

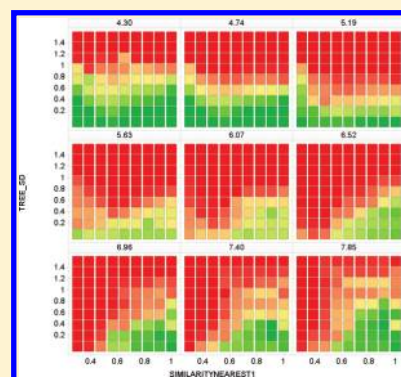
Three Useful Dimensions for Domain Applicability in QSAR Models Using Random Forest

Robert P. Sheridan*

Chemistry Modeling and Informatics, Merck Research Laboratories, Rahway, New Jersey 07065, United States

S Supporting Information

ABSTRACT: One popular metric for estimating the accuracy of prospective quantitative structure–activity relationship (QSAR) predictions is based on the similarity of the compound being predicted to compounds in the training set from which the QSAR model was built. More recent work in the field has indicated that other parameters might be equally or more important than similarity. Here we make use of two additional parameters: the variation of prediction among random forest trees (less variation among trees indicates more accurate prediction) and the prediction itself (certain ranges of activity are intrinsically easier to predict than others). The accuracy of prediction for a QSAR model, as measured by the root-mean-square error, can be estimated by cross-validation on the training set at the time of model-building and stored as a three-dimensional array of bins. This is an obvious extension of the one-dimensional array of bins we previously proposed for similarity to the training set [Sheridan et al. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1912–1928]. We show that using these three parameters simultaneously adds much more discrimination in prediction accuracy than any single parameter. This approach can be applied to any QSAR method that produces an ensemble of models. We also show that the root-mean-square errors produced by cross-validation are predictive of root-mean-square errors of compounds tested after the model was built.



■ INTRODUCTION

In quantitative structure–activity relationship (QSAR) study, a statistical model is generated from a training set of molecules (represented by descriptors) and their activities. The model is then used to predict the activities of molecules not in the training set. A relatively new subfield for QSAR is that of “domain applicability”.^{1–15} The idea behind this subfield is that the prediction accuracy (or “reliability”) of the model cannot be represented by a single number for all possible molecules to be predicted, but that some molecules are more or less likely to be predicted accurately based on the relationship of those molecules to the “domain” of the model. Domain applicability is important because it is now common for users in the pharmaceutical industry to make predictions from QSAR models without having direct knowledge of the molecules that went into making the models. The builders of models must inform the user not only of the prediction for a given molecule but of the reliability of the prediction as well. A paper from our laboratory⁸ was among the first to propose that predictions for molecules close to the training set tend to be more accurate than for compounds that are far away. “Close” means either “has many similar compounds in the training set” or “has a high similarity to the most similar compound in the training set.” One measure of accuracy is the root-mean-square-error (RMSE) between observed and predicted activity (the smaller this value, the more accurate the prediction), and one can show quantitatively how RMSE falls with increasing similarity. This trend does not depend on any specific QSAR method, and the descriptors used to calculate the similarity do not have to be

related to the descriptors used for the QSAR model. Subsequent literature has upheld that observation for many different definitions of similarity and for many different QSAR methods and data sets.^{3,6,7,9,10,12,13} More recent literature has suggested that there are other, perhaps more important, metrics than closeness to the training set.^{2,11}

This paper is the result of investigating parameters in addition to similarity. Here we devise a system of domain applicability that depends on a total of three parameters: similarity to the training set, variation of prediction among random forest trees, and the predicted value itself. We are able to show that using these three parameters together permits much more discrimination than similarity to the training set alone.

■ METHODS

QSAR Method, Descriptors, and Similarity. A wide variety of QSAR methods are described in the literature. One highly regarded method, and the one we will use exclusively here, is random forest.¹⁶ Random forest is an ensemble recursive partitioning method where each recursive partitioning “tree” is generated from a bootstrapped sample of compounds, with random subset of descriptors used at each branching of each tree. Typically we generate 100 trees; adding further trees does not improve prediction accuracy. Two useful features of random forest are that one does not have to do descriptor selection to obtain good results and that predictions appear

Received: January 3, 2012

Published: March 2, 2012

robust to changes in the adjustable parameters. Random forest can do regressions, in which case the input activities and predictions will be floating-point values. It can also do binary classifications, in which case the input activities are categories, e.g. “active vs inactive”, and each prediction represents the probability of being active. The fact that random forest produces an ensemble of models makes it very attractive for the work in this paper.

There are also a large number of possible descriptors in the literature that can be used for QSAR. Generally we prefer substructure descriptors to molecule property descriptors, and we find that in our hands a combination of the Carhart atom pairs (AP)¹⁷ and a donor–acceptor pair (DP)—called “BP” in the work of Kearsley et al.¹⁸ give us the most accurate cross-validated predictions. Both descriptors are of the form:

atom type i – (distance in bonds) – atom type j

For AP, atom type includes the element, number of nonhydrogen neighbors, and number of pi electrons; it is very specific. For DP, atom type is one of seven (cation, anion, neutral donor, neutral acceptor, polar, hydrophobe, and other); it is more general.

In our original paper on domain applicability⁸ we calculated molecular similarity using the AP descriptor and the Dice similarity metric, and we stay with that convention here. This definition of similarity is independent of which descriptors might be important in the QSAR model.

Extrapolation Curve As a One-Dimensional Set of Bins. To understand the subsequent sections it is useful to review our original suggestion⁸ to produce a graph that shows the accuracy of prediction as a function of the similarity of the compound being predicted to the nearest molecule in the training set (henceforth called SIMILARITYNEAREST1). This graph is specific for each QSAR model. Since SIMILARITYNEAREST1 encodes “extrapolation” from the training set (in terms of chemical structure), we call the graph the “extrapolation curve.” Assume we have a large number of compounds not included in the QSAR model, their predictions, and their observed values. The top of Figure 1 shows a scatterplot of the absolute error (i.e., |predicted – observed|) for a large number of compounds vs SIMILARITYNEAREST1. Because SIMILARITYNEAREST1 seldom falls below 0.3, in our implementation we draw a lower limit of SIMILARITYNEAREST1 at 0.3 and include any lower similarity values in the 0.3 bin.

One can divide the range of SIMILARITYNEAREST1 into a number of overlapping “bins”. For example, Figure 1 shows the bin with center “0.5” and with a window of ± 0.1 similarity unit. One calculates the root-mean-square error (RMSE) for compounds in that bin. Do this for several bin centers, and one has generated an extrapolation curve (bottom of Figure 1). RMSE over an entire set of molecules is a standard metric for the “goodness” of QSAR prediction; the extrapolation curve merely displays RMSE for subsets of molecules that differ in SIMILARITYNEAREST1. The lower limit of RMSE in a bin would be zero (as in a perfect reliability). One other point of comparison is the RMSE corresponding to “guessing” (as in poor reliability). We estimate “guessing” as the RMSE over all predictions where the correspondence between observed and predicted values have been randomized.

The set of bins provides a “lookup table.” For the prediction of a new molecule on the QSAR model, one notes its SIMILARITYNEAREST1 to the same training set from which

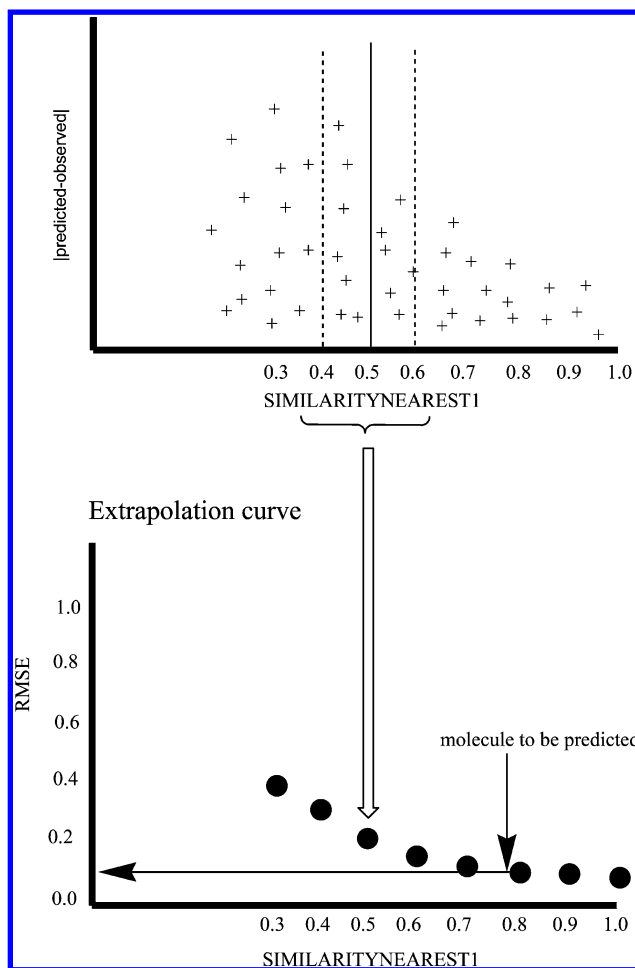


Figure 1. Derivation of the 1-dimensional extrapolation curve from predicted data. In the top picture, each cross represents one prediction. In the bottom picture, each circle represents a bin.

the model was built. In the bottom of Figure 1, the new molecule is closest to the “0.8” bin, and one can take the RMSE for the 0.8 bin as representing the uncertainty, or “error bar” for the prediction, in this example ~ 0.1 . Note that we are not saying |predicted – observed| for that specific molecule will be exactly 0.1, only that the value is a typical absolute error for a molecule in that bin. Obviously one may do sophisticated interpolations between bins, but the table lookup is rapid and accurate enough for most purposes.

In this paper, we use RMSE as our standard method of “prediction goodness.” However, one may imagine alternatives to RMSE. For instance, if one were more interested in correlation than in numerical agreement, one could calculate R^2 between observed and predicted values for each bin (see more in the Discussion section).

Extrapolation Curve by Cross-Validation. In an ideal world the predictions used to build the extrapolation curve would be prospective. One would generate the QSAR model using all available data, wait for enough diverse molecules not already in the model to be tested for the same specific activity, predict those molecules with the model, and process the information into an extrapolation curve as described above. This is difficult in practice, however. We usually want to include all available data in the model, and we need to generate the extrapolation curve at the same time as the model so users can immediately distinguish more vs less reliable predictions.

Also only a few assays have high enough throughput to generate enough prospective measurements in a reasonable time to populate the bins sufficiently. Our original suggestion⁸ to approximate the ideal situation was to generate a large number of predictions at the time of model-building by cross-validation. We can generate as many predictions by cross-validation as necessary to populate the bins sufficiently to achieve reasonable statistics.

Assume we have a QSAR model using method Q and descriptors D generated from a data set T.

- (1) Perform cross-validation:
 - (a) Randomly assign a number n_t of molecules from T to be in the “training set.” Whatever molecules are left are the “test set.” Generate a QSAR model from the training set with method Q and descriptor D.
 - (b) Predict the molecules in the test set with the model. Make a note of the |predicted – observed| for each prediction.
 - (c) Calculate the (AP/Dice) similarity of each molecule in the test set to each molecule in the training set to find the SIMILARITYNEAREST1.
 - (d) Repeat a–c, say, 10 times to build up the number of predictions. The same molecule may be predicted a number of times by models from different cross-validation runs.
- (2) Pool the data (SIMILARITYNEAREST1, |predicted – observed|) for all predictions from all the runs to create the extrapolation curve.

Our original suggestion was to make $n_t = 0.5N$ for all cross-validation runs where N is the number of molecules in the data set. However, our current practice is to vary n_t among the runs, with n_t ranging somewhere between $0.05N$ and $0.5N$ where N is the number of molecules in the data set or 10 000, whichever is smaller. Having a small n_t allows for sampling more molecules with low SIMILARITYNEAREST1. It is a robust observation⁸ that the placement of the curve for RMSE vs SIMILARITYNEAREST1 is not sensitive to n_t . This is what allows us to pool the predictions from multiple values of n_t , even when $n_t \ll N$.

Two More Important Parameters. In searching for more parameters to use in estimating accuracy, we must confine ourselves to those parameters that we can derive for completely novel molecules (i.e., not previously seen in any training set), and for which we do not know the observed activity. We are allowed, however, to use information about how the new molecule relates to the training set for the model. SIMILARITYNEAREST1 is such a parameter. Two new parameters we found to be useful (as will be shown in the Results section) are the following:

- (1) **TREE_SD.** The literature^{2,11} suggests that consistency of prediction from slightly different models (e.g., by different descriptors or by different bagged sets of data) is a very good predictor of the accuracy of prediction. Fortunately, a random forest model is implicitly an ensemble of separate models that make slightly different predictions. One can output the matrix of individual molecules and their predictions on each of the 100 trees in the ensemble. The standard deviation of the predicted activity of a molecule over the trees (TREE_SD), is a measure of the variability of the prediction. Presumably,

the less variation among models, the more accurate the prediction is likely to be.

- (2) **PREDICTED.** Certain ranges of activity may be more easily predicted than others, and the predicted value itself is an important discriminator.

3D Array of Bins. The method of generating predictions by cross-validation is much the same as described above, except that for each prediction we are monitoring a total of three parameters: SIMILARITYNEAREST1, the TREE_SD for that molecule, and PREDICTED, which is the prediction itself. By analogy with the construction of the 1D extrapolation curve, we can calculate RMSE for a 3D array of bins on those three parameters. Currently we are using nine bins in each dimension, covering the minimum and maximum values for those parameters from the data pooled over many cross-validation runs. Each bin has a default window size equal to ± 1 grid spacing in each dimension.

One complication with 3D bins compared to 1D bins is that the number of predictions per bin is much smaller, so we need to increase the total number of predictions to at least 50 000. Smaller data sets may require more runs since the number of predictions per run will be smaller. Bins are populated very unevenly, so even a very large number of total predictions will not populate all bins sufficiently at their default window sizes. Therefore, one practical compromise to control computation time is to widen some of the underpopulated bins at the stage of calculating the RMSE. We incrementally make the window 20% larger in all three dimensions around the bin center until a minimum number of predictions in that bin (e.g., 50) is obtained.

The added dimensionality makes it harder visualize RMSE over a set of 3D bins than 1D bins. In principle, we can display the bins as a grid in three dimensions and display the RMSE as, say, color. However, since much of the grid is obscured any given viewing angle, it is much easier to see trends in RMSE by using a stack of two-dimensional slices. For this paper we are using a “trellis” of nine 2D plots of TREE_SD vs SIMILARITYNEAREST1, each at a specific value (a “slice”) of PREDICTED. In each plot, the RMSE of an individual bin is represented by a color.

As with the 1D bins, the 3D bins constitute a lookup table. For a new molecule being predicted with the model, we calculate the SIMILARITYNEAREST1 relative to the entire data set used for the model, the TREE_SD on the model, and the prediction itself. One identifies the nearest 3D bin to those three parameters, and the error bar of the prediction is the RMSE for that 3D bin.

Data Sets. We will show as examples 12 QSAR data sets listed in Table 1. We tried to include target-specific data sets and ADME data sets of various sizes (~ 1000 to $>100\,000$ molecules) from various sources. Some are from the open literature included for the purposes of “reproduceability,” while others are proprietary data sets from Merck. It is necessary to include some proprietary data in this study. While it is not hard nowadays to find realistically large ($>10\,000$ molecules) data sets from PubChem,¹⁹ these are limited to high-throughput screens and/or confirmation data. Also, some of our later tests require prospective validation, where the activity of molecules is measured after the model (and 3D bins) is built, and dates of testing are easily available in-house but nearly impossible to find in publically available data sets.

Table 1. Data Sets

name	description	source	N	type	RMSE from pooled cross-validation (randomized RMSE)
1A2	−log(IC ₅₀) for CYP 1A2 inhibition	Pubchem AID 1815 refs 20, 21	13 243	regression	0.52 (0.90)
3A4	−log(IC ₅₀) for CYP 3A4 inhibition	in-house	87 312	regression	0.48 (0.76)
AIDS	−log(EC ₅₀) for protection of cells against HIV	ref 22	38 870	regression	0.53 (0.63)
AMES	Ames mutagenicity on Salmonella	ref 23	6512	classification	0.41 (0.56)
CNS	CNS compound or not	ref 24	1685	classification	0.21 (0.50)
FACTORX	−log(IC ₅₀) or −log(K _i) inhibition of human factorX	ChEMBL version 8 TID = 194 ref 25	4784	regression	1.00 (2.16)
HERG	−log(IC ₅₀) for binding to hERG channel	in-house	198 326	regression	0.62 (0.91)
HPLC LOGD	LOGD measured by retention time on HPLC	in-house	101 849	regression	0.74 (1.41)
NK1	−log(IC ₅₀) for binding to substance P receptor	in-house	13 482	regression	0.71 (1.53)
PGP	log(BA/AB) for active transport by <i>p</i> -glycoprotein	in-house	6873	regression	0.37 (0.63)
PXR	induction of pregnane X receptor relative to rifampicin	in-house	67 308	regression	34.9 (48.5)
pyruvate kinase	−log(IC ₅₀) of inhibition of Bacillus pyruvate kinase	Pubchem AID 361 ref 26	51 441	regression	0.65 (0.71)

The distribution of activities for the data sets are in the Supporting Information as Figure S1. In many pharmaceutical data sets, particularly those from high-throughput screening, the distributions are skewed toward low activities. Some of these very low activities may be “qualified data”, for example, “IC₅₀ > 30 μ M.” Most off-the-shelf QSAR methods, including the implementation we use here, do not treat qualified data as anything but a definite value, i.e. >30 μ M would be treated as 30 μ M (or −log(IC₅₀) = −4.5 in units of molar). Generally we find it is necessary to include very inactive molecules in QSAR models so the predictions span the proper range.

RESULTS

Discrimination of Individual Parameters. Figure 2 shows, for some examples, RMSE vs individual 1D bins (green squares). The horizontal line in these plots corresponds to the RMSE for “guessing” (listed in Table 1) for each data set. Table 2 shows the maximum and minimum RMSE among the 1D and 3D bins for all data sets.

The plot of RMSE vs SIMILARITYNEAREST1 in Figure 2 is the same as the “extrapolation curve” from our original publication. It is typical for RMSE to trend downward with increasing SIMILARITY1, as expected from previous work—more similar compounds are better predicted. (One sometimes sees a slight decrease of RMSE at SIMILARITYNEAREST1 = 0.3, which we will discuss later.) It is also typical for RMSE to trend strongly upward with increasing TREE_SD, consistent with the expectation that poorer agreement among random forest trees implies lower accuracy of prediction. The behavior of PREDICTED is more varied. The most common situation is for the maximum RMSE to be at midrange values of PREDICTED, i.e. it is harder to predict midrange than high or low values. The data sets 1A2, 3A4, AIDS, AMES, CNS, FACTORX, and PGP are examples of this behavior. This might be expected a priori for one of the two following reasons:

- (1) Classification data sets such as AMES and CNS have predictions that are in the form of probability. Since the observed values are either 0 or 1, RMSE has to be highest in the region where PREDICTED = 0.5.

- (2) In regression data sets like PGP and NK1, it might be suspected that predicted values at the low and high extremes might have smaller RMSE because there is less range of prediction possible at the extremes.

However, those reasons are not sufficient to explain all cases. We see data sets where very high predictions correspond to a high RMSE; this behavior most often happens where the data set is skewed toward lower observed values. AIDS, HERG, and Pyruvate kinase are examples. The opposite is seen in HPLC LOGD where the lowest predictions correspond to higher RMSE; there the data set is skewed toward higher values.

The relative “discrimination” of parameters can be measured by the range in RMSE. As can be seen in Table 2, TREE_SD is always the most discriminating single parameter, having the largest range in RMSE. TREE_SD is usually much more discriminating than SIMILARITYNEAREST1, consistent with the observations of Dragos et al.² and Sushko et al.¹¹ that variation among models is more important than similarity to the training set. Also, we see the range of RMSE for PREDICTED can be larger than that for SIMILARITYNEAREST1. Table 2 also shows that the range in RMSE in the 3D bins is 9–90% larger (with a mean of 37%) relative to the best single parameter TREE_SD. One can see this for individual examples in Figure 2. Each 1D bin (green) is accompanied by 81 sub-bins (red) based on the other two parameters. It is clear that the red squares cover a much wider range of RMSE than any one set of green squares.

Characteristics of 3D Bins. Figure 3 shows the trellis plots for all data sets. Within each data set, bins are colored such that a deep green represents a RMSE of zero (a perfect prediction for molecules in the bin) and a deep red represents the RMSE corresponding to “guessing.” Pale yellow represents the RMSE for all cross-validated predictions. It is clear for any individual data set that in many places in the 3D space prediction is no better, and sometimes worse, than guessing (red), whereas in other places the predictions can be very good (green); this departs from the overall value (yellow). This again illustrates the idea that accuracy of prediction cannot be summarized by a single number for a given model.

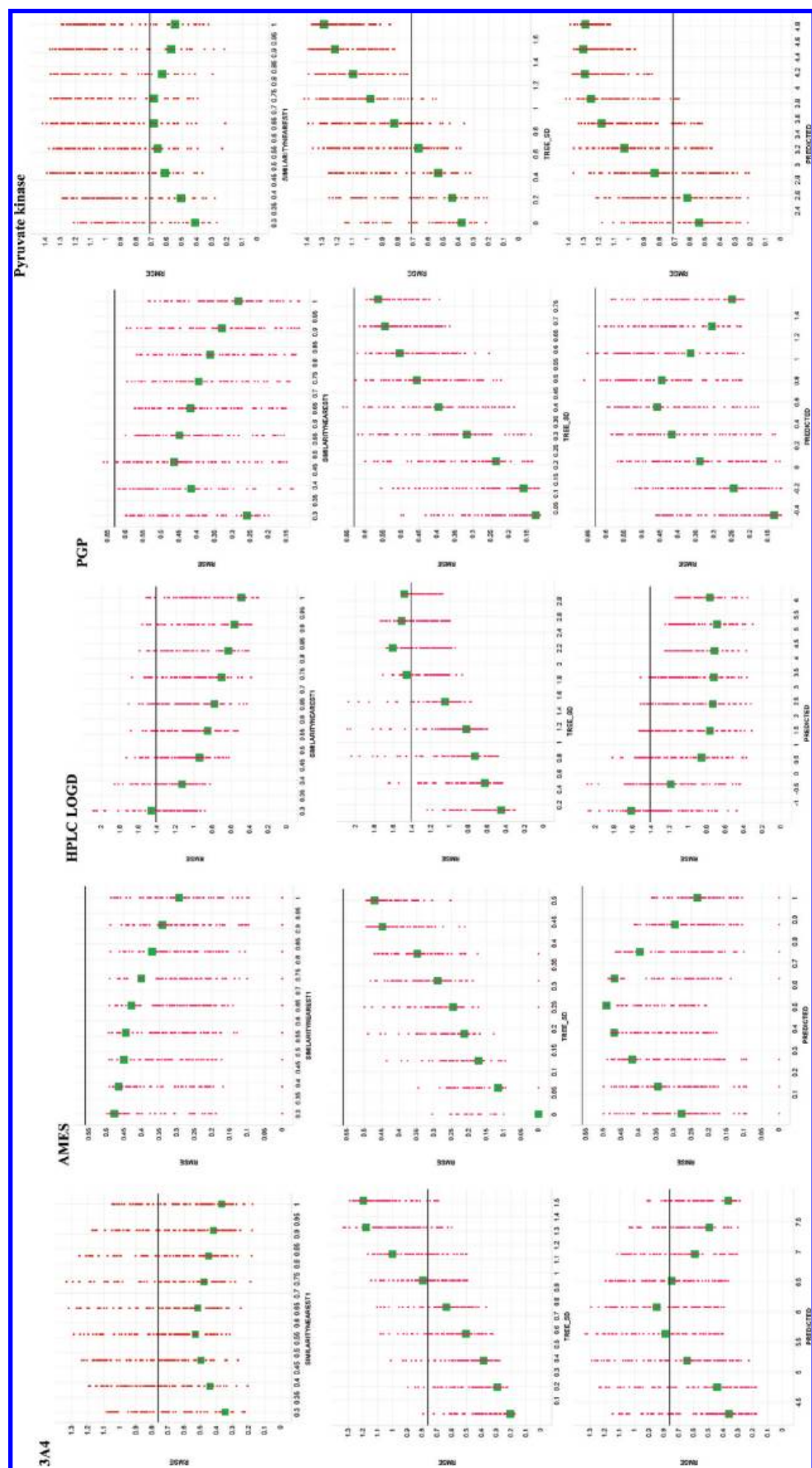


Figure 2. Examples of RMSE in individual one-dimensional bins (green) vs three-dimensional bins (red). The horizontal line represents the value corresponding to guessing, i.e. the RMSE obtained by scrambling the correspondences between observed and predicted.

Table 2. Minimum and Maximum RMSE for Individual Parameter Bins and 3D Bins^a

name	SIMILARITYNEAREST1	TREE_SD	PREDICTED	3DBINS
1A2	0.33–0.57 (0.24)	0.24–1.12 (0.88)	0.28–0.66 (0.38)	0.18–1.14 (0.96)
3A4	0.35–0.53 (0.18)	0.21–1.20 (0.99)	0.36–0.85 (0.49)	0.18–1.33 (1.15)
AIDS	0.51–0.56 (0.05)	0.45–1.75 (1.30)	0.49–1.39 (0.90)	0.38–2.17 (1.79)
AMES	0.29–0.48 (0.41)	0.00–0.47 (0.47)	0.23–0.29 (0.06)	0.00–0.57 (0.57)
CNS	0.07–0.30 (0.23)	0.00–0.45 (0.45)	0.07–0.49 (0.42)	0.00–0.52 (0.52)
FACTORX	0.82–1.52 (0.70)	0.60–1.77 (1.08)	0.69–1.05 (0.36)	0.38–2.43 (2.05)
HERG	0.48–0.69 (0.21)	0.46–1.58 (1.12)	0.47–1.36 (0.89)	0.35–1.61 (1.26)
HPLC LOGD	0.49–1.46 (0.97)	0.45–1.61 (1.16)	0.69–1.62 (0.93)	0.30–2.09 (1.79)
NK1	0.63–1.10 (0.47)	0.41–1.18 (0.77)	0.45–0.82 (0.37)	0.19–1.49 (1.30)
PGP	0.26–0.46 (0.20)	0.13–0.56 (0.43)	0.13–0.46 (0.33)	0.11–0.66 (0.55)
PXR	28.4–39.2 (10.8)	18.1–52.1 (34.0)	18.3–40.3 (22.0)	8.9–59.8 (50.9)
Pyruvate kinase	0.41–0.68 (0.27)	0.38–1.29 (0.91)	0.54–1.31 (0.77)	0.21–1.42 (1.21)

^aRange is in parentheses.

It should be pointed out that the color scale in Figure 3 is fairly conservative in that in practice the lower limit to RMSE would not necessarily be zero, but equal to the experimental error in the activities. For example, a typical error in IC₅₀ would be a factor of 2, or ~ 0.3 on a log scale. Unfortunately, it is hard to know the experimental error for all data sets, especially those from the literature. Setting “deep green” to correspond the experimental error instead of zero would make the plots in Figure 3 look greener overall, but we would not expect the qualitative trends discussed below to change.

Since RMSE generally falls with increasing SIMILARITYNEAREST1 and falls with lower TREE_SD, one might expect to see more green at the lower right of each slice and more red at the upper left. Also, since PREDICTED matters, one would expect some slices to be redder overall than others. CNS, FACTORX, and NK1 are example where all trends work in the expected direction. However, in many other data sets the situation is more complex than a simple additivity of these individual influences. For example, the relative importance of SIMILARITY1 and TREE_SD can vary from slice to slice. For example in 3A4, the slice PREDICTED = 4.30 shows a color gradient from top to bottom with little gradient from left to right, showing TREE_SD is the dominant influence. In the PREDICTED = 7.85 slice, the relative influences of SIMILARITYNEAREST1 and TREE_SD appear more equal. We see a similar phenomenon in 1A2, PGP, and PXR. AMES is an example of the opposite situation, where TREE_SD is more dominant at higher PREDICTED and SIMILARITYNEAREST1 and TREE_SD are more equal at lower PREDICTED.

That fact that SIMILARITYNEAREST1 is not always influential probably accounts for the dip in RMSE we sometimes see at SIMILARITYNEAREST1 = 0.3 for some data sets, a phenomenon we had observed before, but for which we had no explanation in the realm of that single parameter. In PGP, for example, the molecules where SIMILARITYNEAREST1 = 0.3, also have low TREE_SD and low PREDICTED, which in this case would tend toward lower RMSE instead of the higher RMSE we would expect from low SIMILARITYNEAREST1 alone.

Some models are very poorly predictive overall, with most slices being very red. AIDS and Pyruvate kinase are examples. However, there are still some good predictions where TREE_SD and PREDICTED are both low. SIMILARITYNEAREST1 is less influential in those examples. Generally poorer prediction happens in those data sets where the observed activity is biased toward lower values, i.e. there are

many very inactive molecules and few moderate or active molecules, a situation that occurs often in high-throughput data sets. This is reflected in the skewness of the distribution; the poorly predicted data sets tend to have skewness > 2 (see Supporting Information Figure S1). This type of bias causes observed (and presumably predicted) activities to cluster around a particular (low) value, and the chance of predicted and observed matching purely by chance is high. Hence the value of RMSE corresponding to guessing tends to be lower, and the plots look redder overall.

Overall, we can say that the relative importance of individual parameters and the interaction between parameters varies from data set to data set in a way not obvious a priori.

Prospective Prediction. In our original paper from 2004,⁸ we could not validate the idea that RMSE in prospective prediction (i.e., after the model was built) mirrors the RMSE generated by cross-validation, because at the time it was very hard to generate the thousands of prospective predictions needed for useful statistics. However, the throughput of assays is now much greater, and there are five Merck data sets where enough prospective testing has been done that we can process >6000 randomly selected new molecules on existing models: 3A4, HERG, HPLC LOGD, PGP, and PXR. For each new molecule, we know the true $|\text{predicted} - \text{observed}|$. We also know the SIMILARITYNEAREST1, TREE_SD, and PREDICTED, so we can assign that molecule to a bin. We can calculate the RMSE for each bin for the new molecules. On the other hand, each bin has an RMSE that was assigned to it by cross-validation when the model was built. Figure 4 shows the RMSE from the prospective validation vs the RMSE expected from the cross-validation for those five examples. There are four types of bins shown distinguished by color: the 3D bins (red squares) and three 1D bins for the individual parameters. There are more red squares than those of other colors because there are potentially $9 \times 9 \times 9$ 3D bins to be occupied but at most only 9 1D bins. Since the number of prospective predictions is limited, only the bins that contain enough new molecules (here ≥ 50) to make a sensible calculation of RMSE are shown.

There are three things one can say from Figure 4. First, there is a reasonable correlation of prospective and cross-validation RMSE, which is necessary to say the RMSE from cross-validation is a valid metric of “accuracy of prediction in QSAR.” (We re-emphasize that we are not attempting to estimate the $|\text{predicted} - \text{observed}|$ of an individual molecule, only a typical value for a molecule in a specific bin.) Second, the points are not far from the diagonal, meaning the RMSE from cross-validation

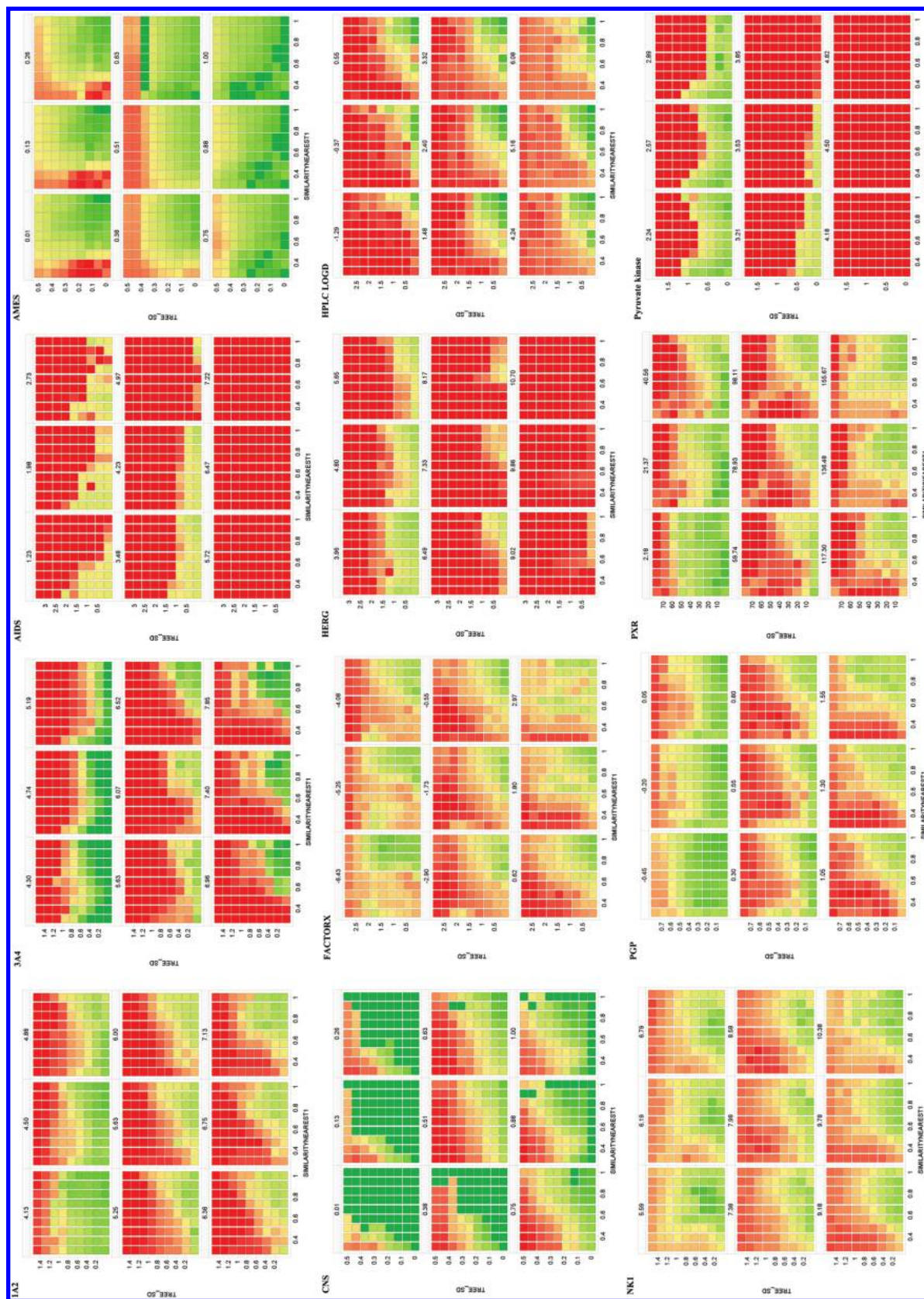


Figure 3. Trellised plots of the three-dimensional space consisting of SIMILARITYNEAREST1, TREE_SD, and PREDICTED for the data sets. Each of nine slices shows TREE_SD vs SIMILARITYNEAREST1 for a given value of PREDICTED. Each small square within each slice represents a bin. Colors are by RMSE with deep green corresponding to zero (extremely accurate predictions) and deep red representing guessing, i.e. the RMSE obtained by scrambling the correspondences between observed and predicted. An intermediate pale yellow represents the RMSE over all cross-validated predictions.

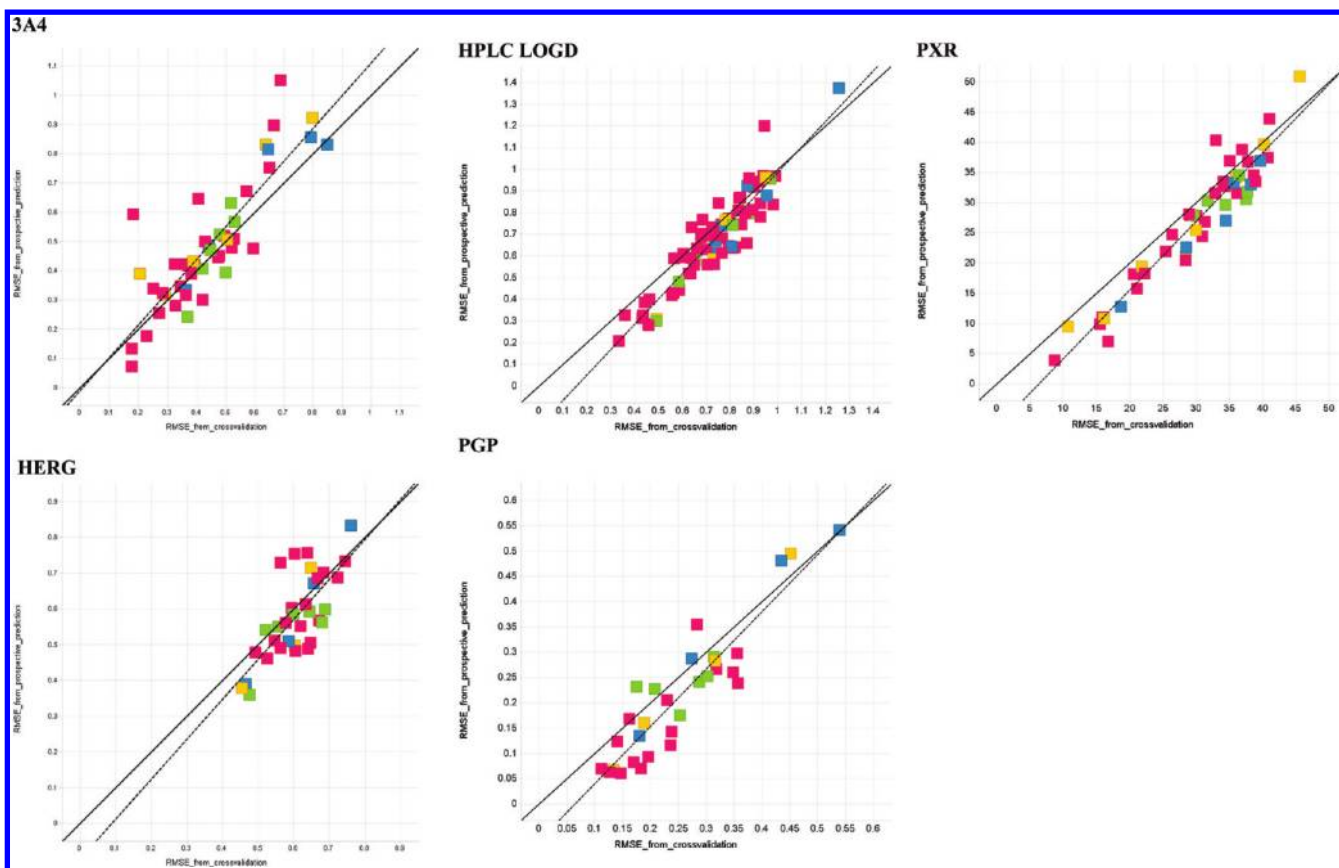


Figure 4. RMSE from prospective prediction vs the RMSE expected from the bins generated by cross-validation. Each square represents one bin. Only those bins that contain ≥ 50 predictions are shown (to ensure reasonable accuracy in RMSE): (color convention) red = 3DBINS, green = SIMILARITYNEAREST1 1D bins, orange = TREE_SD 1D bins, blue = PREDICTION 1D bins. The solid line is the diagonal, and the dashed line is the best linear fit.

gives semiquantitative predictions of prospective RMSE. Finally, all types of bins are coincident, meaning all types of bins are reasonable predictors of accuracy of prediction. However, the usefulness of the 3D bins comes from increased discrimination, i.e. they cover a larger range in RMSE. In particular we note that the SIMILARITYNEAREST1 bins (green squares) have small coverage of RMSE compared to the other types of bins, consistent with Table 2.

DISCUSSION

This paper extends our original method of estimating accuracy of prediction defined by RMSE, which is based on a 1D bin system using SIMILARITYNEAREST1 as the only parameter. The 3D bin system proposed here is based on three parameters: SIMILARITYNEAREST1, TREE_SD, and PREDICTED. The rate-limiting step in estimating RMSE is finding SIMILARITYNEAREST1, especially when the training set for the QSAR model is large. TREE_SD and PREDICTED come along with no additional computational cost, so using three parameters is not more expensive per individual compound than using one parameter.

In this paper we showed that in some data sets TREE_SD and to a lesser extent PREDICTED are more clearly more important as single parameters than SIMILARITYNEAREST1. It was expected from the literature that TREE_SD might be important, but we are not aware of any previous work that suggested that PREDICTED would also be important, although

in retrospect is not unreasonable that some activity values would be harder to predict than others, especially when the data set is skewed to higher or lower values. One can achieve more discrimination between more reliable and less reliable predictions using the three parameters than any single parameter. It is not surprising that adding parameters would allow more discrimination. The point is that each of the parameters studied here can make significant contributions, that the parameters couple in unexpected ways, and that the discrimination is substantially better than in the “similarity only” paradigm. We are not claiming that the system presented here is necessarily optimal. One might find a better definition of similarity to the training set, for example, or identify a parameter even more discriminating than the ones studied here. On the other hand, adding other parameters much beyond three might not be practical, at least using bins. Properly sampling a four- or five-dimensional space may well become too expensive, not to mention much harder to visualize.

While here we used only one combination of substructure descriptor, we have seen similar results with molecular property descriptors or combinations of molecular property descriptors and substructure descriptors.

Finally, in this paper we were able to show that our method of estimating RMSE via cross-validation at model-building time is a reasonable approximation of RMSE from prospective prediction for either 1D or 3D bins, something we could not demonstrate before.

The limitations for using 3D bins compared to 1D bins are the following:

- (1) We need to run more cross-validations, or have fewer bins over the range of the parameters, to generate better statistics per bin.
- (2) It is harder to visualize 3D RMSE information than 1D information and harder to explain the interaction of multiple parameters to chemists.
- (3) We are limited in the types of “goodness” metrics one may apply to the bins. In this paper and our previous paper⁸ we used RMSE, presented as an error bar on the prediction, as a reliability metric. However, since 2004 we have used R^2 between predicted and observed values as an alternative metric in our extrapolation plots using SIMILARITYNEAREST1 (unpublished data). R^2 has the advantage that all data sets can be put into a single “goodness” range 0–1 (1 being perfect prediction and 0 being guessing), whereas the meaning of the value of RMSE depends on the individual data set. Unfortunately, since one of the parameters for 3D bins is PREDICTED, the range of predicted values in any given bin is very narrow and R^2 , which measures correlation, will be very low even when the predictions closely match the observations. In contrast, the metric Q^2 does not depend on having a large range of predicted values in each bin. Q^2 effectively normalizes RMSE by the standard deviation in observed activities thus bringing different data sets into the same range. (In this case, the standard deviation would be for all observed activities, not just those in a particular bin.)
- (4) While SIMILARITYNEAREST1 and PREDICTED can apply to any QSAR method, TREE_SD applies only to methods such as random forest that produce an ensemble of models. Since random forest is accepted as one of the most predictive QSAR methods, and we use it for almost all of our in-house work, this is not a severe limitation. It should be possible to apply an analogous parameter to QSAR methods that produce ensembles of, say, neural networks. For methods that do not produce ensembles, e.g. SVM, one may use SIMILARITYNEAREST1 and PREDICTED to obtain more discrimination than SIMILARITYNEAREST1 alone.

It is possible to question our approach of using a lookup table. One alternative might be to find an equation, or perhaps a recursive partitioning tree, that expresses |predicted – observed| or RMSE as a function of parameters, in a sense making a parametrized model of the reliability of a separate QSAR model. With these alternative approaches one could potentially have a much “smoother” prediction of RMSE, and one might not have to worry so much about sampling certain regions of parameter space, so we might incorporate more parameters, thus getting around some of the limitations we discussed above. The appeal to us of the lookup table is that we can directly view the empirical RMSE at a given place in parameter space, and not the “fit” of an equation to that information. Also one can visualize the influence of the individual parameters and how they couple. Ultimately, the consumer of QSAR predictions and reliability measures are synthetic chemists, and it is usually easier to convey information to them as a picture rather than an equation.

There are two ways to make QSAR more useful: improving QSAR methods so predictions are more accurate overall and improving methods for better discriminating reliable from

unreliable predictions for specific molecules with specific models, i.e. domain applicability. The current state of the field is that we have reached a plateau in the first area and it is hard to improve over random forest and SVM which are considered the current standards.^{16,27,28} Developing new descriptors is also possible, but again, improvement over older descriptors is hard to show. More progress is being made in domain applicability, and several useful parameters have been proposed. From this work, it is clear that using more than one parameter simultaneously is useful to capture the accuracy of QSAR predictions.

■ ASSOCIATED CONTENT

§ Supporting Information

Figure S1: Distribution of activities for data sets. Worked public domain example: original activities, cross-validated parameters, 3Dbins. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: sheridan@merck.com.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

The author thanks Joseph Shpungin for parallelizing random forest so that it can handle very large data sets. The QSAR infrastructure used in this work depends on the MIX modeling infrastructure, and the author is grateful to other members of the MIX team. A large number of Merck biologists, over many years, generated the data for examples used in this paper.

■ REFERENCES

- (1) Baskin, I. I.; Kireeva, N.; Varnek, A. The one-class classification approach to data description and to models applicability domain. *Mol. Inf.* **2010**, *29*, 581–587.
- (2) Dragos, H.; Gilles, M.; Varnek, A. Predicting the predictability: a unified approach to the applicability domain problem of QSAR models. *J. Chem. Inf. Model.* **2009**, *49*, 1762–1776.
- (3) Dimitrov, S.; Dimitrova, G.; Pavlov, T.; Dimitrova, N.; Patlewicz, G.; Niemela, J.; Mekenyan, O. A stepwise approach for defining the applicability domain of SAR and QSAR models. *J. Chem. Inf. Model.* **2005**, *45*, 839–849.
- (4) Ellison, C. M.; Sherhod, R.; Cronin, M. T. D.; Enoch, S. J.; Madden, J. C.; Judson, P. N. Assessment of methods to define the applicability domain of structural alert models. *J. Chem. Inf. Model.* **2011**, *51*, 975–985.
- (5) Gua, R.; Van Drie, J. H. Structure-activity landscape index: quantifying activity cliffs. *J. Chem. Inf. Model.* **2008**, *48*, 646–658.
- (6) He, L.; Jurs, P. C. Assessing the reliability of a QSAR model's predictions. *J. Mol. Graph. Model.* **2005**, *23*, 503–523.
- (7) Kuhne, R.; Ebert, R.-E.; Schuurman, G. Chemical domain of QSAR models from atom-centered fragments. *J. Chem. Inf. Model.* **2009**, *49*, 2660–2669.
- (8) Sheridan, R. P.; Feuston, B. P.; Maiorov, V. N.; Kearsley, S. K. Similarity to molecules in the training set is a good discriminator for prediction accuracy in QSAR. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1912–1928.
- (9) Schroeter, T. B.; Schwaighofer, A.; Mika, S.; Laak, A. T.; Suelzle, D.; Ganzer, U.; Heinrich, N.; Muller, K.-R. Estimating the domain of applicability for machine learning QSAR models: a study on aqueous solubility of drug discovery molecules. *J. Comp. Aid. Mol. Des.* **2007**, *21*, 651–664.

- (10) Sprous, D. G. Fingerprint-based clustering applied to define a QSAR model use radius. *J. Mol. Graph. Model.* **2008**, *27*, 225–232.
- (11) Sushko, I.; Novotarskyi, S.; Körner, R.; Pandey, A. K.; Cherkasov, A.; Li, J.; Gramatica, P.; Hansen, K.; Schroeter, T.; Müller, K.-R.; Xi, L.; Liu, H.; Yao, X.; Oberg, T.; Hormozdiari, F.; Dao, P.; Sahinalp, C.; Todeschini, R.; Polishchuk, P.; Artemenko, A.; Kuz'min, V.; Martin, T. M.; Young, D. M.; Fourches, D.; Muratov, E.; Tropsha, A.; Baskin, I.; Horvath, D.; Marcou, G.; Muller, C.; Varnek, A.; Prokopenko, V. V.; Tetko, I. V. Applicability domain for classification problems: benchmarking of distance to models for Ames mutagenicity set. *J. Chem. Inf. Model.* **2010**, *50*, 2094–2111.
- (12) Tetko, I. V.; Sushko, I.; Pandey, A. K.; Zhu, H.; Tropsha, A.; Papa, E.; Oberg, T.; Todeschini, R.; Fourches, D.; Varnek, A. Critical assessment of QSAR models of environmental toxicity against *Tetrahymena pyriformis*: focusing on applicability domain and overfitting by variable selection. *J. Chem. Inf. Model.* **2008**, *48*, 1733–1746.
- (13) Weaver, S.; Gleeson, M. P. The importance of the domain of applicability in QSAR modeling. *J. Mol. Graph. Model.* **2008**, *26*, 1315–1326.
- (14) Soto, A. J.; Vazquez, G. E.; Strickert, M.; Ponzoni, I. Target-driven subspace mapping methods and their applicability domain estimation. *Mol. Inf.* **2011**, *30*, 779–789.
- (15) Tetko, I. V.; Bruneau, P.; Mewes, H.-W.; Rohrer, D. C.; Poda, G. I. Can we estimate the accuracy of ADME-Tox predictions? *Drug Discovery Today*. **2006**, *11*, 700–707.
- (16) Svetnik, V.; Liaw, A.; Tong, C.; Culberson, J. C.; Sheridan, R. P.; Feuston, B. P. Random forest: a classification and regression tool for compound classification and QSAR modeling. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1947–1958.
- (17) Carhart, R. E.; Smith, D. H.; Ventkataraghavan, R. Atom pairs as molecular features in structure-activity studies: definition and application. *J. Chem. Inf. Comput. Sci.* **1985**, *25*, 64–73.
- (18) Kearsley, S. K.; Sallamack, S.; Fluder, E. M.; Andose, J. D.; Mosley, R. T.; Sheridan, R. P. Chemical similarity using physiochemical property descriptors. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 118–27.
- (19) PubChem, <http://pubchem.ncbi.nlm.nih.gov/> (accessed Oct 1, 2011).
- (20) Sun, H.; Veith, H.; Xia, M.; Austin, C. P.; Huang, R. Predictive models for cytochrome P450 isozymes based on quantitative high throughput screening data. *J. Chem. Inf. Model.* **2011**, *51*, 2474–2481.
- (21) National Center for Biotechnology Information. PubChem BioAssay Database; AID=1815, Source=Scripps Research Institute Molecular Screening Center, <http://pubchem.ncbi.nlm.nih.gov/assay/assay.cgi?aid=1815> (accessed Oct 1, 2011).
- (22) http://dtp.nci.nih.gov/docs/aids/aids_data.html (accessed Oct 1, 2011).
- (23) Hansen, K.; Mika, S.; Schroeter, T.; Sutter, A.; ter Laak, A.; Steger-Hartmann, T.; Heinrich, N.; Muller, K.-R. Benchmark dataset for in silico prediction of Ames mutagenicity. *J. Chem. Inf. Model.* **2009**, *49*, 2077–2081.
- (24) Adenot, M.; Lahana, R. Blood-brain barrier permeation models: discrimination between CNS and non-CNS drugs including P-glycoprotein substrates. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 239–248.
- (25) ChEMBL database. <https://www.ebi.ac.uk/chembl/> (accessed Feb 14, 2012).
- (26) National Center for Biotechnology Information. PubChem BioAssay Database; AID=361, Source=Scripps Research Institute Molecular Screening Center, <http://pubchem.ncbi.nlm.nih.gov/assay/assay.cgi?aid=361> (accessed Oct 1, 2011).
- (27) Bruce, C. L.; Melville, J. L.; Pickett, S. D.; Hirst, J. D. Contemporary QSAR classifiers compared. *J. Chem. Inf. Model.* **2007**, *47*, 219–227.
- (28) Svetnick, V.; Wang, T.; Tong, C.; Liaw, A.; Sheridan, R. P.; Song, Q. Boosting: an ensemble learning tool for compound classification and QSAR modeling. *J. Chem. Inf. Comput. Sci.* **2005**, *45*, 786–799.