

# OptiSim: An Extended Dissimilarity Selection Method for Finding Diverse Representative Subsets<sup>‡</sup>

Robert D. Clark<sup>†</sup>

Tripos, Inc., 1699 South Hanley Road, St. Louis, Missouri 63144

Received April 7, 1997<sup>®</sup>

Compound selection methods currently available to chemists are based on maximum or minimum dissimilarity selection or on hierarchical clustering. Optimizable *K*-Dissimilarity Selection (OptiSim) is a novel and efficient stochastic selection algorithm which includes maximum and minimum dissimilarity-based selection as special cases. By adjusting the subsample size parameter *K*, it is possible to adjust the balance between representativeness and diversity in the compounds selected. The OptiSim algorithm is described, along with some analytical tools for comparing it to other selection methods. Such comparisons indicate that OptiSim can mimic the representativeness of selections based on hierarchical clustering and, at least in some cases, improve upon them.

## INTRODUCTION

The advent of combinatorial chemistry and high throughput screening has made the ability to identify “good” subsets in large libraries of compounds very important, whether the libraries in question are realized or virtual. Traditionally such subsets have been created using expert systems—i.e., having a medicinal or pesticide chemist select compounds manually based on a series 2D structures. This approach is labor-intensive and can be rather dependent on the expert used. Moreover, it is neither routinely practical nor robust for more than 300–1000 compounds, and then only when the library in question includes one or more homologous series.

Currently available alternative methods include maximum dissimilarity selection, minimum dissimilarity selection, and hierarchical clustering, among others.<sup>1</sup> Each of these methods can be effective, but each has some intrinsic limitations. They are described in some detail below, and *Optimizable K-Dissimilarity Selection* (OptiSim<sup>2</sup>) is introduced as a generalized extension of the dissimilarity-based methods.

**Maximum Dissimilarity.** The computational algorithms currently most often used for selecting compounds operate by maximizing the diversity of the selected subset with respect to a set of descriptors and some associated (dis-)similarity measure.<sup>1,3,4</sup> The basic algorithm is straightforward and takes as parameters a minimum acceptable dissimilarity (redundancy) threshold *R* and a maximum selected subset size *M*<sub>max</sub>:

1. Select a compound at random from the dataset of interest and create a list of candidate compounds from the remainder of the dataset.

2. Examine the pool of candidates and identify the candidate which is most dissimilar to those which have already been selected.

3. Is the dissimilarity of that best candidate less than *R* (redundancy test)? If so, quit; otherwise, add it to the selection set and remove it from the pool of candidates.

4. If this is the third selection, return the first two selections to the pool of candidate compounds. (The first selection was chosen randomly, and the second selection is strongly biased by the first. Transferring them back into the candidate pool reduces the effect of the initial random selection.)

5. Are *M* selections in hand? If so, quit.

6. Is the pool of candidates empty? If so, quit; otherwise, go to step 2.

A related method applies a genetic algorithm across the entire set to maximize the diversity of the selected subset.<sup>5</sup>

Maximally diverse subsets are, by definition, biased toward inclusion of outliers. In some situations, this is a very useful property, but medicinal chemists tend to avoid outliers in making their own selections because they may not “look like” drugs. In some cases, outliers in corporate databases are outliers for good reason—difficulty of synthesis or toxicity, for example—which reduces their value as potential leads. Moreover, a maximally diverse subset may not be adequately representative of the biochemical diversity in a dataset.

One justification for maximizing diversity is based on experimental design considerations commonly employed for analyzing quantitative structure/activity relationships (QSARs),<sup>6</sup> where outliers are important because they have the greatest statistical leverage. That leverage is critically dependent on how well biochemical response can be approximated as a linear function of the descriptors being used; however, internal data points are essential for characterizing quadratic response functions and become more important still in more complex systems. The libraries of interest here are usually much broader in scope than are those used for QSAR studies. Response functions are correspondingly more complex, so outliers lose most of that statistical leverage.

**Minimum Dissimilarity.** A complementary approach to compound selection, employed in “leader”-based partitioning methods,<sup>1</sup> can be characterized as minimum dissimilarity selection. The algorithm takes the same two parameters as maximum dissimilarity selection—a minimum dissimilarity threshold *R* and *M*<sub>max</sub>, the maximum number of compounds to select—but applies them differently:

1. Select a compound at random from the dataset of interest and create a list of candidate compounds from the

<sup>†</sup> Phone: (314)-647-1099. Fax: (314)-647-9241. E-mail: bclark@tripos.com. Internet: <http://www.tripos.com>.

<sup>‡</sup> Keywords: molecular diversity, dissimilarity selection, clustering, combinatorial chemistry, compound selection, diversity analysis, chi square, and representative.

<sup>®</sup> Abstract published in *Advance ACS Abstracts*, September 15, 1997.

remainder of the dataset.

2. Scan the pool of candidates and remove any for which the dissimilarity to the new selection is less than  $R$ .
3. Is the pool of candidates empty? If so, quit.
4. Select a compound at random from the pool of candidates.
5. Have  $M_{\max}$  compounds been selected? If so, quit.
6. Go to step 2.

Minimum dissimilarity selection tends to be order dependent.<sup>1</sup> This can be alleviated by setting  $M_{\max}$  very high, so that the algorithm runs to exhaustion. For most datasets, doing so with an empirically justified value for  $R$  (e.g., 0.15–0.2 for Tanimoto dissimilarity of 2D fingerprints<sup>7,8</sup>) will return an undesirably large number of selections. Typically, several minimum dissimilarity runs must be made to find a value of  $R$  which will return the desired number of selections when minimum dissimilarity is run to exhaustion.

If a reasonable value of  $M_{\max}$  and an empirically determined radius  $R$  are used, minimum dissimilarity will return a representative subset differing from random selection only in that there will be no redundant selections, as defined by  $R$ . Such a selection will be representative but may not be diverse enough to satisfy some chemists.

**Hierarchical Clustering.** In agglomerative hierarchical clustering, the most similar pair of clusters are consolidated at each level, starting from one singleton cluster for each compound in the set. Selecting from clusters obtained using Ward's method or complete linkage<sup>1,9</sup> returns subsets which are both representative and diverse, in that each compound is represented by some neighbor and the representatives are distributed across the full breadth of the dataset. Medicinal chemists generally find selections based on hierarchical clustering intuitively appealing and natural, especially after they have eliminated "oddball" clusters and singleton compounds from the selection list. Indeed, they sometimes make their own selections by manually clustering structures. By examining the dissimilarity between the most similar clusters at each step, one can identify the most "natural" level of clustering near the number of compounds one wishes to select—i.e., levels at which the last clusters consolidated were substantially more similar than any remaining clusters are to each other.

Hierarchical clustering is not always a practical option, however. The speed of the classical technique scales with the cube of the size  $N$  of the dataset<sup>1,9</sup> and so becomes slow for large libraries. In addition, memory requirements generally restrict direct applications to relatively small datasets ( $\leq 2000$  compounds). Faster approaches, including reciprocal nearest neighbors (RNN),<sup>9</sup> are available which relieve the memory limitations and can usually reduce scaling problems dramatically.<sup>1</sup> Unfortunately, the scaling benefits can only be fully realized when centroids are well-defined and well-behaved in the metric space being explored, which is not the case for some important metrics of interest—in particular, for Tanimoto similarities between 2D fingerprints.<sup>7,8</sup>

**Optimizable  $K$ -Dissimilarity.** Holliday and Willett have pointed out<sup>4</sup> that their own<sup>10</sup> and other<sup>3</sup> dissimilarity-based selection methods were actually specific manifestations of a more general, unified method, much as Lance and Williams<sup>11</sup> had done for hierarchical clustering. In much the same way, maximum and minimum dissimilarity selection can be reformulated as limiting cases of a single, more

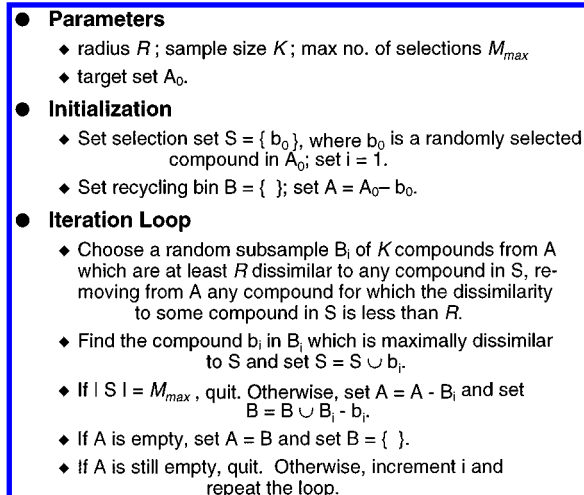


Figure 1. The Optimal Dissimilarity Selection algorithm.

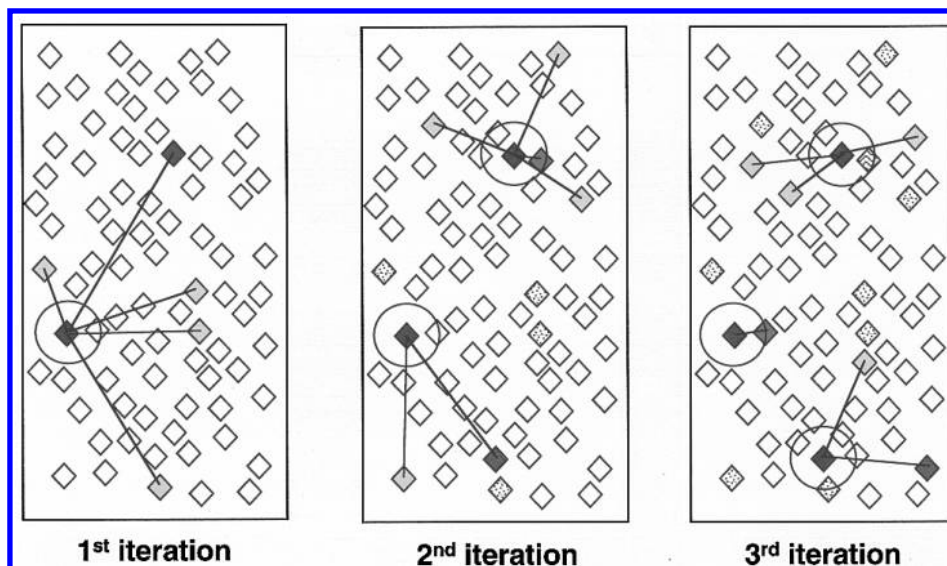
general algorithm. The generalization entails introduction of a parameter  $K$  which defines a subsample size at each iteration:

1. Select a compound at random from the dataset of interest and create a list of candidate compounds from the remainder of the dataset. Create an empty recycling bin and subsample set.
2. Remove a compound at random from the candidate pool. If it has a dissimilarity less than  $R$  with respect to any of those already selected, discard it. Otherwise, add it to the subsample.
3. Repeat step 2 until the subsample includes  $K$  compounds or until the candidate pool is exhausted.
4. If there are fewer than  $K$  compounds in the subsample and the candidate pool is exhausted, then remove all compounds from the recycling bin, put them into the candidate pool, and go to step 2.
5. If the subsample is empty, quit.
6. Examine the subsample and identify the "best" compound, e.g., the one most dissimilar to those already selected.
7. Remove the best compound from the subsample and add it to the selection set.
8. Remove from the subsample those compounds which were not selected and put them into the recycling bin.
9. Have  $M_{\max}$  compounds been selected? If so, quit. Otherwise, return to step 2 to start a fresh subsample.

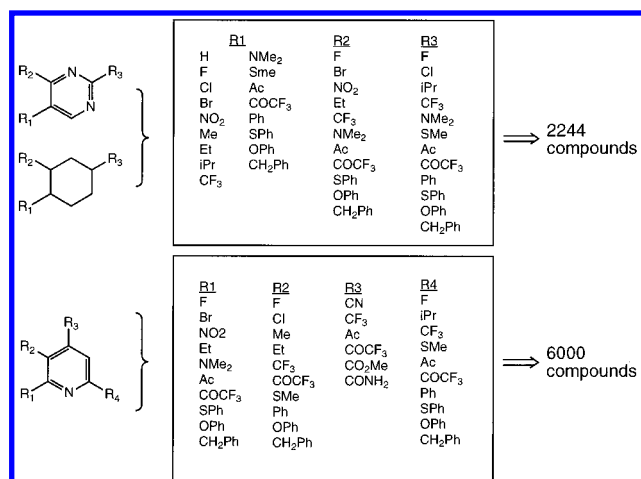
A more concise, mathematical version of the algorithm is given in Figure 1. The results of its application for three iterations are illustrated schematically in Figure 2, where blue symbols indicate selected compounds, green indicates the new selection made at each iteration, and yellow symbols show subsample compounds considered but not selected.

Clearly, maximum and minimum dissimilarity selection algorithms can be seen simply as extreme instances of this *Optimizable  $K$ -Dissimilarity Selection* algorithm—OptiSim, for short. For maximum dissimilarity,  $K$  is effectively  $N$ , the number of elements in the dataset: all compounds are considered as candidates at each step. The only substantive change involved is that the first two selections are not dropped (step 4 in the maximum dissimilarity algorithm). Minimum dissimilarity is simply the special case of  $K = 1$ .

By choosing an intermediate value of  $K$ , one strikes a balance along the continuum between the diversity of



**Figure 2.** Illustrative results from the first three iterations of an OptiSim run ( $R = 6$  mm (original scale),  $K = 5$ ,  $M_{\max} \geq 3$ ). Symbols representing compounds which have already been selected are blue. Subsamples are indicated by yellow or green symbols, with green indicating the compounds selected from each subsample. The circles around each blue symbol show the minimum dissimilarity radius  $R$ ; red indicates compounds examined at each iteration which are excluded as too similar to those which have already been selected. Stippled symbols indicate candidates in the recycling bin, i.e., which have already been considered for selection once.



**Figure 3.** Design of the parent libraries contributing to the combinatorial dataset.

maximum dissimilarity and the representativeness of minimum dissimilarity—hence the inclusion of “Optimizable” in the name. As shown below, it turns out that one can in some respects also mimic selection based on hierarchical clustering.

## METHODOLOGY

In evaluating a compound selection method, it is important to differentiate between the performance of the method itself and the validity of the particular descriptor to which the method is applied. To avoid confounding these two issues here, independent structural and scalar descriptor sets were constructed. Skewed, nonuniform distributions were chosen in both cases to reflect the kinds of substructure present in many “real world” applications in a controlled form and to maximize the discriminating power of the  $\chi^2$  statistic used in comparing selection methods.

**Combinatorial Dataset Design.** The Legion<sup>12</sup> combinatorial builder module in SYBYL<sup>12</sup> was used to create a homologous set of libraries, each with a pyrimidine, a cyclohexane, or a pyridine at the core (Figure 3) and an analogous pattern of substitution around each ring. The

pyrimidine and cyclohexane libraries consisted of 2244 compounds each, whereas the pyridine library was made up of 6000 compounds. A composite library was built from the 6000-compound pyridine dataset, 500 randomly selected cyclohexanes, and 100 randomly selected pyrimidines. The final dataset of 1000 compounds—892 pyridines, 92 cyclohexanes, and 16 pyrimidines—was created by randomly selecting from among the 6600 compounds in the composite library. The Tanimoto dissimilarity  $T^*$  (equivalent to the Soergel distance<sup>13,14</sup>) was used to assess the dissimilarity between any two compounds **a** and **b** in the dataset:

$$T^*(a,b) = 1 - T(a,b) = (|a \cup b| - |a \cap b|) / (|a \cup b|)$$

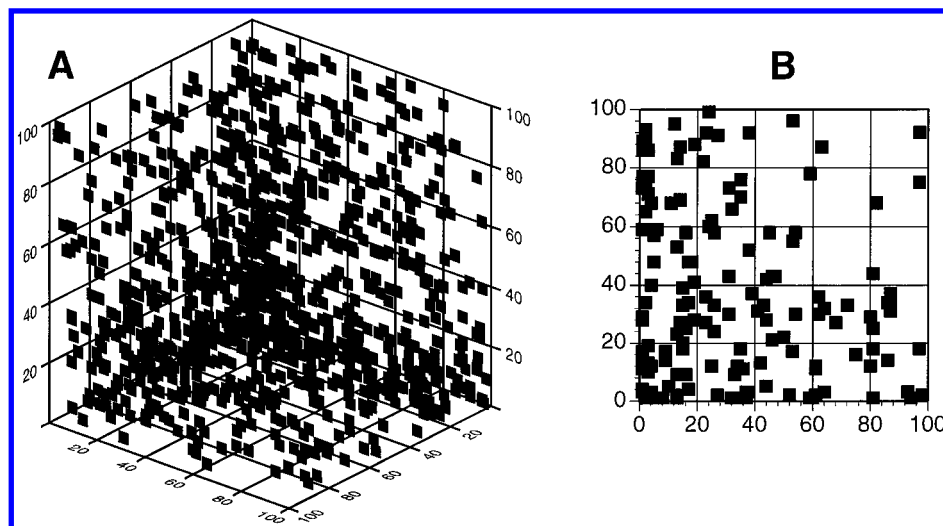
where the vertical bars denote cardinality and **a** and **b** correspond to the fingerprint bit sets for **a** and **b**, respectively. Standard UNITY<sup>12</sup> 2D fingerprints were used for evaluating dissimilarities.

Three scalar descriptors were generated for the combinatorial dataset by drawing three numbers for each compound from a uniform random population of reals between 0 and 1, then squaring each number, and multiplying it by 99. Adding 1 to the integer part of each value produced three descriptors for each compound with values between 0 and 100. These were distributed independently of the 2D structure and, hence, of the corresponding fingerprints. The skewed distribution of values resulting from the squaring operations produced a gradient of density in a cubical descriptor space 100 units on a side (Figure 4), running out from the concentration of points near the origin. Dissimilarity was evaluated in terms of Euclidean distance for these scalar descriptors.

**$\chi^2$  Comparisons.** Subsets generated by Optimal Dissimilarity Selection at various subsampling levels  $K$  were characterized nonparametrically in terms of the  $\chi^2$  statistic

$$\chi^2 = \sum (O_i - E_i)^2 / E_i \quad (1)$$

where  $O_i$  denotes the observed count obtained from the  $i$ th of  $c$  clusters and  $E_i$  denotes the count expected from that



**Figure 4.** Distribution of scalar descriptors for the combinatorial dataset: (A) data for all 1000 compounds is displayed in three dimensions and (B) distribution for 150 compounds in the XY plane.

cluster<sup>15</sup> (the term “cluster” is used here because the comparisons of most immediate interest are between cluster-based selection and OptiSim results; “category” or “class” would be equally appropriate).

The  $\chi^2$  statistic with respect to a **random** sample is a measure of how representative a particular selection set is. In general, the expected count for random sampling is given by

$$\bar{E}_i(\text{random}) = bMn_i/N \quad (2)$$

where  $b$  is the number of trials;  $M$  is the number selected per trial;  $n_i$  is the number of compounds in the  $i$ th cluster; and  $N$  is the total number of entities being selected from. For random sampling, one expects to select most often from the most populous cluster and to select proportionately less often from smaller clusters. The larger a selection set's value of  $\chi^2$  (random), the more it diverges from being representative of the dataset as a whole.

Note that the OptiSim algorithm explicitly precludes reselecting any compound, whereas the random selection distribution is for sampling with replacement. As a result, the selections are not strictly independent of each other, and eq 2 is not exact. This is not a problem if the number of selections per trial does not greatly exceed the number of clusters. It is then necessary, however, to block trials if one is to keep the expected number of selections for each cluster large enough to avoid having to make a continuity correction to the  $\chi^2$  statistic calculated in (1).<sup>15</sup>

If selection is perfectly **uniform** across clusters, each cluster will be equally likely to be sampled. How uniformly a selection is distributed across  $c$  clusters, then, is a measure of how similar a result is to cluster-based selection. The  $\chi^2$  (uniform) statistic is therefore calculated from

$$\bar{E}_i(\text{uniform}) = bM/c \quad (3)$$

Again, perfect concordance gives a  $\chi^2$  of 0. The smaller  $\chi^2$  (uniform) is, the better the result mimics cluster-based selection.

**Scaling.** As noted above, a random selection distributed perfectly proportionately across all clusters will have a  $\chi^2$  (random) of 0, and a perfectly uniform selection will have a  $\chi^2$  (uniform) of 0. Both results, however, are quite unlikely

when selections are made independently. For either statistic, the mean  $\chi^2$  expected by chance is equal to the degrees of freedom ( $df = c - 1$ ). Scaling by this population mean makes it easier to compare experiments which involve different numbers of clusters (i.e., categories of classification), since it makes the expected result equal to 1 no matter how many clusters are involved.

**Hierarchical clustering** was carried out in the SYBYL QSAR module or in Selector.<sup>12</sup> Complete linkage<sup>9</sup> was used as the method of choice because it produces compact clusters of relatively uniform diameter.<sup>16,17</sup> This minimizes dissimilarities within clusters, which makes selections from complete linkage clusterings generally more representative than those based on other linkage methods. In some cases, group average hierarchical clusterings were run for comparison.

By its nature, hierarchical clustering can return any desired level of clustering. “Natural” levels for the artificial descriptor sets were determined by examining the relative dissimilarities between clusters consolidated at each stage of agglomeration. Those levels were chosen which gave a workable number of clusters while minimizing dissimilarities within clusters relative to dissimilarities between clusters.

## RESULTS AND DISCUSSION

**Scalar Descriptors.** Hierarchical clustering on the three scalar descriptors generated for the 1000 compounds in the combinatorial dataset resolved it into ten clusters when the complete linkage method was used. Two clusters split one corner of the cubical descriptor space; the rest cover the remaining corners and the center. The ten clusters obtained were made up of 283, 167, 134, 130, 90, 78, 45, 29, 24, and 20 compounds, respectively. Not surprisingly, the largest cluster is near the origin, and the smallest is at the opposite vertex of the cubical boundaries of the descriptor space.

Optimal Dissimilarity Selection was then applied to the dataset 20 times at each of seven different subsampling rates  $K$ , selecting ten compounds ( $M_{\max} = 10$ ) each time;  $R$  was set to 10. The same random number was used to “seed” the algorithm for one trial at each subsampling rate, so any biases due to the initial selection were spread equally across all values of  $K$ . The number of selections from each cluster



**Table 1.** Divergence of Subsets Selected from the Combinatorial Dataset from Random and Uniform Distributions across Clusters<sup>c</sup>

$K^a$	scalars (df = 9)		fingerprints (df = 20)	
	random <sup>b</sup>	uniform	random	uniform
1	1.23 ± 0.43 <sup>c</sup>	7.21 ± 0.14	0.89	15.05
5	2.06 ± 0.50	3.56 ± 0.51	4.04	7.36
10	2.58 ± 0.11	2.44 ± 0.11	11.15	5.56
15	5.04 ± 0.66	1.98 ± 0.24	13.75	3.94
25	5.88 ± 1.16	1.93 ± 0.38	20.32	3.38
35	7.61 ± 2.89	2.02 ± 0.07	26.74	2.08
1000 <sup>d</sup>	7.59 ± 0.72	0.67 ± 0.15	49.84	2.05

<sup>a</sup> Subsampling rate at each iteration; see Figure 1 for definition.<sup>b</sup> Reference distribution. <sup>c</sup> Mean ± SEM for two blocks. <sup>d</sup> All unselected compounds considered at each iteration. <sup>e</sup> Values cited are in terms of scaled chi square ( $\chi^2/\text{degrees of freedom}$ ).

was summed across a block of ten such trials, so that each block included 100 selections for each subsample size. The  $\chi^2$  statistics obtained were then averaged across two such blocks. The total number of scalar selections was therefore

$$(10/\text{trial}) \times (10 \text{ trials/block}) \times (2 \text{ blocks}) \times (7 \text{ subsampling rates}) = 1400$$

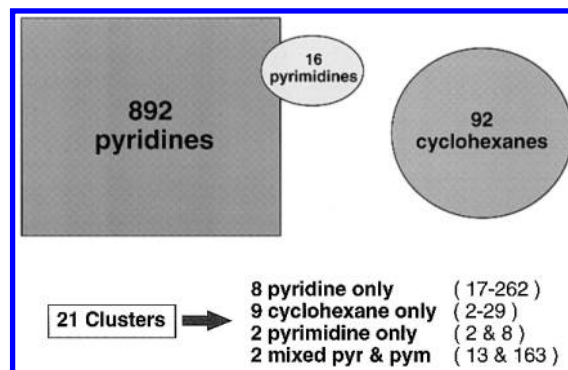
in 140 trials, with a total of 200 selections made for each subsampling rate.

The distributions of these selections across clusters were then compared to the totals expected for random sampling and to the totals expected for uniform selection from the ten clusters. Table 1 shows the  $\chi^2$  values obtained as averages across two blocks for each value of  $K$ . The  $\chi^2$  (random) for minimal dissimilarity selection ( $K = 1$ ) is not significantly different from that expected by chance (1.23 vs 1.0), but it rises steadily with increasing  $K$  as the selected subsets grow more diverse and less representative. The  $\chi^2$  (uniform) profile, on the other hand, falls quickly with increasing  $K$  to a plateau value of about 2 for  $K = 15$ –35. The maximal dissimilarity extreme ( $K = 1000$ ) produces a significantly more uniform distribution (0.67 vs ~2.0), which reflects the fact that a cluster is, in this case, located at each of the first nine linear D-optimal points of the descriptor space—the corners plus the center.

Using the group average method gave eight clusters, one at each corner of the cubical descriptor space. The size distribution for these clusters were similar to those obtained using complete linkage, as were the  $\chi^2$  profiles obtained for OptiSim selections with  $M_{\text{max}} = 8$ .

Note: maximum dissimilarity selection ( $K = 1000$ ) performs unusually well here with respect to hierarchical clustering in part because the dimensionality (3) is much smaller than  $M_{\text{max}}$  (10), and in part because the peak population density near the origin coincides with one extreme of the descriptor space.

**Combinatorial Fingerprints.** Hierarchical clustering of the combinatorial dataset using 2D fingerprints and complete linkage gave 21 clusters as a “natural” level (see discussion above of hierarchical clustering). These included 19 “pure” clusters made up of 262, 152, 75, 64, 64, 45, 43, or 17 pyridines; 29, 21, 11, 10, 6, 5, 5, 3, or 2 cyclohexanes; and 8 or 2 pyrimidines. Two mixed clusters contained both pyridines and pyrimidines—161 and 2 or 9 and 4, respectively. With few exceptions, shared structural features within each cluster were quite easy to identify: one cluster was

**Figure 5.** Relationships among compound classes in the combinatorial dataset as shown by hierarchical clustering on UNITY 2D fingerprints using complete linkage.<sup>9</sup>

made up entirely of nitrocyclohexanes, for example, another of thiophenylcyclohexanes, and another of 3-phenoxy-pyridines.

The relationships between the sublibraries shown by this clustering pattern are illustrated schematically in Figure 5. Note that the areas of each set shown in Figure 5 indicate the *number* of compounds of that class which are in the dataset, not the degree of diversity within each class. Because they are drawn from homologous combinatorial libraries, the degree of structural variety found within each class is similar. This is reflected in the similar numbers of pyrimidine and cyclohexane clusters; that fewer pyrimidine clusters were identified simply reflects the low population of this class.

The group average method was also applied, but it gave a more skewed and less useful distribution of cluster sizes. In particular, more than half of the compounds (561) were assigned to a single cluster of pyridines, and many small clusters were produced, including four pyrimidine singletons.

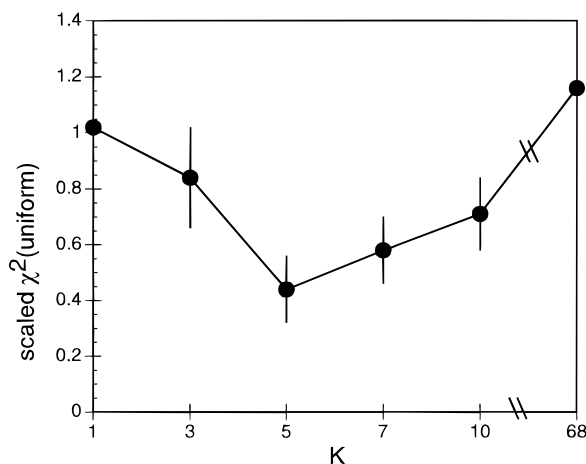
Twenty-one compounds were selected in each trial ( $M_{\text{max}} = 21$ ), and distributions across clusters were summed across 10 trials for each block at each of seven values of  $K$ ;  $R$  was set to 0.15.<sup>7,8</sup> Again, there was one trial at each subsampling rate for each random number seed used; the seeds used were different from those used for analyzing the associated scalar descriptors. In this case, blocks were not replicated. Hence the total number of trials was

$$(21/\text{trial}) \times (10 \text{ trials/block}) \times (7 \text{ subsampling rates}) = 1470$$

in 70 trials, with a total of 210 selections made at each of seven values of  $K$ .

The values of scaled  $\chi^2$  found with respect to random and uniform complete linkage clusters are shown in Table 1. Again, the OptiSim selections move away from being purely representative and begin to resemble cluster-based selections even at low values of  $K$ . Note that under this high-dimensional, non-Euclidean metric the limiting scaled  $\chi^2$  (uniform) is 2.0.

Maximum dissimilarity ( $K = 1000$ ) returned the most uniformly distributed selection set. Note, however, that the number of superclusters (3—one for each combinatorial core structure) is small compared to  $M_{\text{max}}$  (here, 21). The distribution of compounds chosen using maximum dissimilarity selection can be uncharacteristically uniform in such a situation.



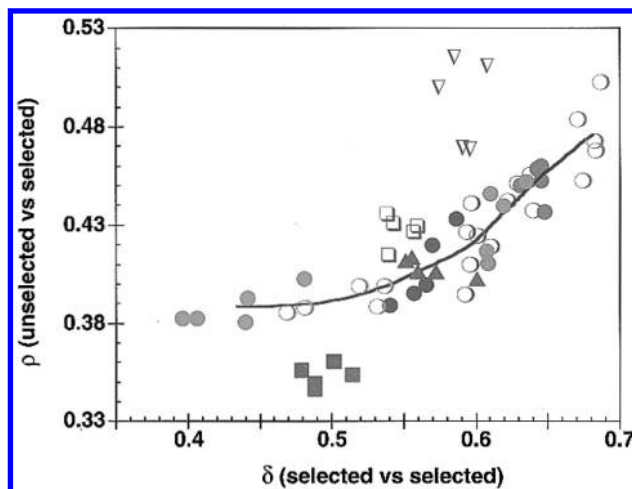
**Figure 6.** Divergence of OptiSim selections from uniform sampling across pharmacological classes for the Mannhold–Rekker dataset<sup>18</sup> as a function of subsample size  $K$ . The  $\chi^2$  statistic is scaled by division by its 5 degrees of freedom.

**Mannhold–Rekker Dataset.** The combinatorial dataset described above is very structured and artificial. It was created so deliberately in an effort to keep evaluation of the selection algorithm—OptiSim—separate from considerations of the appropriateness of the particular metric being used. Nonetheless, it is important to get some idea of how well the methodology performs with more realistic datasets. Mannhold *et al.*<sup>18</sup> compiled a database of 68 structurally diverse compounds from six pharmacological classes for evaluating different ways of predicting a compound's octanol/water partition coefficient  $P$  from its structure. The dataset includes 15 class I antiarrhythmics, 13 phenothiazines, 12 class III antiarrhythmics, 11  $\beta$ -blockers, 9 benzamides, and 8 potassium channel openers. OptiSim was used to select six compounds from the dataset in nine trials with  $K$  set to 1, 3, 5, 7, 10, or 68. Again, standard UNITY 2D fingerprints were used to evaluate the dissimilarity between molecules.

In this case, the uniform reference distribution was based on pharmacological classes. Because there is relatively little variation in population between classes in this dataset, however, the analysis applied above to the combinatorial dataset is not very informative, particularly *vis à vis* random sampling. Instead,  $\chi^2$  (uniform) was calculated for each trial, and the results averaged across the nine trials for each subsampling rate  $K$ . The results obtained are shown in Figure 6; selection of one example from each pharmacological class was best reproduced at  $K = 5-7$ .

Note that here, where the number of completely distinct structural classes in the dataset (six) is comparable to  $M_{\max}$ , intermediate values of  $K$  outperform maximum dissimilarity. In fact, OptiSim performs somewhat better with respect to the pharmacological classes in this dataset than does cluster-based selection based on 2D fingerprints and complete linkage (data not shown), in part because the variation in structures within classes is uneven. Bear in mind that this dataset is rather small, however, so the results obtained are best regarded as qualitative.

**Parametric Measures of Diversity and of Representativeness.** The  $\chi^2$  statistic is a good measure of similarity for validation work, but it requires a reference distribution and so is of limited usefulness when no hierarchical classification is available for comparison. If OptiSim is to



**Figure 7.** Relationship between the average dissimilarity  $\rho$  between unselected compounds and sets of 21 compounds selected by OptiSim as a function of the average dissimilarity  $\delta$  among selected compounds. One of two sets of 100 compounds randomly selected from the combinatorial dataset was used to calculate  $\rho$  for each trial. The data shown are for UNITY 2D fingerprints from the combinatorial dataset, with each point representing an average of across two trials: (●, ○) data points from OptiSim selections, color-coded by  $K = 1, 2, 5, 10, 15, 25, 35$  or 1000, with filled and open symbols alternating; (■, □) selections made from hierarchical complete linkage<sup>9</sup> clusters; (▲, ▽) selections made from hierarchical group average<sup>9</sup> clusters; (□, ▽) one compound was randomly selected from each cluster; (■, ▲) one compound was randomly selected from among the three most central<sup>20</sup> compounds in each cluster; and (—) spline curve drawn through the mean values of  $\rho$  and  $\delta$  for each value of  $K$ .

be used as an alternative to hierarchical clustering, more readily accessible measures will be needed to know when the optimal balance between representativeness and diversity has been obtained for a given subsample size in any particular application. Such measures also provide insight into why  $\chi^2$  (uniform) for OptiSim fails to go to 1 as  $K$  increases.

It is convenient to use averages when characterizing large datasets, because the law of large numbers guarantees that a good estimate of the average can usually be obtained from a random sample of the entire dataset.<sup>19</sup> Let  $S$  be the set of  $M$  compounds selected and  $U$  be a set of  $n$  compounds chosen at random from among the compounds which were not selected. Then the average dissimilarity  $\delta$  between each compound in  $S$  and the other compounds in  $S$  is a measure of diversity, whereas the average dissimilarity  $\rho$  between compounds in  $U$  and those in  $S$  is a measure of the representativeness of  $S$ . Here, we are using the minimum pairwise dissimilarity criterion<sup>3,4</sup> for evaluating the dissimilarity between each compound and the reference set, so

$$\delta = (1/M) \sum_{i=1}^M \min(T^*(s_i, s_j) : 1 \leq j \leq M, i \neq j) \quad (4)$$

$$\rho = (1/n) \sum_{i=1}^n \min(T^*(u_i, s_j) : 1 \leq j \leq M) \quad (5)$$

Figure 7 shows a plot of  $\rho$  as a function of  $\delta$  for OptiSim selections based on fingerprints for the combinatorial dataset; results at each of seven levels of subsample size  $K$  are shown (circles in Figure 7). For comparison, data are also shown for selections based on hierarchical clusters obtained using either complete linkage or the group average method (square

and triangular symbols, respectively). For the cluster-based selections, one compound was taken at random from each cluster (open square and triangular symbols) or from among the three most central<sup>20</sup> compounds in each cluster (closed square and triangular symbols).

The scatter in  $\rho$  and  $\delta$  among OptiSim selection sets obtained at the same value of  $K$  reflects the stochastic nature of the algorithm. The variance in  $\delta$  falls sharply as  $K$  increases, with a standard deviation (SD) of 0.042, 0.017, and 0.007 at  $K = 1, 5$ , and  $35$ , respectively. The variance in  $\rho$ , on the other hand, is primarily determined by the randomly chosen *unselected* compounds (**U**) with which the selection sets are compared and so decreases only slightly with increasing  $K$ . For any particular dataset and a given  $M_{\max}$ , there is a characteristic limiting value for  $\delta$ . In this case,  $\delta_{\max}$  is slightly less than 0.7. Note that  $\rho$  and  $\delta$  are both smallest (on average) when  $K = 1$ . Moreover, the expected values for  $\rho$  and  $\delta$  for OptiSim selections will be equal if and only if **S** and **U** are the same size,  $R = 0$ , and  $K = 1$ .

Diversity ( $\delta$ ) increases with increasing  $K$  but does so at the cost of a decrease in representativeness (increase in  $\rho$ ). As expected, maximum dissimilarity ( $K = 1000$ ) gives more diverse but less representative selections than does minimum dissimilarity ( $K = 1$ ). Note, however, that most of the increase in diversity is realized for subsamples which are relatively small and still quite representative ( $K \leq 15$ , in this instance); most of the loss in representativeness (increase in  $\rho$ ) occurs at higher values of  $K$ .

Among the cluster-based selection methods examined here, taking central compounds from clusters obtained using complete linkage (filled squares in Figure 7) gave the most representative but least diverse selections. As anticipated, using the group average method gave more diverse but less representative selections than did complete linkage. For both cluster-based selection methods, random sampling from clusters decreased representativeness but increased diversity, particularly for complete linkage.

Selecting central compounds from group average clusters (filled triangles in Figure 7) returned selection sets with distributions qualitatively similar to those obtained with OptiSim for subsample sizes  $K = 5$  and  $10$ . The latter selections were actually slightly better—both more representative (lower  $\rho$ ) and more diverse (higher  $\delta$ )—than were selections drawn randomly from complete linkage clusters (open squares), whereas  $K = 10$  to  $35$  gave better results than did random selection from group average clusters (open triangles).

Except at  $K = 1$ , all OptiSim selections were less representative but more diverse than were selections of central compounds from complete linkage clusters. This *qualitative* difference in how the selections are distributed accounts for the failure of  $\chi^2$  (uniform) to fall much below 2 for this dataset when complete linkage clustering is used to obtain the reference distribution. Note, however, that selection sets may be equally diverse and representative yet still be quite different from each other.

Qualitatively similar results were obtained for the scalar descriptors, though in this case the differences in representativeness across the different selection methods were much smaller (data not shown).

Ward's method is an alternative hierarchical clustering method which uses a probabilistic, least-squares rationale

to maximize distances between clusters while minimizing distances within clusters under Euclidean metrics.<sup>9,17</sup> Evidently the stochastic nature of the OptiSim algorithm imparts similar properties to its selection sets. This is potentially a quite useful generalization, since Ward's method is not applicable to metric spaces in which centroids of clusters are not well-defined—in particular, it is not directly applicable to Tanimoto coefficients based on 2D fingerprints or other bit set descriptors.

## VARIATIONS AND EXTENSIONS OF THE BASIC ALGORITHM

The heart of the Optimizable K-Dissimilarity Selection algorithm lies in taking a series of random subsamples from a dataset and selecting the best candidate from each subsample, where "best" is defined by some criterion which is a function of those compounds (or, more generally, elements) which have been selected in previous steps. Several useful variations on the algorithm employed here immediately suggest themselves.

**Other Evaluation Criteria.** OptiSim has been defined here solely in terms of maximum minimum pairwise dissimilarity to those compounds already selected as the criterion by which the best compound is to be selected from each subsample. The highest average dissimilarity<sup>4</sup> could just as well be used, and the cosine coefficient<sup>10</sup> could be substituted for the Tanimoto coefficient. In fact, the algorithm is generalizable to any set of selection criteria or decision rules.

**Sampling with Replacement.** As set out in Figure 1, each OptiSim subsample is drawn from the candidate pool for consideration and then is set aside until all candidates have been examined once—i.e., the dataset is sampled without replacement. If samples are drawn with replacement of those compounds which do not get selected, no recycling bin is required. A given setting of  $K$  will then return a more representative but less diverse subset because sampling of more sparsely populated (outlier) regions will be avoided. The tradeoff is that a particular level of diversity among the compounds selected will only be approached at higher values of  $K$ , which can become computationally expensive.

**Excluding Redundant Compounds from the Subsample.** The implementation described in Figure 1 tests each compound for redundancy before putting it into a subsample. An alternative approach is to select  $K$  compounds from the candidate pool and apply the redundancy test only to the best compound; if it is redundant—that is, if it is too similar to those which have already been selected—no selection from that subset is made. This approach can be made faster "up front" than the version of OptiSim set out in Figure 1 for some descriptors and similarity measures but will be correspondingly slower at later selection steps. In addition, the balance between representativeness and diversity will be shifted toward making more representative selections, just as subsampling with replacement will.

**Dataset Clustering.** As demonstrated here, OptiSim selection sets behave in many ways like selection sets based on hierarchical clustering. Indeed, OptiSim selections can be used as centers (i.e., as leaders<sup>1</sup>) for efficiently clustering large datasets on the basis of a secondary similarity radius,  $R'$ , or by assignment of each compound to the most similar center. Moreover, selected compounds can themselves be

submitted to hierarchical clustering. Under this scenario, the OptiSim selections will be true centers of their  $R'$  neighborhoods<sup>8</sup> for any metric, so their hierarchy will perforce accurately reflect hierarchical relationships across the entire dataset.

### CONCLUSION

OptiSim is a generalized dissimilarity selection algorithm which includes the established methods of maximum and minimum dissimilarity selection as special cases. By varying the subsample size  $K$ , one can adjust the balance between how representative the selected subset is and how diverse it is. Intermediate settings can mimic the results obtained for selection based on hierarchical clustering and, at least in some cases, improve upon them.

### ACKNOWLEDGMENT

The author would like to thank John Begemann and Paul Weber of Tripos, Inc. for their support of the work described here; Prof. Peter Willett of the University of Sheffield for his many helpful comments on this work; and the reviewers, particularly for their suggestion of the term " $K$ -dissimilarity selection". Nicole Van Opdenbosch, Dick Cramer, Bill Langton, Trevor Heritage, Tad Hurst, and Jon Swanson, also of Tripos, Inc., also provided helpful suggestions for the manuscript.

### REFERENCES AND NOTES

- (1) Barnard, J. M.; Downs, G. M. Clustering of chemical structures on the basis of two-dimensional similarity measures. *J. Chem. Inf. Comput. Sci.* **1992**, 32, 644–649.
- (2) Registered trademark; patent pending.
- (3) Lajiness, M.; Johnson, M. A.; Maggiora, G. M. Implementing drug screening programs using molecular similarity methods. In *QSAR: Quantitative Structure-Activity Relationships in Drug Design*; Fauchere, J. L., Ed.; Alan R. Liss, Inc.: New York, 1989; pp 173–176.
- (4) Holliday, J. D.; Willett, P. Definitions of "dissimilarity" for dissimilarity-based compound selection. *J. Biomol. Screening* **1996**, 1, 145–151.
- (5) Agrafiotis, D. K. Stochastic algorithms for maximizing molecular diversity. *3rd Electronic Computational Chemistry Conference* **1996**.
- (6) See, for example: Brannigan, L. H.; Duewer, D. L. Experimental design in the Development of biologically active compounds. *Pharmacochem. Libr.* **1991**, 16, 553–556.
- (7) Brown, R. D.; Martin, Y. C. Use of structure-activity data to compare structure-based clustering methods and descriptors for use in compound selection. *J. Chem. Inf. Comput. Sci.* **1996**, 36, 572–584.
- (8) Patterson, D. E.; Cramer, R. D.; Ferguson, A. M.; Clark, R. D.; Weinberger, L. E. Neighborhood behavior: a useful concept for validation of molecular diversity descriptors. *J. Med. Chem.* **1996**, 39, 3049–3059.
- (9) Murtagh, F. A survey of recent advances in hierarchical clustering algorithms. *Comput. J.* **1983**, 26, 354–359.
- (10) Holliday, J. D.; Ranade, S. S.; Willett, P. A fast algorithm for selecting sets of dissimilar molecules from large chemical databases. *Quant. Structure-Activity Rel.* **1996**, 14, 501–506.
- (11) Lance, G. N.; Williams, W. T. A. A general theory of classificatory sorting strategies. I. Hierarchical systems. *Comput. J.* **1967**, 9, 373–380.
- (12) Legion, SYBYL, UNITY, and Selector are available from Tripos, Inc., 1699 S. Hanley Road, St. Louis, MO 63144.
- (13) Gower, J. C. Measures of similarity, dissimilarity, and distance. In *Encyclopedia of Statistical Sciences*; Kotz, S., Johnson, N. L., Eds.; John Wiley & Sons: New York, 1985; Vol. 5, pp 397–405.
- (14) Cheng, C.; Maggiora, G.; Lajiness, M.; Johnson, M. Four association coefficients for relating molecular similarity measures. *J. Chem. Inf. Comput. Sci.* **1996**, 36, 909–915.
- (15) Conover, W. J. *Practical Non-parametric Statistics*, 2nd ed.; John Wiley & Sons: New York, 1980; pp 143–170.
- (16) Cramer, R. D.; Clark, R. D.; Patterson, D. E.; Ferguson, A. M. Bioisosterism as a molecular diversity descriptor: steric fields of single topomeric conformers. *J. Med. Chem.* **1996**, 39, 3060–3069.
- (17) Kaufman, L.; Rousseeuw, P. J. In *Finding Group in Data: An Introduction to Cluster Analysis*; Wiley-Interscience: New York, 1990; pp 230–243.
- (18) Mannhold, R.; Rekker, R. E.; Sonntag, C.; ter Laak, A. M.; Dross, K.; Polymeropoulos, E. E. Comparative evaluation of the predictive power of calculation procedures for molecular lipophilicity. *J. Pharm. Sci.* **1995**, 84, 1410–1419.
- (19) Mood, A. M.; Graybill, F. A.; Boes, D. C. *Introduction to the Theory of Statistics*, 3rd ed.; McGraw-Hill: New York, 1974; p 232.
- (20) Molecular Diversity Manager Manual, Version 6.3; Tripos, Inc.: St. Louis, 1996; pp 212–213.

CI970282V