1 **Quantitative Predictions from Chemical Read-Across and Their**

2 **Confidence Measures**

3

4

5

6

7

8

9

10 **Arkaprava Banerjee, Mainak Chatterjee, Priyanka De and Kunal Roy\***

11

12 *Drug Theoretics and Cheminformatics Laboratory, Department of Pharmaceutical*

13 *Technology, Jadavpur University, Kolkata 700032, India.*

14 *URL: https://sites.google.com/site/kunalroyindia*

15

16

17

18

19

20

21

22 \*Corresponding author; *E-mail: kunal.roy@jadavpuruniversity.in*

23

24

**Abstract**

*In silico* modeling new approach methodologies (NAMs) are viewed as a promising starting point for filling the existing gaps in safety and ecosafety data. Read-across is one of the most widely used alternative tools for hazard assessment, aimed at filling data gaps. However, there are no systematic studies or recommendations on the measures to identify the quality of read-across predictions for the data points without any experimental response data. Recently, we have reported a new similarity-based read-across algorithm for the prediction of toxicity (biological activity in general) of untested compounds from structural analogues (the tool available from https://sites.google.com/jadavpuruniversity.in/dtc-lab-software/home). Three similarity estimation techniques such as, Euclidean distance-based similarity, Gaussian kernel function similarity, and Laplacian kernel function similarity are used in this algorithm. As the confidence of predictions for untested compounds is an important information, we have addressed this issue here by consideration of several similarity and error – based criteria. The role of these measures in discriminating high and low residual query compounds is studied in three different approaches: (a) comparison of means of a measure for high and low residual groups; (b) development of classification models for absolute residuals to identify the contributing measures; (c) application of the sum of ranking differences (SRD) approach to identify the measures closer to the reference rank defined by the absolute residuals. Finally, the frequency of occurrences of different measures in the three approaches is compared. The results from three data sets with 10 divisions of source and target compounds in each case indicate that weighted standard deviation of the predicted response values appear to be the most deterministic feature for the reliability of predictions followed by different similarity-based features. The derived reliability measures will provide a greater confidence to the quality of quantitative predictions from the chemical read-across tool for new query compounds.

**Keywords:** Read-across; Similarity; Prediction; Residual; Discriminant function

**Introduction**

Computational prediction tools are designed and developed as an alternative to experimental biological activity/toxicity tests in order to potentially minimize the need for animal testing, reduce the associated cost and time required for such experimental studies, and improve the quality and availability of data from activity/toxicity prediction and risk/safety assessment [1, 2]. More importantly, *in silico* tools can estimate activity/toxicity of virtual compounds even before their synthesis thus minimizing the cost involved in the synthesis and testing of less potential or less prioritized chemicals. This can help design industrial chemicals/drug candidates with better toxicity/pharmacokinetic profile and prioritize them for experimental testing. Computational methods of toxicity predictions are accepted as tools to bridge data gaps by regulatory agencies like Organization of Economic Cooperation and Development (OECD), European Chemicals Agency (ECHA), Food and Drug Administration (FDA), etc [3-6].

Among various *in silico* techniques for data gap filling, quantitative structure-activity relationship (QSAR) modeling is a popular method [7]. QSAR is a statistical model building process requiring sufficient number of data points for meaningful model development. In addition, in most of the cases, the data points available are required to be split into training and test sets for validation purpose in order to comply with the requirements as recommended by the OECD (https://www.oecd.org/chemicalsafety/risk-assessment/validationofqsarmodels.htm). Thus, a portion of the available experimental data cannot be used for model building and are kept aside for model validation. In case of small data sets, such waste may lead to statistically less reliable model development. Read-across, a chemical similarity-based grouping technique [8], can better address the situation as it does not rely on statistical model development. It is a non-animal alternative data gap filling method that provides information for biological activity/toxicological risks of *target* compounds

77   derived from known activity/toxicity data of *source* compound(s) with a *similar* property or

78   chemical profile. It is one of the most important contemporary *in silico* approaches which is

79   majorly applied in the ecotoxicological data generation, data gap filling, and regulatory

80   decision making. The qualitative read-across approach is most popular and widely used by the

81   regulatory authorities, although the use of quantitative read-across methods has also been seen

82   in the recent past. The query chemicals are mostly termed as the target chemicals whereas the

83   chemical analogues with known toxicity data are called source chemicals. In common practice,

84   read-across predictions are obtained by analogue and category approaches. The analogue

85   approach essentially takes a single source chemical for the prediction, whereas more than one

86   source chemicals are used in the category approach; thus it is more robust and reliable one.

87   Easy algebraic calculations are used in the quantitative read-across algorithm which makes it a

88   computationally less exhaustive process. Apart from that, this method is also an effective

89   approach for the prediction of toxicity of small datasets due to the use of simple calculation

90   (independent of statistical operations). The weighted average of toxicity data (**equation 1**) of

91   chemical analogues is a way for the prediction of untested chemicals.

92
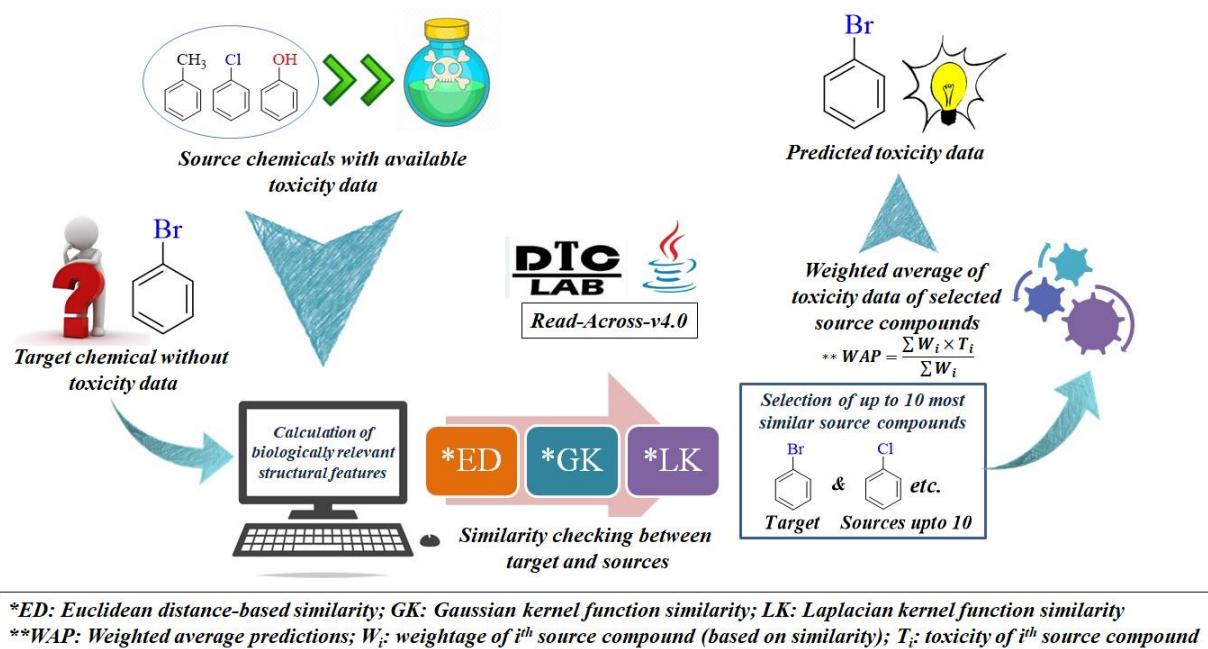$$Weighted\ average = \frac{\sum W_i \times X_i}{\sum W_i} \quad (1)$$

93   where, $W_i$ is the weightage of $i^{th}$ source compounds which is calculated based on the similarity

94   with the target compound; $X_i$ is the toxicity of the corresponding source compound $i$.

95   For a successful read-across operation, the identification of chemical category and the

96   associated uncertainty of this identified category is very important to claim the reliability of

97   predictions. The major objective of the read-across technique is to provide prediction data that

98   is thought to be (more or less) equivalent to the omitted standard experimental assay, and hence

99   this has been applied mainly for toxicity/ecotoxicity data gap filling of chemicals in a

100   regulatory context, However, these new approach methods (NAMs) are finding applications in

101   several other regulatory frameworks, including in the assessment of impurities and degradation

102  products of pharmaceuticals, assessment of plant protection product metabolites, extractables

103  from personal protective and medical devices, food-contact substances, and cosmetics [9].

104  Structural similarity and similar properties, fate and/or activities between the source and target

105  chemicals provide a convenient means of identifying likely analogues and are thus used as a

106  basis for justifying read-across [10]. Apart from only the structural similarity consideration,

107  one should additionally consider physico-chemical properties, reactivity and metabolism, and

108  mechanistic similarity for the precision of predictions [11]. In this direction, the researchers

109  may perform the grouping based on changes in structural aspects and physico-chemical

110  properties and possible fates, degradation and/or the mode of metabolism. Furthermore,

111  identification of the experimental data gaps in physico-chemical characterization, exposure and

112  hazard assessments within the defined groups/categories should also be done [10, 11]. For

113  regulatory acceptance, a read-across prediction should be robust, reliable and easily explicable.

114  Two important aspects of any read-across predictions are the degree of similarity between

115  target(s) and source substance(s) and defining the uncertainties in the read-across predictions

116  [12]. It is generally accepted that the reliability of a read-across prediction depends on the

117  aspects of the defined similarity and the type and degree of uncertainty associated with the

118  particular read-across. Therefore, addressing these two elements in an unambiguous manner is

119  of utmost necessity. Although there are several reports on read-across predictions for different

120  toxicity and ecotoxicity endpoints [13-16], there are no systematic studies and

121  recommendations on the measures of reliability of quantitative read-across predictions for new

122  query compounds. The current manuscript addresses this gap and explores the important

123  measures that may be used to identify the quality of quantitative read-across predictions in

124  absence of experimental data.

125

126  Development of novel algorithms for read-across predictions is a topic of contemporary

127  research in regulatory toxicology. Extensive research has not yet been done for developing

128   algorithms of quantitative read-across predictions. This is especially interesting when limited

129   experimental data for the endpoint of interest is available. Recently, we have developed a read-

130   across tool that predicts the endpoint data of query chemicals based on chemical similarity to

131   the available source compounds using the Euclidean, Gaussian kernel or Laplacian kernel-

132   based similarity functions **(Figure 1)** [17]. We have also applied this tool for prediction of

133   nanotoxicity data of three different data sets showing better quality of predictions than the

134   previously reported read-across and QSAR predictions [17]. However, it is indeed important to

135   know the reliability of read-across predictions for new query compounds without having

136   experimental response values thus not allowing a comparison of predictions with the observed

137   responses. There must be some measures and features that would provide us with confidence or

138   uncertainty of predictions in such cases. The present communication tries to explore the factors

139   governing the reliability of predictions for new query chemicals using the read-across

140   prediction tool.



*ED: Euclidean distance-based similarity; GK: Gaussian kernel function similarity; LK: Laplacian kernel function similarity
**WAP: Weighted average predictions; $W_i$: weightage of $i^{th}$ source compound (based on similarity); $T_i$: toxicity of $i^{th}$ source compound

**Figure 1.** General workflow of Read-across predictions

**Materials and Methods**

The read-across predictions have been done using the Read-Across-v4.0 tool available from

https://sites.google.com/jadavpuruniversity.in/dtc-lab-software/home. In this tool, for each

query compound, up to 10 close source compounds are selected based on the similarity

measure (Euclidean, Gaussian kernel and Laplacian kernel–based similarity), and read-across

predictions are made using a weighted average approach [17]. In the output, weighted average

prediction ($\overline{x_{wtd}}$) values along with weighted standard deviation ($s_{weighted}$) and weighted

standard error ($(s_{\bar{x}})_{weighted}$) are reported [18]. The details are available in Supplementary

Materials SI-1. The user can choose the number of close source compounds to be used by the

tool, optimize the hyper-parameters for Gaussian and Laplacian kernel functions and provide

with the distance and similarity threshold values.

We have used in the current study three different data sets recently used by us for QSAR

modeling: (1) acute contact toxicity of plant protection products against honey bees [19], (2)

Bobwhite Quail ecotoxicity data [20], and androgen receptor binding affinity [21]. We have

used the same physicochemical features as reported in the original QSAR reports for the

present study. For each data set, we have used the original division pattern (training and test

sets) in addition to nine additional new divisions made using a variety of approaches like sorted

response, Kennard-Stone, k-medoids, and random division [22] maintaining the similar

training-test size ratio and ensuring diversity in composition. We have used here the term

"training set" for the whole set of source compounds and the term "test set" for the whole set of

query compounds. For the original division of each data set, optimization of hyper-parameters

and distance and similarity threshold settings were done based on a sub-training set and a

validation set. The optimized settings were used for the original division and nine additional

divisions for read-across predictions. As our objective of this study is not to obtain the best

predictions from a given data set, and it is rather to explore the features indicating the quality

170 of quantitative predictions, we have not done optimization of the settings separately for each

171 division.

172 The read-across tool generates, in addition to read-across predictions, various similarity and

173 error measures such as standard deviation and coefficient of variation of the activity of similar

174 source compounds for each query compound, average and standard deviation of similarity

175 levels and their coefficient of variation of similar close compounds to each query compound,

176 maximum similarity level to positive and negative compounds (based on the whole "training

177 set" response mean), a concordance measure indicating similarity to positive, negative or both

178 classes of close source compounds [23], etc. as detailed in **Table 1**.

179

180 **Table 1.** List of similarity and various error measures generated for each query compound

181 during read-across predictions

| Measure | Description | Comment | Formula |
|---|---|---|---|
| SD_activity (s$_{weighted}$) | Standard deviation of the (observed) response values of the selected close source compounds for each query compound | Dispersion measure | $$s_{weighted} = \sqrt{\frac{\sum_{i=1}^{n} w_i (x_i - \overline{x_{wtd}})^2}{\sum_{i=1}^{n} w_i} \times \frac{n}{n-1}}$$ where, $\overline{x_{wtd}} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$ , $w_i$ is the respective weight for the response $x_i$, $n$ is the number of data points used in computation of the average. |
| CV_activity | Coefficient of variation of the response | Relative Error measure | $$CV_{activity} = \frac{s_{weighted}}{\overline{x_{wtd}}}$$ |
| Average similarity | Mean similarity to the close source | Similarity measure | $$Similarity_{average} = \frac{\sum_{i=1}^{n} Similarity_i}{n}$$ |

8

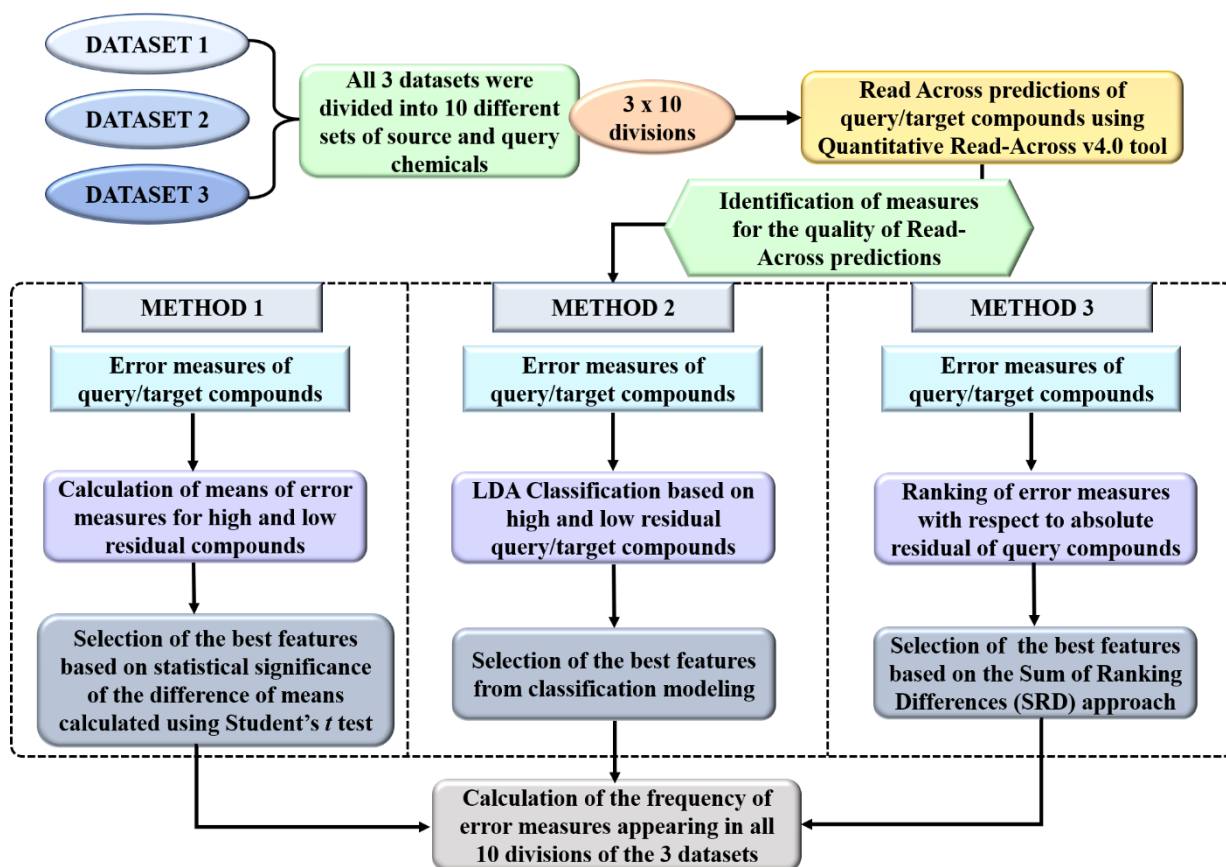| | compounds for each query compound | | |
|---|---|---|---|
| SD_similarity | Standard deviation of the similarity values of the selected close source compounds for each query compound | Dispersion measure | $$S_{similarity=}\sqrt{\frac{\sum_{i=1}^{n}(Similarity_i - \overline{Similarity})^2}{n-1}}$$ where $\overline{Similarity} = Similarity_{average}$ |
| MaxPos | Maximum Similarity level to Positive close source set compounds (based on the "training set" observed mean) | Similarity measure | |
| MaxNeg | Maximum Similarity level to Negative close source set compounds (based on the "training set" observed mean) | Similarity measure | |
| AbsDiff *or* Abs(MaxPos-MaxNeg) | Absolute difference between MaxPos and MaxNeg | Similarity measure | $$AbsDiff = |MaxPos - MaxNeg|$$ |
| *g* | This is a concordance measure | Similarity measure | $$g = 1 - 2 \times |PosFrac - 1/2|$$ where *PosFrac* is the fraction of the close source compounds belonging to the Positive Class based on the "training set" response mean as the threshold [23]. |

182

183    We have attempted to understand the role of the above measures in determining the quality of

184    quantitative read-across predictions for new query compounds. We have done this analysis

185    using three different strategies by studying:

186

187    1. The frequency of measures showing statistically significant differences between the

188    corresponding means of high residual and low residual target set compounds

189

190    2. The frequency of measures appearing important in the developed classification models for

191    high and low residual target compounds

192

193    3. The frequency of measures found important to rank target compounds based on their

194    absolute residuals in the Sum of Ranking Difference approach.

195

196

197    The above three strategies were applied to ten different divisions (source and target

198    compounds) of three different data sets (Figure 2).

**Figure 2.** General workflow of the current study.

1. **Study 1. Comparison of means of high residual and low residual groups:** For each set of division of each data set, we have compiled the read-across predictions of the query set compounds along with various similarity and error measures (as in Table 1) in addition to the observed and predicted response values, then ranked the query compounds in the descending order based on predicted residuals, and finally, identified two sets of 10 compounds with the highest and lowest predicted residual values. We have then compared the means of residuals and different similarity and error measures of the two sets of compounds (high residual and low residual compounds) to identify the important similarity and error measures showing significant difference between the two groups of compounds. The *t* test for comparison of means of two groups [24] was used for this purpose (Supplementary Materials SI-1) with the Gaussian distribution assumption of both the classes.

213

2. **Study 2. Linear discriminant analysis of graded residuals using error and similarity measures:** We used the compiled data of residual values along with different similarity and error measures as described under Study 1 above and graded the data points as positive (1) or negative (0) based on the mean residual value of the corresponding "training set". Then we used the graded response as the dependent variable (Y) and different similarity and error measures and the predicted activity as the independent variables (X) for developing linear discriminant analysis (LDA) models using stepwise variable selection with the F-to-enter 4 and F-to-remove 3.9 setting (in most of the cases) using SPSS statistics software [25]. The LDA tries to maximize the variance between the classes while minimizing the within-class variance, using a linear discriminant function [26]. This also assumes that data in every class are described by a Gaussian probability density function with the same covariance. A linear discriminant function, which is a linear combination of the independent (X) variables, divides the feature space by a hyperplane decision surface. Although we understand that it is overoptimistic to model predicted residuals or errors in this approach and hence, we do not aim at obtaining a perfect classifier, this exercise will definitely help identifying important measures and throwing a light on the reliability of predictions.

3. **Study 3. Application of the sum of ranking differences (SRD) to identify the important measures for ranking the query compounds based on their quality of predictions**

The sum of ranking differences [27] is a useful way to compare metrics, methods, models, methods, analytical techniques, etc. in a general manner. We have used this method to compare the performance of various similarity metrics to understand the quality of read-across predictions. Here, the cases (query compounds) to be ranked are placed in the rows and the metrics in the columns of an input matrix. Then, the results

239   of each metric for each case are ranked in the order of increasing magnitude. The

240   difference between the rank of the metric results and the rank of the known or standard

241   results (here absolute residuals) is then computed. This is followed by the calculation of

242   the sum of absolute values of the differences for all metrics. A lower value of SRD

243   (close to 0) indicates a better metric. The closeness of SRD values indicates the

244   similarity of the metrics, whereas large variation indicates dissimilarity. A permutation

245   test is used for the validation of the SRD method which uses a recursive algorithm for

246   the computation of the discrete distribution for a small number of objects (n<14) or the

247   normal distribution if the number of objects is large. The theoretical distribution is

248   visualized for random numbers and it can be used to identify SRD values for metrics

249   that are far from being random.  A random resampling with sevenfold cross-validations

250   is also applied to validate the obtained results, and the results are presented with a Box-

251   Whiskers plot of the cross-validated SRD values [28].  The SRD runs were made using

252   the program available from http://aki.ttk.hu/srd/.
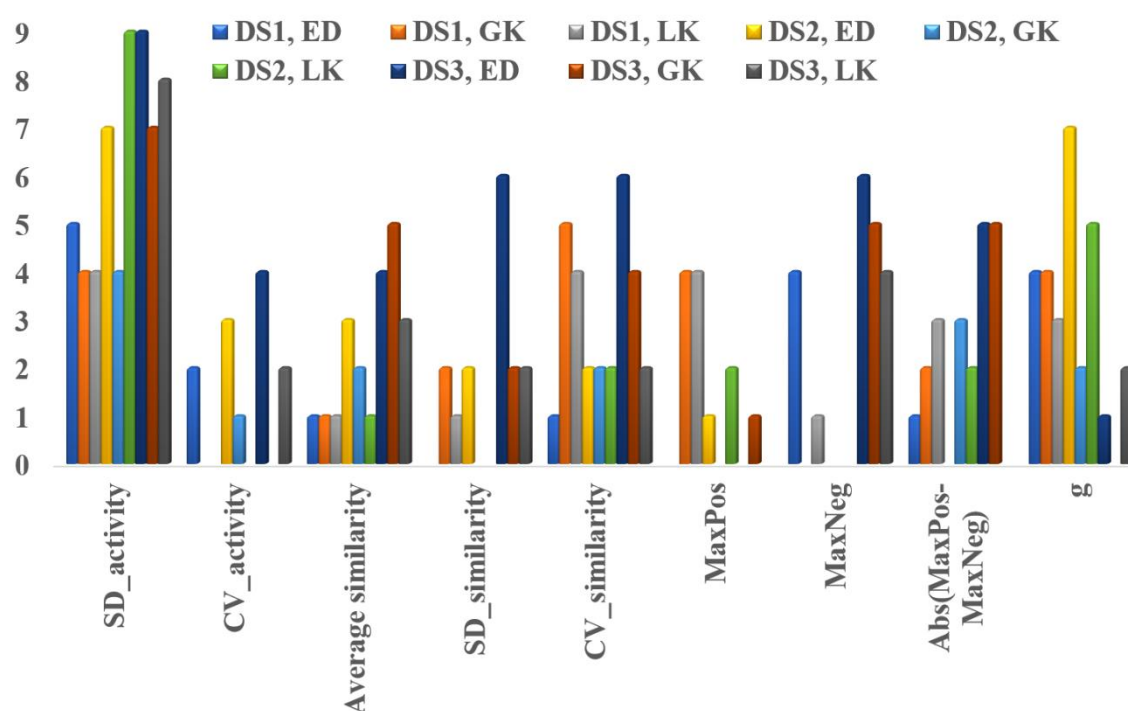
253

254   The results obtained from Studies 1 to 3 above are compared and discussed to conclude

255   on the features responsible for the reliability of quantitative read-across predictions.

256   Please note that the objective of the present analysis is not making new predictions for

257   the data sets being considered or comparing them to the previously reported analysis.

258   We try to explore here various features that may be useful in determining the

259   uncertainty of quantitative predictions from the read-across tool for new query
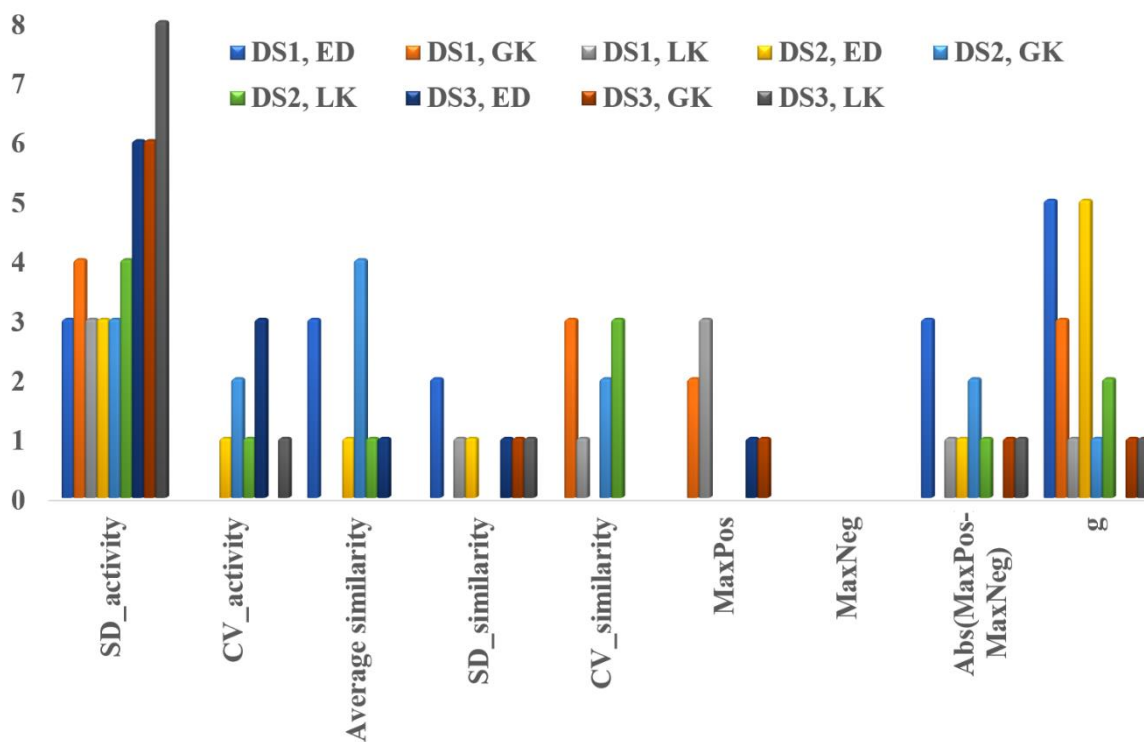
260   compounds.

261

262   **Results and Discussion**

263   We present here the results obtained from the two strategies of our analysis. The frequency of

264   occurrences of different measures for discriminating high and low residual compounds at $p <=$

265   0.05 is shown in **Figure 3** while that for developed LDA models for predicting the class of

266   high or low residuals is shown in **Figure 4** (in addition to **Supplementary Materials SI-1**).

267   The frequency of occurrences of different measures for correctly ranking high and low residual

268   compounds as per the SRD analysis is shown in **Figure 5** (also see **Supplementary Materials**

269   **SI-1).** The details of the results and raw data are available in **Supplementary Materials SI-2**

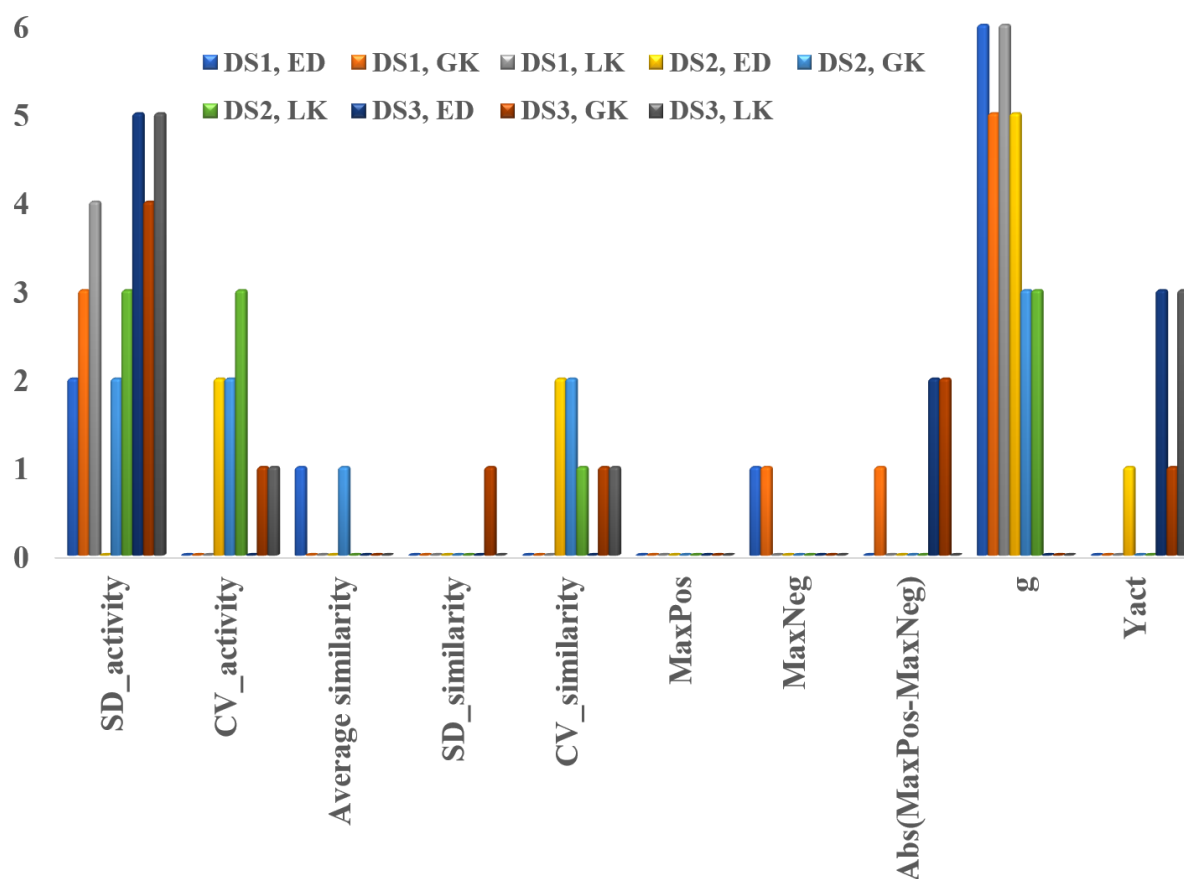270   and on request from the authors.

271



272   **Figure 3.** Frequency of occurrences of different dispersion and similarity measures in

273   differentiating high and low residual compounds (out of 10 trials for each division of each data

274   set) (DS1 = data set 1, DS2 = Data set 2, DS3 = Data set 3, ED = Euclidean distance, GK =

275   Gaussian kernel, LK = Laplacian kernel).

**Figure 4.** Frequency of occurrences of different dispersion and similarity measures in the developed LDA models for predicting the class of high or low residuals (out of 10 trials for each division of each data set) (DS1 = data set 1, DS2 = Data set 2, DS3 = Data set 3, ED = Euclidean distance, GK = Gaussian kernel, LK = Laplacian kernel).
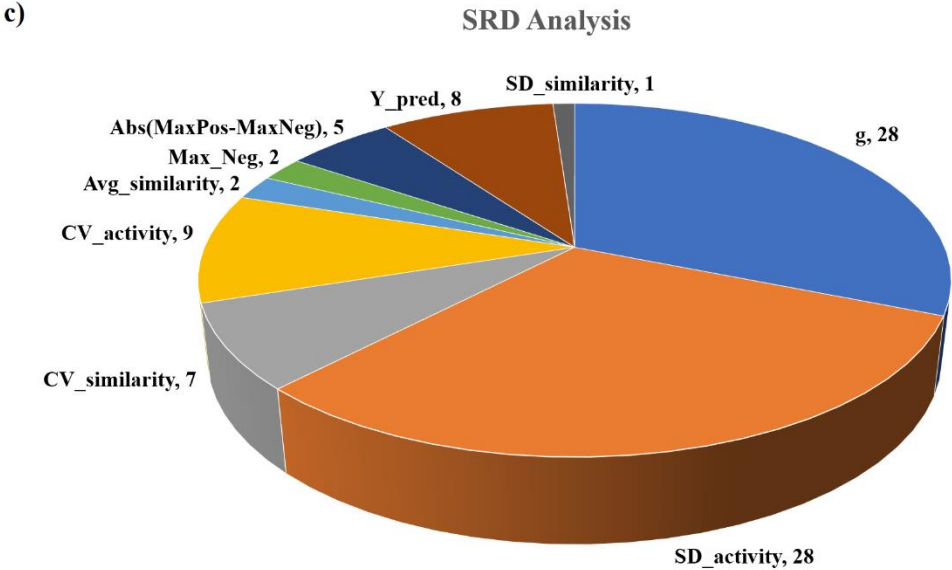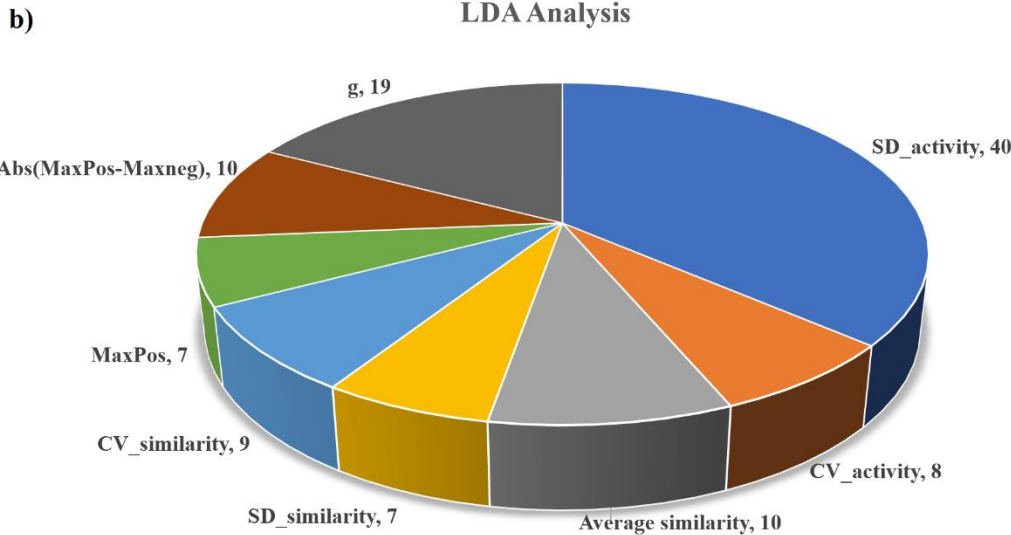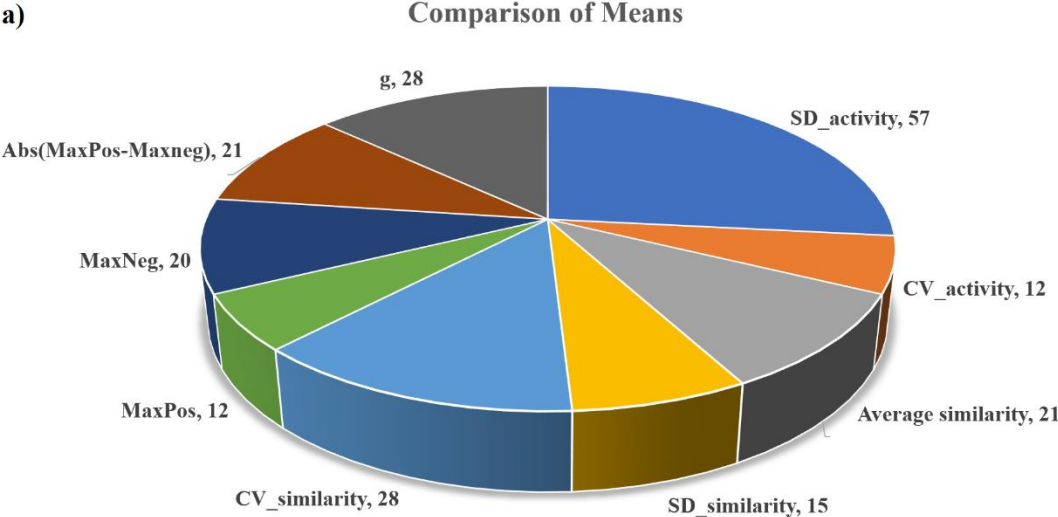
15

**Figure 5.** Frequency of occurrences of different dispersion and similarity measures for correctly ranking high and low residual compounds as per the SRD analysis ((DS1 = data set 1, DS2 = Data set 2, DS3 = Data set 3, ED = Euclidean distance, GK = Gaussian kernel, LK = Laplacian kernel).

**Study 1.**

Figure 1 shows the difference in arithmetic means of different similarity and error measures along with the predicted residuals between the two groups of compounds (high residual and low residual compounds) in the query sets of different divisions of the three data sets as has been found significant at p<0.05 based on the *t* test for comparison of means. It is obvious from the results (**Supplementary Materials SI-2**) that the residual means show statistically significant difference between the two groups in all cases. Among the other measures, we have

295 presented here only those which show statistically significantly difference at $p <= 0.05$ (**Figure**

296 **3).** For Data set 1, in case of the Euclidean distance - based read-across predictions,

297 SD_activity, MaxNeg and $g$ occur 4 times or more (out of 10 trials); in case of the Gaussian

298 kernel-based predictions, SD_activity, CV_similarity, MaxPos and $g$ occur 4 times or more

299 (out of 10 trials); in case of Laplacian kernel - based predictions, SD_activity, CV_similarity

300 and MaxPos occur 4 time or more (out of 10 trials). It is to be noted that SD_activity occurs as

301 the most influential feature considering all similarity-based read-across prediction methods

302 while CV_similarity and MaxPos occur in case of two similarity-based prediction methods. For

303 Data set 2, SD_activity and $g$ emerge as the most significant discriminating features: in case of

304 Euclidean distance - based predictions, both of them occur 7 times out of 10 trials; in case of

305 Gaussian kernel- based predictions, SD_activity occurs 4 times, in case of Laplacian kernel-

306 based predictions, SD_activity occurs 9 times while $g$ occurs 5 times out of 10 trials. For Data

307 set 3, in case of Euclidean distance- based predictions, SD_activity appears nine times

308 followed by predicted response (8 times), observed response, CV_similarity, SD_similarity and

309 MaxNeg (6 times each), absolute difference between MaxPos and MaxNeg (5 times),

310 CV_activity and average similarity (4 times each); in case of Gaussian kernel - based

311 predictions, the most frequently appearing measure is SD_activity (7 times) followed by

312 observed and predicted responses, MaxNeg, average similarity and absolute value of difference

313 between MaxPos and MaxNeg (5 times each) and CV_similarity (4 times); in case of Laplacian

314 kernel- based predictions, the most frequently appearing measures are SD_activity and

315 predicted response (8 times each), observed response (6 times), MaxNeg (4 times).

316 Interestingly, observed and predicted response also show statistically significantly different

317 means in considerable number of trials for Data set 3.

318 A close analysis of the results from three data sets (**Figure 6a**) shows that SD_activity is the

319 most frequently appearing feature for all three data sets. SD_activity corresponds to the

320 dispersion of the observed responses of close source compounds from which a target

17

321    compound is predicted. If this dispersion is higher, the reliability of predictions for the query

322    compound will also be lower as obvious from the high residual values in such cases. The next

323    important measures are $g$ and CV_similarity, each of which occurs for 28 times for the three

324    data sets. While $g$ appears to be important for Datasets 1 and 2, it is not so for Data set 3 where

325    CV_similarity is more important. It appears that either $g$ (concordance measure) or

326    CV_similarity level (along with the average similarity level) is very important in

327    discriminating the high and low residual chemicals. The absolute difference between MaxPos

328    and MaxNeg (especially for Data set 3) is also found somewhat important.

**a)**

**Comparison of Means**



**b)**

**LDA Analysis**



**c)**

**SRD Analysis**



329

330    **Figure 6.** Frequency of occurrences of different important error and similarity measures found

331    from (a) comparison of means between high and low residual compounds; (b) LDA models; (c)
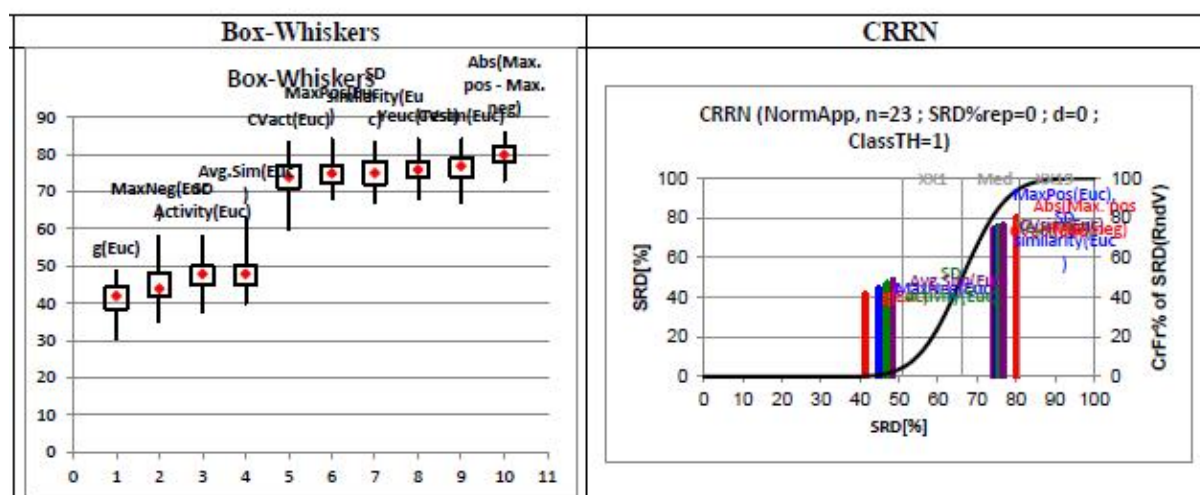
332    SRD analysis

333

334    **Study 2.**

335    The trends obtained from the LDA analysis cannot be expected to be identical with the

336    previous analysis (Study 1) which was based on the quantitative residual values of two sets of

337    samples (high and low residual compounds) drawn from the individual query set while Study 2

338    was performed for the graded residuals of the whole query set. In spite of this, one clear trend

339    we find from the frequency of occurrences of different measures (**Figure 6b**) that SD_activity

340    and *g* occur most in the obtained LDA equations followed by average similarity, CV_activity,

341    CV_similarity and Abs(MaxPos-MaxNeg).

342    **Study 3.**

343    Study 3 applies the sum of ranking differences approach in identifying the most suitable

344    measures that rank the query compounds most similar to the ranking based on the absolute

345    residuals (reference).

346

347  **Figure 7.** Sevenfold cross-validated SRD results (Box-Whiskers plot) and normalized SRD

348  values (between 0 and 100) compared to random ranking (CRRN or comparison of ranks with

349  ranking numbers) for division 1 of Data set 1.

350

351  **Figure 7** shows that ranking of the cases as per $g$(ED) is the closest to that of the reference

352  ranking using absolute residuals for Division 1 of Data set 1. The remaining images of other

353  divisions and other data sets are available in **Supplementary Materials SI-1**. Figure 5 shows

354  the frequency of occurring different measures as determinant of rank order similar to the

355  reference ranking. In line with the observations from Studies 1 and 2, SD_activity and g occur

356  most frequently as the most influential indicator of the quality of predictions followed by

357  CV_activity, $Y_{Pred}$ and CV_similarity (**Figure 6c**).

358

359  Considering the results from all three studies as discussed above **(Figure 6)**, for read-across

360  predictions of new query compounds, the quality of predictions is thought to be dependent on

361  the following factors:

362   1. First line diagnostic measures: Dispersion of activity of close source compounds

363      (SD_activity and CV_activity)

364   2. Second line diagnostic measures: Similarity measures ($g$, CV_similarity and average

365      similarity and Abs(MaxPos-MaxNeg)

366   3. Level of the predicted activity

367  Other measures as discussed above would be related to one or more of the above-mentioned

368  measures.

369  ***Dispersion of response values of close source compounds:*** If the dispersion (in the form of

370  standard deviation and coefficient of variation, but mainly standard deviation) of the response

371  values of the source compounds is high, the precision level of the prediction for a query

372  compound will be low. For example, in Data set 1, Division 3 (Laplacian kernel method), the

21

373   10 compounds showing highest predicted residuals have mean SD_activity value of 0.847 vs.

374   the 10 compounds having lowest predicted residuals showing SD_activity of 0.600. This

375   difference is significant at p <0.05.  This is an inversely proportional relationship suggesting

376   the requirement of selection of optimum number of close source compounds to avoid a high

377   dispersion or error value.

378   **Similarity measures**: "*g*" is a concordance measure to indicate whether the close source

379   compounds selected belong to the same class or the other (based on the mean response of the

380   original list of source compounds as the threshold) [23]. If all of them belong to either positive

381   or negative class, the concordance is higher, and the residual is expected to be low. For

382   example, Division 1 (Euclidean distance method) of Dataset 2 shows that the 10 compounds

383   having high residual values have the mean *g* value of 0.44 compared to 0.18 for the 10

384   compounds having low residual values. This difference is significant at $p < 0.05$. When all (or

385   most of the) close source compounds belong to the same class (either positive or negative), it is

386   expected that the predicted value will be precise and it will not at least misclassify the

387   prediction for the query chemical.

388   The average similarity level to the close source compounds is also an important indicator of

389   reliability. From Division 9 of Data set 2, it is seen that the group of 10 query compounds

390   having high residuals have the average similarity level (Euclidean) of 0.882 compared to 0.946

391   of the group of 10 query compounds having low residual values. If the similarity level

392   increases, reliability of predictions also increases. This is to be noted here that the similarity

393   level of the close source compounds for any query compound is usually higher in case of the

394   Euclidean distance-based approach than the Gaussian kernel based similarity followed by

395   Laplacian kernel based similarity. For this reason, the difference in similarity between the two

396   classes is significant at lower confidence level for the latter two cases. In addition, the

397   interpretation of average similarity is also dependent on data structure (distribution of positive

398  and negative compounds in the data set while the classification is based on the training set

399  response mean as the threshold).

400  The coefficient of variation of the similarity also plays an important role. For example, in

401  Division 2 (Euclidean distance method) of Data set 3, the average CV_similarity level of the

402  set of 10 compounds having high residual values is 0.066 compared to 0.041 for the set of

403  compounds with low residuals. As the CV of similarity values increases, the reliability of

404  predictions decreases. Similar results are also seen for SD_similarity, but its significance is

405  observed in lower number of cases.

406  We may note that the trend mentioned here with regard to the average similarity with respect to

407  high and low residual compounds may also be opposite if the level of dispersion of similarity

408  of close source compounds is high. This happens when the number of close source compounds

409  belonging to either positive or negative class is similar (i.e., *PosFrac* is close to 0.5). For

410  example, in case of the Euclidean distance-based similarity of Division 7 of Data set 1, the

411  dispersion of the similarity values of close similar source compounds is relatively higher and

412  the SD_similarity value of lower residual compounds is actually higher than the high residual

413  compounds, while the average similarity value of lower residual compounds is thus lower than

414  the high residual compounds. This depends on the data structure showing the relative number

415  of close similar source compounds belonging to either positive or negative class and in such

416  cases, SD_activity is the main determining factor for the quality of predictions.

417  The absolute difference between maximum similarity to positive compounds and maximum

418  similarity to negative compounds is also found important in several cases. This difference may

419  be thought to be a perplexity measure. It may be expected that for low residual compounds,

420  this difference may be higher for a more deterministic prediction as observed in case of

421  Laplacian kernel-based similarity of Division 4 of Data set 1. Here, the absolute difference

422  value for the low residual compounds is 0.138 compared to 0.021 in case of high residual

423  compounds. However, the opposite trend is found in Data set 3 where in case of Division 2

424    (Euclidean distance-based similarity), the absolute difference value for low residual

425    compounds is 0.055 compared to 0.123 in case of high residual compounds. In this case, the

426    SD_similarity value is lower for low residual compounds, while in case of Data set 1, Division

427    4 (Laplacian similarity), the SD_similarity value is higher for low residual compounds. This

428    explains the observed difference in the impact of absolute difference value which is in turn

429    dependent on the data structure.

430    Other similarity-based measures like maximum similarity to positive compounds and

431    maximum similarity to negative compounds are also found important in some cases. But their

432    significance depends on the data structure and they are related to other similarity measures

433    already discussed.

434

435    ***Level of predicted values:*** In some cases, especially in case of Data set 3, either or both of

436    observed or/and predicted response values show statistically significant differences between

437    high and low residual compounds. For example, in case of Division 1 (Euclidean distance

438    method) of Data set 3, the average predicted value of the high residual compounds is -1.266

439    while that for the low residual compounds is -2.174. This difference is significant at $p < 0.05$.

440    This indicates that a compound predicted to be lower active has more confidence of predictions

441    than a compound predicted as higher active. The uncertainty level of higher level of

442    quantitative predictions is also higher.

443    Based on the results obtained from the three data sets with their 10 division pattern, we propose

444    here at a preliminary level a set of diagnostic thresholds of different similarity measures (based

445    on Euclidean based similarity) to identify the quality of quantitative predictions (**Table 2**). The

446    first and some of the rest criteria as mentioned in **Table 2** are expected to be met for reliable

447    predictions. Apart from the above, a compound predicted to be more active will have in general

448    less confidence level. However, the indicated thresholds may be more refined in the future with

449    the availability with additional results with other data sets.

24

450

451 **Table 2.** Desired level of different dispersion/similarity measures for good reliability of

452 quantitative read-across predictions (based on Euclidean distance-based similarity)

| Sl. | Dispersion/Similarity measure | Desired range | Reliability |
|---|---|---|---|
| 1. | SD_activity (Euclidean) | <=0.75 | Very good (All criteria met); |
| 2. | $g$ (Euclidean) | <=0.4* | Good (Criterion 1 and at least one |
| 3(a) | Average similarity (Euclidean) | >=0.85 | of the rest but not all); |
| 3(b) | CV_similarity (Euclidean) | <=0.05 | Moderate (Any one met); Bad (None of the criteria met) |

453 *Corresponds to *PosFrac* >= 0.8 or *PosFrac* <= 0.2

454

## Overview and Conclusion

456 In absence of experimental data for toxicity or property of any query chemical, a chemical

457 similarity-based approach is an ideal alternative to bridge the data gaps. Chemical read-across

458 has emerged as a proven method for efficient prediction in this regard which is also recognized

459 and accepted by different regulatory bodies like OECD, US EPA, etc. and regulations like

460 REACH [29]. Although chemical read-across may quickly predict the target property or

461 toxicity of the query chemicals, in absence of the experimental values, it may be challenging to

462 attach a level of uncertainty to the compound-specific predictions. We have discussed this

463 aspect in the context of the Read-Across-v4.0 tool developed by us, but in general the

464 principles should be applicable to other chemical read-across predictions also. From the present

465 analysis, dispersion of the response values of selected close source compounds (specifically

466 standard deviation) emerges to be the most deterministic feature for the reliability of

467 predictions. In the discussed tool, read-across predictions are made using a weighted average

468 approach. Naturally, weighted standard deviation and weighted standard error values are also

469  reported. Based on this, a confidence interval of each predicted value may be presented as

470  below:

471  $95\% \ confidence \ interval \ of \ read-across \ predictions \ = \ weighted \ average \ +$

472  $t_{95\%} \times \frac{s_{weighted}}{\sqrt{n}}$ (2)

473  Apart from the dispersion measures, chemical similarity metrics like concordance measure *g*,

474  which indicates whether the close source compounds belong to either a definite class (positive

475  or negative, leading to more reliability) or a mixed class (less reliability), average similarity

476  level (higher reliability for higher similarity level) and coefficient of variation of similarity (a

477  greater value leads to lower reliability) have been found to important contributing factors. The

478  difference between the maximum similarity levels of query compounds to positive and

479  negative source compounds is also found important in some cases depending on the data

480  structure. The interpretation of the similarity-based measures depends on the data structures. In

481  case of a high dispersion of similarity of close source compounds to a query compound  and/or

482  equal proportion of close positive and negative source compounds for a query compound, the

483  dispersion of observed responses is the main deterministic measure for the reliability of

484  predictions. We have also made a preliminary recommendation about the desired values of

485  different dispersion/similarity measures for good reliability of read-across predictions;

486  however, this may be refined further with the availability of additional results.

487  Finally, a higher range of predicted response values has been found to be associated with

488  higher uncertainty of predictions in some cases. It appears that a compound is predicted to be

489  less active with more certainty than a compound predicted to be higher active.

490  The dispersion and similarity features as listed above may be considered to ascertain the level

491  of confidence during quantitative read-across predictions of query compounds without having

492  experimental response values. These measures will definitely enhance usability of chemical

493  read-across quantitative predictions in absence of observed data. The similarity and error-based

494  measures discussed here are also suitable for a novel kind of modeling (quantitative read-

495  across structure-activity relationship or q-RASAR) which is discussed elsewhere [30].

496

497  **Conflict of interest**

498  Declared none.

499

500  **Acknowledgements**

505

506  **References**

507  [1]    S. Kar, K. Roy, Predictive toxicology using QSAR: A perspective, J. Indian Chem. Soc. 87

508         (2010) 1455–1515.

509  [2]    S. Kar, K. Roy, Risk assessment for ecotoxicity of pharmaceuticals – an emerging issue, Expert

510         Opin. Drug Saf. 11 (2012) 235–274. doi:10.1517/14740338.2012.644272.

511  [3]    A.B. Raies, V.B. Bajic, In silico toxicology: computational methods for the prediction of

512         chemical toxicity, WIREs Comput. Mol. Sci. 6 (2016) 147–172. doi:10.1002/wcms.1240.

513  [4]    S. Kar, H. Sanderson, K. Roy, E. Benfenati, J. Leszczynski, Ecotoxicological assessment of

514         pharmaceuticals and personal care products using predictive toxicology approaches, Green

515         Chem. 22 (2020) 1458–1516. doi:10.1039/C9GC03265G.

516  [5]    S. Klatte, H.C. Schaefer, M. Hempel, Pharmaceuticals in the environment – A short review on

517         options to minimize the exposure of humans, animals and ecosystems, Sustain. Chem. Pharm. 5

518         (2017) 61–66. doi:10.1016/j.scp.2016.07.001.

519    [6]    F. Mansour, M. Al-Hindi, W. Saad, D. Salam, Environmental risk analysis and prioritization of

520           pharmaceuticals in a developing world context, Sci. Total Environ. 557–558 (2016) 31–43.

521           doi:10.1016/j.scitotenv.2016.03.023.

522    [7]    A. Cherkasov, E.N. Muratov, D. Fourches, A. Varnek, I.I. Baskin, M. Cronin, J.

523           Dearden, P. Gramatica, Y.C. Martin, R. Todeschini, V. Consonni, V.E. Kuz'Min, R.

524           Cramer, R. Benigni, C. Yang, J. Rathman, L. Terfloth, J. Gasteiger, A. Richard, A.

525           Tropsha, QSAR modeling: Where have you been? Where are you going to?, J. Med.

526           Chem. 57 (2014) 4977–5010. doi:10.1021/JM4004285

527    [8]    E. Berggren, P. Amcoff, R. Benigni, K. Blackburn, E. Carney, M. Cronin, H. Deluyker, F.

528           Gautier, R.S. Judson, G.E.N. Kass, D. Keller, D. Knight, W. Lilienblum, C. Mahony, I. Rusyn,

529           T. Schultz, M. Schwarz, G. Schüürmann, A. White, J. Burton, A.M. Lostia, S. Munn, A. Worth,

530           Chemical safety assessment using read-across: Assessing the use of novel testing methods to

531           strengthen the evidence base for decision making, Environ. Health Perspect. 123 (2015) 1232–

532           1240. doi:10.1289/ehp.1409342.

533    [9]    S. Kovarich, L. Ceriani, M. Fuart Gatnik, A. Bassan, M. Pavan, Filling data gaps by read-

534           across: A mini review on its application, developments and challenges, Mol. Inform. 38 (2019)

535           1800121. doi:10.1002/minf.201800121.

536    [10]   A. Gajewicz, K. Jagiello, M.T.D. Cronin, J. Leszczynski, T. Puzyn, Addressing a bottle neck

537           for regulation of nanomaterials: Quantitative read-across (Nano-QRA) algorithm for cases

538           when only limited data is available, Environ. Sci. Nano. 4 (2017) 346–358.

539           doi:10.1039/c6en00399k.

540    [11]   A. Gajewicz, Development of valuable predictive read-across models based on "real-life"

541           (sparse) nanotoxicity data, Environ. Sci. Nano. 4 (2017) 1389–1403. doi:10.1039/c7en00102a.

542    [12]   T.W. Schultz, P. Amcoff, E. Berggren, F. Gautier, M. Klaric, D.J. Knight, C. Mahony, M.

543           Schwarz, A. White, M.T.D. Cronin, A strategy for structuring and reporting a read-across

544           prediction  of  toxicity,  Regul.  Toxicol.  Pharmacol.  72  (2015)  586–601.

545           doi:10.1016/j.yrtph.2015.05.016.

546    [13]   G. Schüürmann, R.U. Ebert, R. Kühne, Quantitative read-across for predicting the acute

547    fish toxicity of organic compounds, Environ. Sci. Technol. 45 (2011) 4616–4622.

548    doi:10.1021/ES200361R.

549    [14]    B. van Ravenzwaay, S. Sperber, O. Lemke, E. Fabian, F. Faulhammer, H. Kamp, W.

550    Mellert, V. Strauss, A. Strigun, E. Peter, M. Spitzer, T. Walk, Metabolomics as read-

551    across tool: A case study with phenoxy herbicides, Regul. Toxicol. Pharmacol. 81

552    (2016) 288–304. doi:10.1016/J.YRTPH.2016.09.013.

553    [15]    R. Kühne, R.U. Ebert, P.C. Vonderohe, N. Ulrich, W. Brack, G. Schüürmann, Read-

554    across prediction of the acute toxicity of organic compounds toward the water flea

555    Daphnia magna, Mol. Inform. 32 (2013) 108–120. doi:10.1002/MINF.201200085.

556    [16]    S.J. Enoch, M.T.D. Cronin, T.W. Schultz, J.C. Madden, Quantitative and mechanistic

557    read across for predicting the skin sensitization potential of alkenes acting via Michael

558    addition, Chem. Res. Toxicol. 21 (2008) 513–520. doi:10.1021/TX700322G.

559    [17]    M. Chatterjee, A. Banerjee, P. De, A. Gajewicz-Skretna, K. Roy, A novel quantitative read-

560    across tool designed purposefully to fill the existing gaps in nanosafety data, Environ. Sci. Nano

561    9 (2022) 189–203. doi:10.1039/d1en00725d.

562    [18]    P.R. Bevington, D.K. Robinson, Data reduction and error analysis, McGraw-Hill, New York,

563    1969.

564    [19]    R.K. Mukherjee, V. Kumar, K. Roy, Chemometric modeling of plant protection products

565    (PPPs) for the prediction of acute contact toxicity against honey bees (A. mellifera): A 2D-

566    QSAR approach, J. Hazard. Mater. 423 (2022) 127230. doi:10.1016/j.jhazmat.2021.127230.

567    [20]    R.K. Mukherjee, V. Kumar, K. Roy, Ecotoxicological QSTR and QSTTR modeling for the

568    pediction of acute oral toxicity of pesticides against multiple avian species, Environ. Sci.

569    Technol. 56 (2022) 335–348. doi:10.1021/acs.est.1c05732.

570    [21]    A. Banerjee, P. De, V. Kumar, S. Kar, K. Roy, Quick and efficient quantitative predictions of

571    androgen receptor binding affinity for screening endocrine disruptor chemicals using 2D-QSAR

572    and chemical read-across. ChemRxiv (2022). https://doi.org/10.26434/chemrxiv-2022-gcrjg

573

574 [22]    K. Roy, S. Kar, R. Das, Understanding the Basics of QSAR for Applications in Pharmaceutical

575          Sciences and Risk Assessment, Academic Press, New York, 2015.

576 [23]    J. Wu, S. D'Ambrosi, L. Ammann, J. Stadnicka-Michalak, K. Schirmer, M. Baity-Jesi,

577          Predicting chemical hazard across taxa through machine learning, Environ. Int., 163

578          (2022) 107184. doi:10.1016/j.envint.2022.107184

579 [24]    G. Snedecor, W. Cochran, Statistical Methods, 8th ed., Iowa State University Press, Ames, IA,

580          1989.

581 [25]    SPSS Statistics - India, IBM (2022). https://www.ibm.com/in-en/products/spss-statistics

582          (accessed April 6, 2022).

583 [26]    H. van de Waterbeemd, Discriminant analysis for activity prediction, in: H. van de Waterbeemd

584          (Ed.), Chemom. Methods Mol. Des., VCH, Weinheim, Germany, 1995: pp. 283–293.

585 [27]    A. Rácz, A. Gere, D. Bajusz, K. Héberger, Is soft independent modeling of class

586          analogies a reasonable choice for supervised pattern recognition?, RSC Adv. 8 (2017)

587          10–21. doi:10.1039/C7RA08901E.

588 [28]    K. Héberger, Sum of ranking differences compares methods or models fairly, Trends

589          Anal. Chem. 29 (2010) 101–109. doi:10.1016/J.TRAC.2009.09.009.

590 [29]    H Foth, A.W. Hayes, Background of REACH in EU regulations on evaluation of

591          chemicals, Hum.      Exp.      Toxicol. 27      (2008)      443-461.

592          doi:10.1177%2F0960327108092296

593 [30]    A. Banerjee, K. Roy, First report of q-RASAR modeling towards an approach of easy

594          interpretability and efficient transferability, ChemRxiv Cambridge Open Engag. (2022).

595          doi:10.26434/chemrxiv-2022-0qclt.

596

597

598