text{crit}}$	c	d

5

6

FIGURE CAPTIONS

1

2 Figure 1

3 Relative frequencies of the number of data sets with different number of influential points (L)
4 identified by the COOK-2 method versus the $PRESS/RSS$ of the data sets. The black squares
5 denote three percentiles of the corresponding F-distribution. Parameters: $n=10$, $p=5$ and
6 $w=0.05$

7 Figure 2

8 Relative frequencies of the number of data sets with different number of influential points (L)
9 identified by the DFBETAS method versus the $PRESS/RSS$ of the data sets. The black squares
10 denote three percentiles of the corresponding F-distribution. Parameters: $n=10$, $p=5$ and
11 $w=0.05$

12 Figure 3

13 Relative frequencies of the number of resampled experimental sets with different number of
14 influential points (L) identified by the DFFITS-1 method versus the $PRESS/RSS$ of the data
15 sets. The black squares denote three percentiles of the corresponding F-distribution.
16 Parameters: $n=10$, $p=3$

17 Figure 4

18 Relative frequencies of the number of resampled experimental sets with different number of
19 influential points (L) identified by the COOK-1 method versus the $PRESS/RSS$ of the data
20 sets. The black squares denote three percentiles of the corresponding F-distribution.
21 Parameters: $n=10$, $p=3$

22 Figure 5

23 Dependence of f_i on the choice of the critical F -value for the resampled experimental data of
24 Zhang et al. $n=10$, $p=3$.

25

26 Figure 6

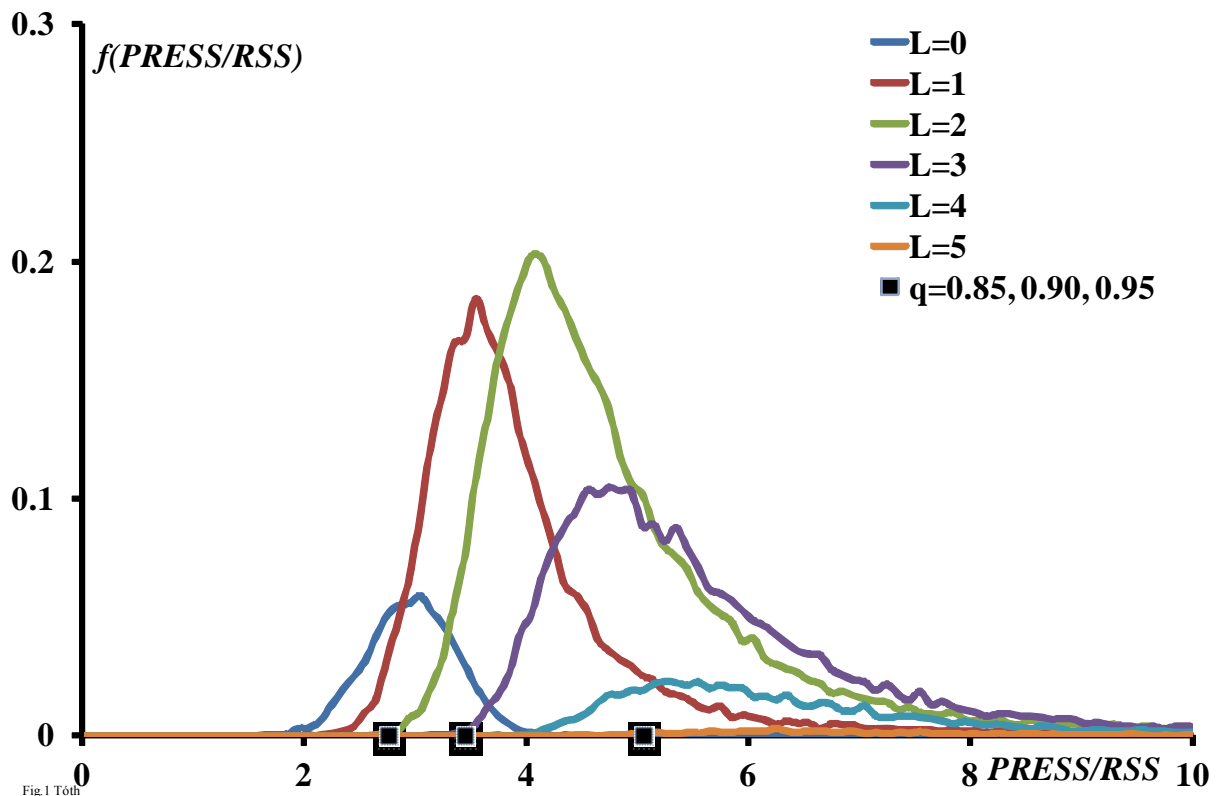
1 Frequencies for the percentiles calculated from F-distributions at the PRESS/RSS values of
2 247 QSAR models. Blue (black): df_{PRESS} and $df_{PRESS}=n-p$ Red (gray): estimation of the
3 pseudo degrees of freedom, df_{PRESS} and $df_{RSS}=n-4*p$ for PLS models.

4

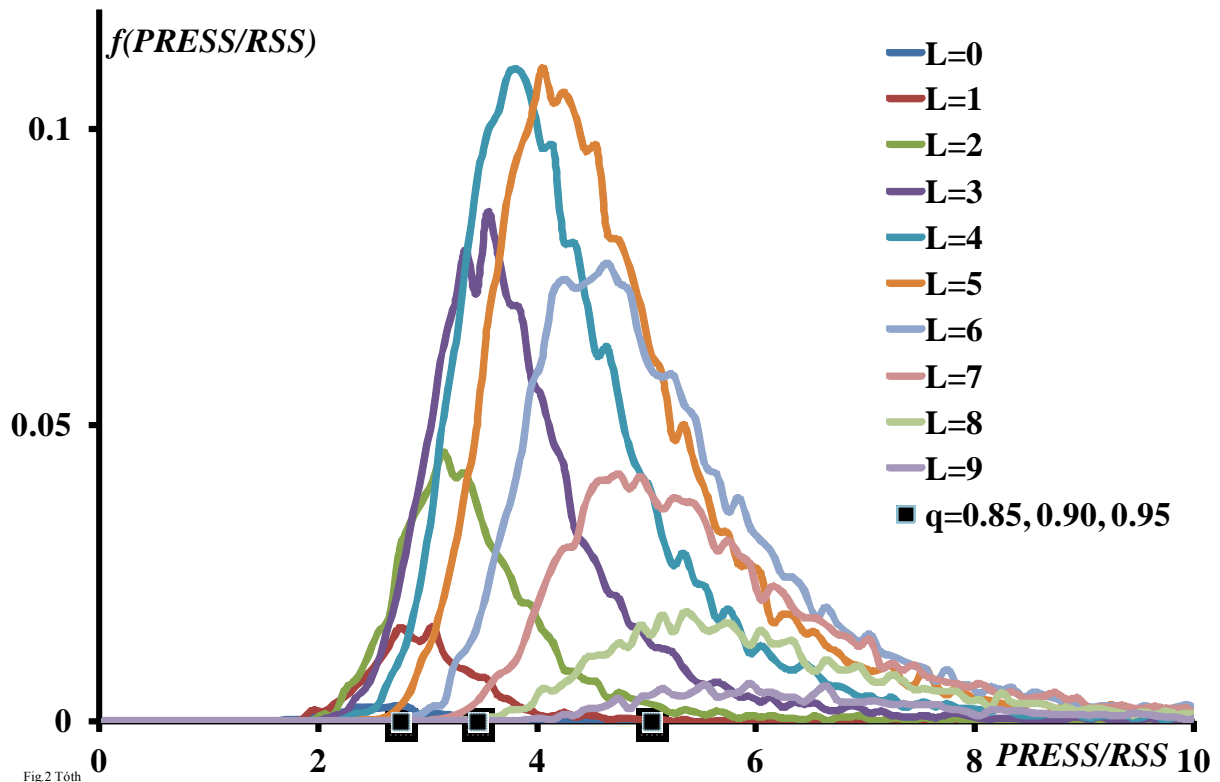
5 Figure 7

6 Scatter plot of bootstrap PRESS/RSS values *versus* standard ones. $n=20$, $m=5$, $w=0.05$

7



1
 2 Figure 1 Relative frequencies of the number of data sets with different number of influential
 3 points (L) identified by the COOK-2 method versus the $PRESS/RSS$ of the data sets. The
 4 black squares denote three percentiles of the corresponding F-distribution. Parameters: $n=10$,
 5 $p=5$ and $w=0.05$
 6



1 Fig.2 Tóth

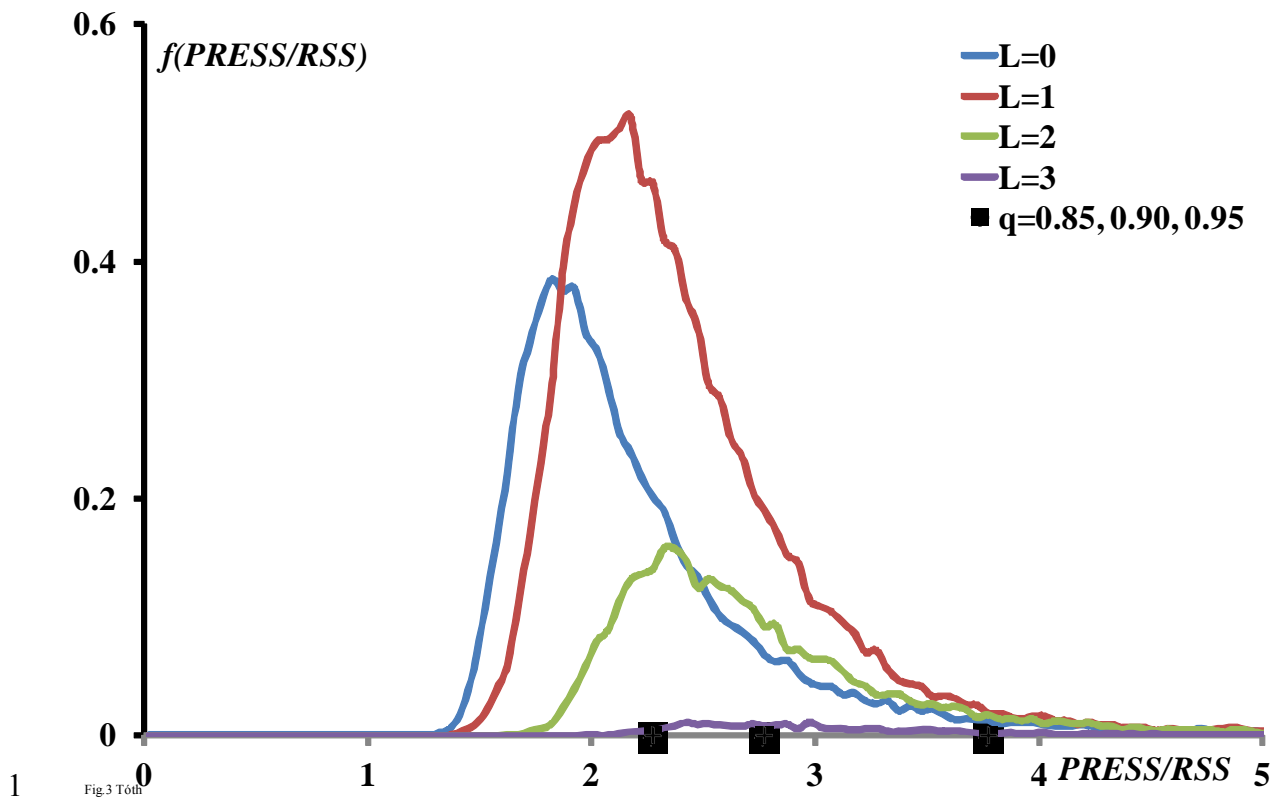
2 Figure 2 Relative frequencies of the number of data sets with different number of influential

3 points (L) identified by the DFBETAS method versus the $PRESS/RSS$ of the data sets. The

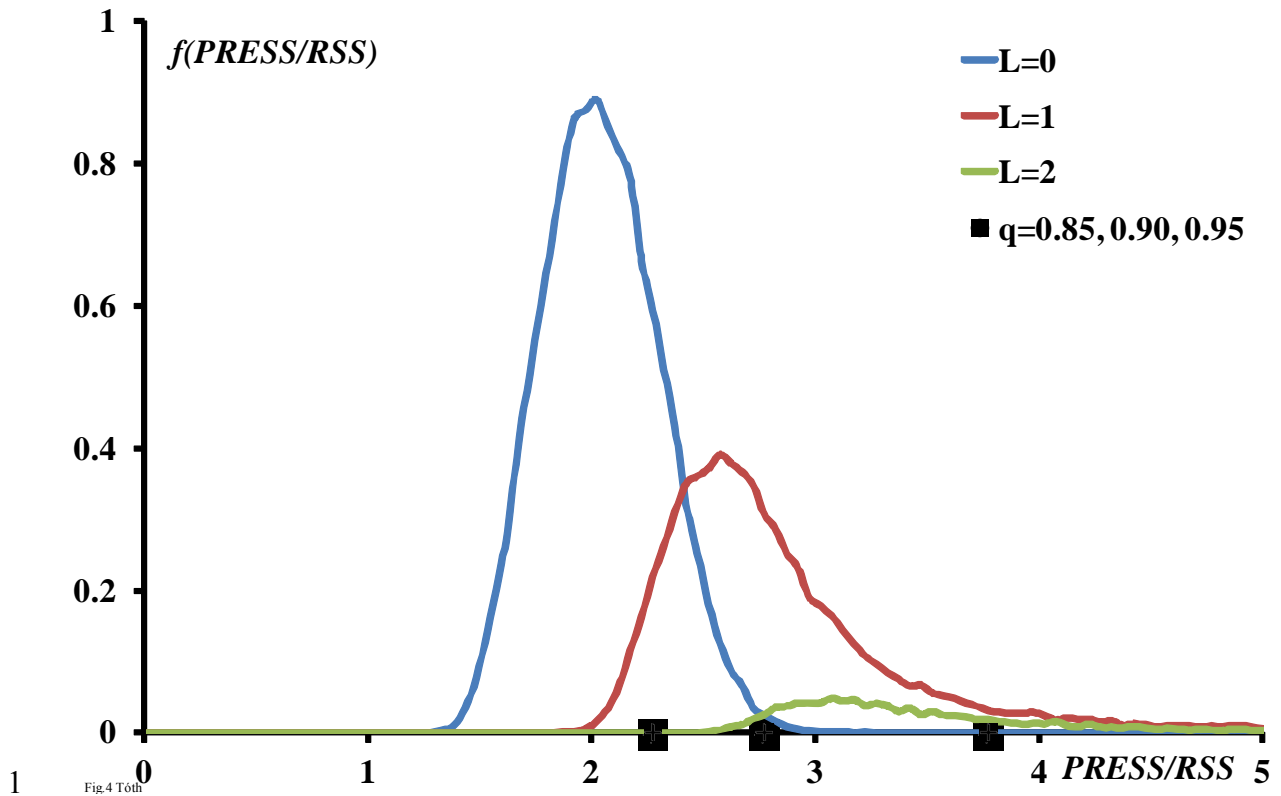
4 black squares denote three percentiles of the corresponding F-distribution. Parameters: $n=10$,

5 $p=5$ and $w=0.05$

6

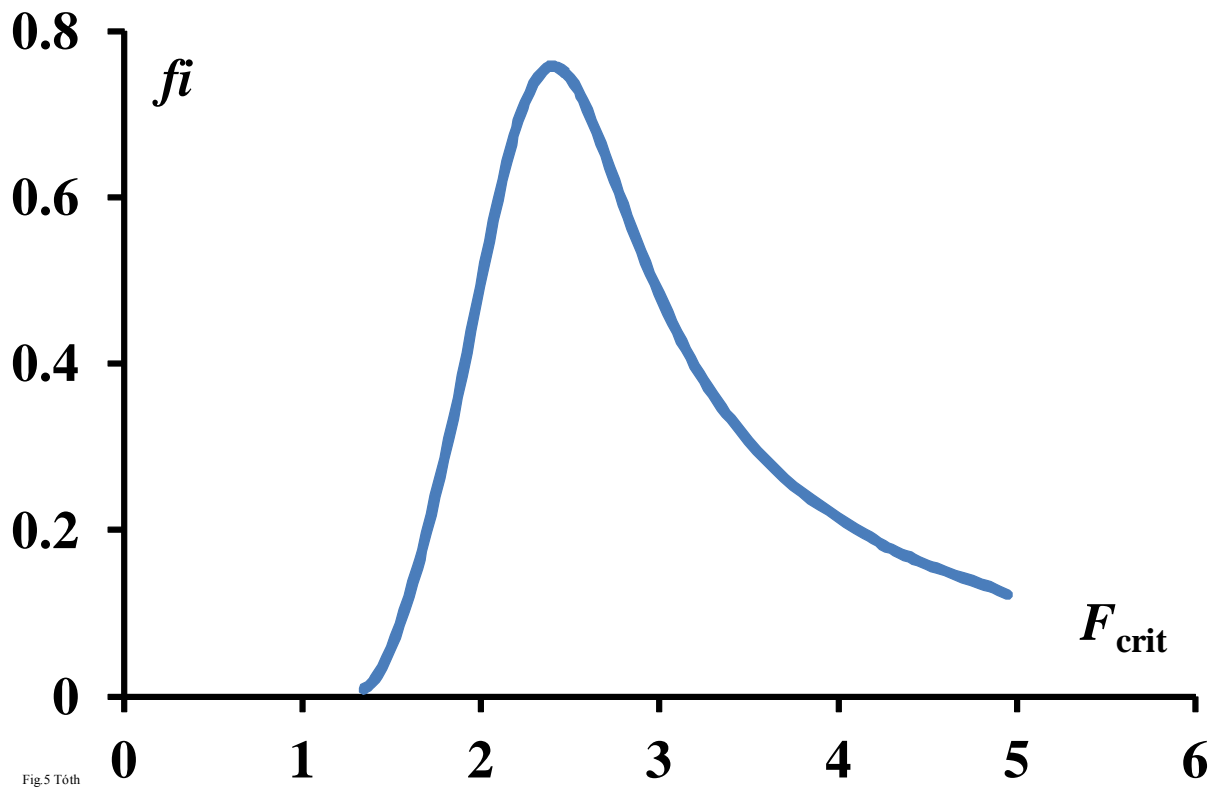


1 Figure 3 Relative frequencies of the number of resampled experimental sets with different
 2 number of influential points (L) identified by the DFFITS-1 method versus the $PRESS/RSS$ of
 3 the data sets. The black squares denote three percentiles of the corresponding F-distribution.
 4 Parameters: $n=10, p=3$

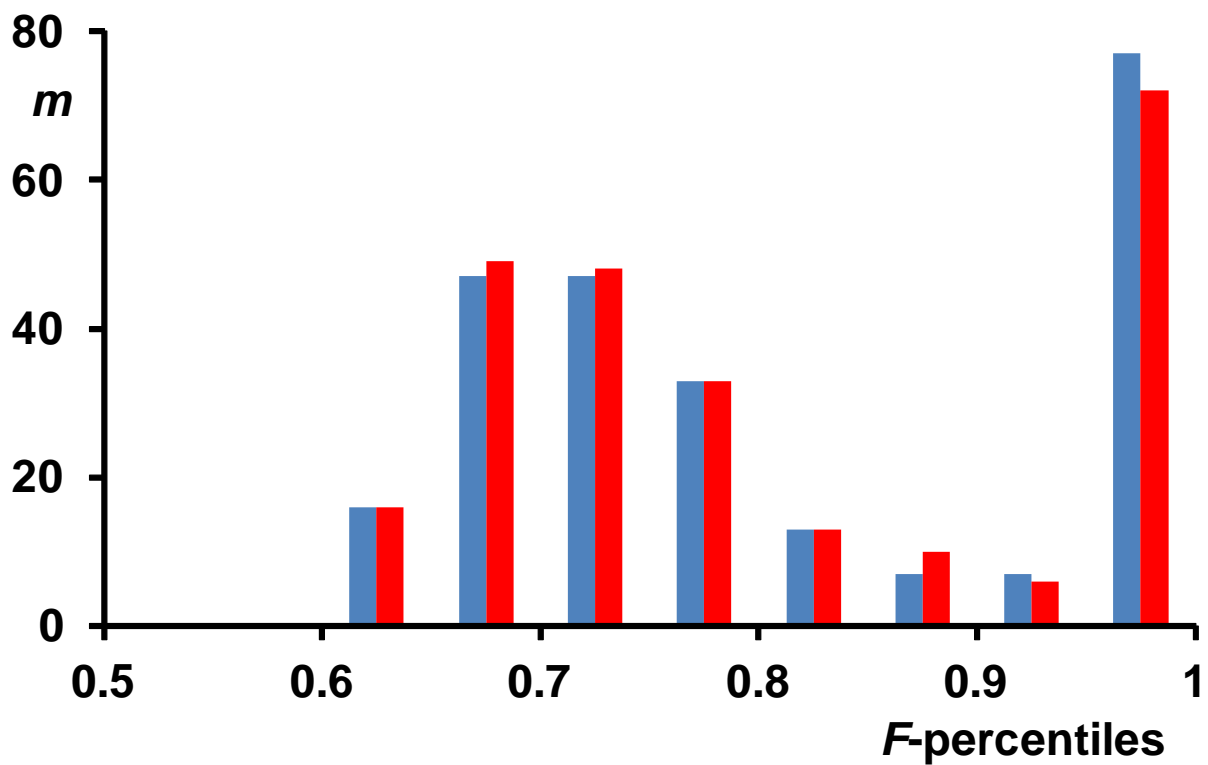


1 Figure 4 Relative frequencies of the number of resampled experimental sets with different
 2 number of influential points (L) identified by the COOK-1 method versus the $PRESS/RSS$ of
 3 the data sets. The black squares denote three percentiles of the corresponding F-distribution.
 4 Parameters: $n=10, p=3$

6



1 Fig.5 Tóth
 2 Figure 5 Dependence of f_i on the choice of the critical F -value for the resampled experimental
 3 data of Zhang et al. $n=10, p=3$.
 4



1 Fig.6 Tóth

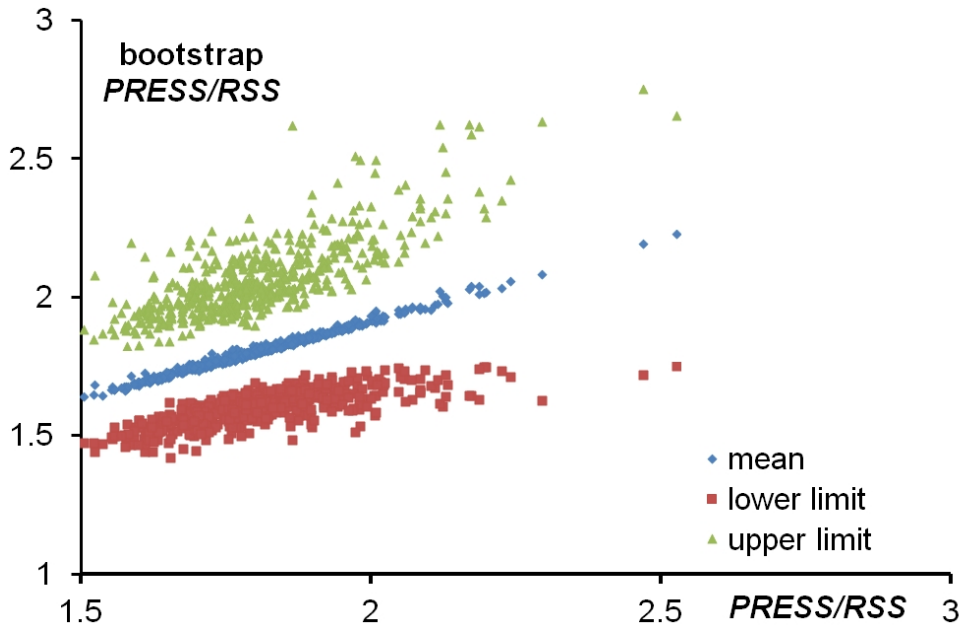
2 Figure 6 Frequencies for the percentiles calculated from F -distributions at the $PRESS/RSS$
 3 values of 247 QSAR models. Blue (black): df_{PRESS} and $df_{PRESS=n-p}$ Red (gray): estimation of
 4 the pseudo degrees of freedom, df_{PRESS} and $df_{RSS=n-4*p}$ for PLS models.

5

1 Figure 7

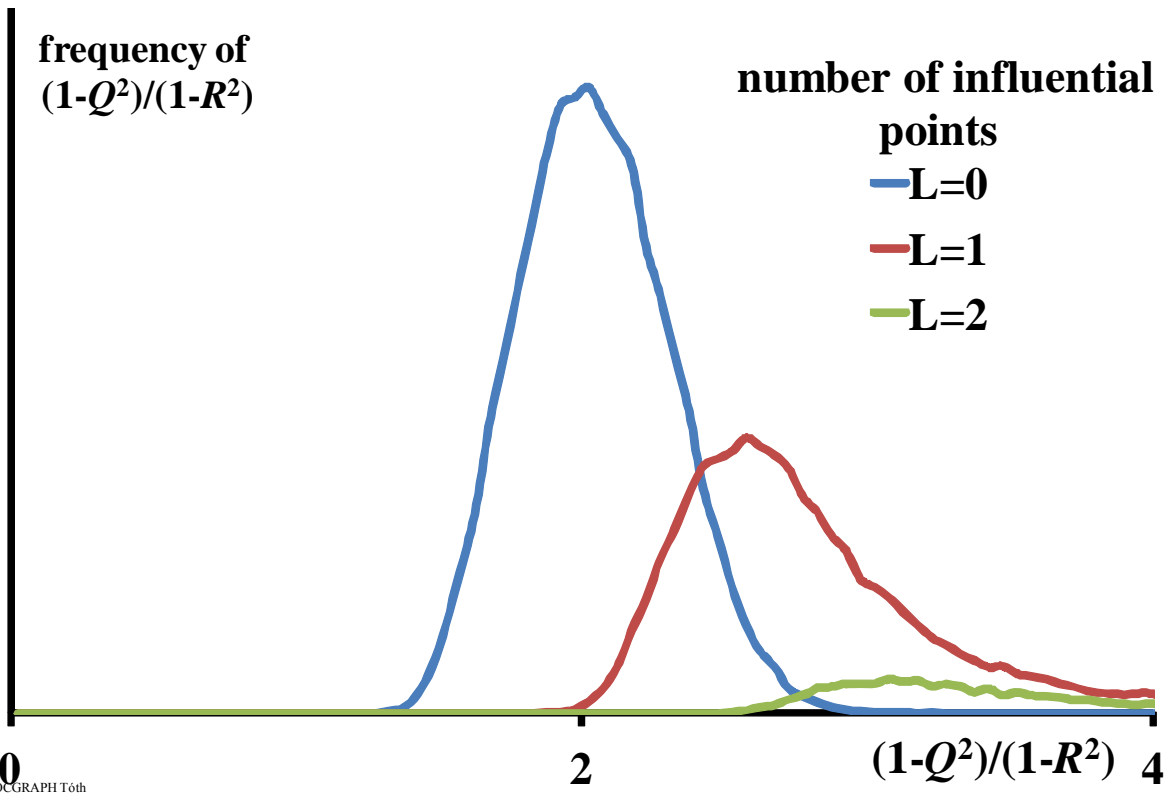
2 Scatter plot of bootstrap PRESS/RSS values *versus* standard ones. $n=20, m=5, w=0.05$

3



4

5



- 1 TOUGRAPH Tóth
- 2 Graphical Abstracts