

Molecular Basis Sets—A General Similarity-Based Approach for Representing Chemical Spaces

Akshay S. Raghavendra^{†,‡} and Gerald M. Maggiora^{*,‡,§}

Electrical & Computer Engineering, BIO5 Institute, and College of Pharmacy, University of Arizona,
Tucson, Arizona 85721

Received December 12, 2006

A new method, based on generalized Fourier analysis, is described that utilizes the concept of “molecular basis sets” to represent chemical space within an abstract vector space. The basis vectors in this space are abstract molecular vectors. Inner products among the basis vectors are determined using an *ansatz* that associates molecular similarities between pairs of molecules with their corresponding inner products. Moreover, the fact that similarities between pairs of molecules are, in essentially all cases, nonzero implies that the abstract molecular basis vectors are nonorthogonal, but since the similarity of a molecule with itself is unity, the molecular vectors are normalized to unity. A symmetric orthogonalization procedure, which optimally preserves the character of the original set of molecular basis vectors, is used to construct appropriate orthonormal basis sets. Molecules can then be represented, in general, by sets of orthonormal “molecule-like” basis vectors within a proper Euclidean vector space. However, the dimension of the space can become quite large. Thus, the work presented here assesses the effect of basis set size on a number of properties including the average squared error and average norm of molecular vectors represented in the space—the results clearly show the expected reduction in average squared error and increase in average norm as the basis set size is increased. Several distance-based statistics are also considered. These include the distribution of distances and their differences with respect to basis sets of differing size and several comparative distance measures such as Spearman rank correlation and Kruskal stress. All of the measures show that, even though the dimension can be high, the chemical spaces they represent, nonetheless, behave in a well-controlled and reasonable manner. Other abstract vector spaces analogous to that described here can also be constructed providing that the appropriate inner products can be directly evaluated as is the case in this work, a problem that is well-known in kernel-based machine learning.

1. INTRODUCTION

The notion of chemical space has been gaining wider recognition due to rapid growth in the size and availability of compound databases.¹ Chemical spaces provide an intuitive and conceptual basis for understanding many relationships among diverse sets of compounds.² Characterizing chemical spaces is, however, a nontrivial task since an invariant description of them does not exist—different molecular representations yield different chemical spaces that may not be related to each other in simple ways. Thus, the representation problem is an issue of paramount importance in the construction of useful chemical spaces.

The present work describes a general, similarity-based approach to the representation of chemical spaces, which is analogous to the *generalized Fourier series* methods used in vector space theory.³ Here, however, the vectors are associated with molecules. First, a linearly independent set of *abstract molecular basis vectors* is selected that approximately spans the chemical space of interest. The set is then orthonormalized by a nonsingular, linear transformation derived from the metric matrix of the basis vectors. The

transformation, called *symmetric orthogonalization*,⁴ ensures that the new basis set is as close as possible, in a least-squares sense, to the original basis.⁵ In addition, carrying out the orthonormalization procedure in this manner provides a check on the presence of approximate basis set dependencies, which can then be removed before the final molecular basis set is determined.

A novel feature of the current method is the manner in which an orthonormal basis set is obtained. First, a metric matrix⁶ is constructed, whose elements are inner products taken over the set of abstract molecular vectors. Since the detailed form of the molecular vectors is unknown, calculation of the elements of the metric matrix would not seem to be possible. This apparent difficulty can, however, be circumvented by applying an *ansatz* that equates the inner product between a pair of molecular vectors with the molecular similarity between the corresponding molecules. The norm of the molecular vectors is unity since the similarity of a molecule with itself is also unity. This procedure is completely in accord with recent works in kernel-based classifier methods that adopt a similar strategy based on the *direct* construction of the inner products without the *explicit* need for feature vectors.^{7–9} Importantly, molecular similarities satisfy the properties of mathematical kernels including positive definiteness.^{7–9} A benefit of this approach is that any bona fide similarity method can be used to

* Corresponding author e-mail: maggiora@pharmacy.arizona.edu.

[†] Electrical & Computer Engineering.

[‡] BIO5 Institute.

[§] College of Pharmacy.

evaluate the similarities (vide infra). For example, Tanimoto similarity coefficients¹⁰ can be evaluated with MACCS keys,¹¹ SMILES strings,¹² or other similar types of fingerprints.¹³ In addition, graph similarity measures,² or 3-D similarity measures such as those developed in MIMIC,¹⁴ Flex-S,¹⁵ and ROCS,¹⁶ to name a few, can also be used. Importantly, the similarity measure employed does influence the nature of the chemical space induced by that measure. For example, most 2-D similarity measures do not take explicit account of stereochemistry. Thus, differences among molecules due purely to stereochemistry will not be expressed, and pairs of stereoisomers will appear as identical with a similarity equal to one. This can, of course, be ameliorated using 3-D similarity methods.^{17,18}

Another potential benefit of the current approach is that it can be applied to problem domains other than those associated with chemical spaces, the only requirement being that some suitable form of similarity measure be available. Thus, everything from image and shape analysis, document processing, polynucleotide and protein sequence analysis, and biological systems analysis, to name a few, can, in principle, be treated with the method described in this work.⁷ Moreover, the spaces associated with very general inner products used in kernel-based methods can also be treated by the present approach, opening the possibility that addition types of analysis can be carried out for these methods.

Once a suitable set of orthonormal molecular basis vectors has been determined, molecules within a chemical space can be represented as generalized Fourier series.¹⁹ Since the generalized Fourier coefficients of these expansions are given as inner products, they can also be evaluated using the same *ansatz* as before but with respect to the appropriate molecular similarities. An orthonormal basis possesses a number of properties that facilitate its use in the type of application described here. For example, each of the squared generalized Fourier coefficients is equivalent to the fraction of the molecular vector represented by that basis vector. Thus, summing the squared coefficients over the entire basis set yields the total fraction of the molecular vector described by the set of basis vectors. Since the norm of all of the molecular vectors is unity, the difference between unity and the sum of squared Fourier coefficients is the *absolute* squared basis-set error or “representational error”.²⁰ That such an error occurs is not surprising and is due to the incompleteness of the molecular basis set,³ a practical limitation in applications of generalized Fourier series. To our knowledge, this feature is not shared by other cheminformatic methods used to characterize chemical space. For example, in popular dimensionality reduction methods such as principal components analysis²¹ the error is usually attributed to the reduced dimension of the representation, but even if the full dimension is used, an inherent error remains since the feature space employed may not be sufficient to represent the problem exactly.

In addition to being able to assess the representational error, the current method has several other advantages. For example, as noted above, the squares of the generalized Fourier coefficients have a clear-cut interpretation as the fraction of a given molecular vector associated with that coefficient. Thus, this affords a reasonably straightforward interpretation of the contribution of each molecular basis

vector to the overall description of molecules in the chemical space.

The approach presented here is similar to the *principal coordinate* method described by Gower nearly 40 years ago.²² It is also reminiscent of two earlier methods developed by Klein and co-workers²³ and by Oprea and Gottfries.²⁴ In each of these methods, a relatively small reference set of molecules is chosen as a “basis” for representing the chemical space of interest, in a similar fashion to what is done here, but with two significant differences. First, in the latter two works, the *components* of the molecular vectors in the reference set and in the training/test set are explicitly given in terms of well-known molecular descriptors. In the SIBAR method of Klein et al.,²³ the elements in a given row of the data matrix (typically denoted by “**X**” in most statistical applications) are determined by computing the Euclidean distance of the training/test molecule, corresponding to that row, to all of the molecules in the reference set—the data matrix is then employed to make predictions of biological activity using partial least squares. In the ChemGPS method of Oprea and Gottfries,²⁴ however, the molecules in chemical space are projected onto the principal components generated by a set of ChemGPS reference molecules. Both methods contrast with the current, more general, approach, which uses an *ansatz* that directly relates the inner product between two molecular vectors to their intermolecular similarity. Thus, the components of the molecular vectors in this work are never explicitly determined, which is a distinct advantage since any type of positive definite similarity measure can be used, a feature that significantly extends the applicability of the method.²⁵ Second, unlike in the approach proposed here, the reference vectors used by the two methods are typically nonorthogonal since their corresponding correlation or covariance matrices are, in general, nondiagonal. Thus, the coordinates of the training/test set molecules with respect to the reference set cannot be determined simply as projections on the nonorthogonal basis but must account explicitly for the skew of the basis vectors defining the coordinate system. Not accounting for this can lead to significant errors in the representation if the coordinate axes are significantly skewed (i.e., if the variables are highly correlated).

Most other feature-based methods also use vector representations that are nonorthogonal, as measured by their corresponding nonzero covariances. Hence, the contribution of each term to the overall description of a molecule in chemical space is not possible, unless of course the feature space is orthonormalized, which is typically not the case.²⁶ Another advantage of the current approach is that the units of the coordinate values corresponding to each molecular basis vector are identical. Hence, there is no need to scale their values employing any one of several approaches typically used.²⁷ Moreover, in all of the current feature-based methods, it is not possible to estimate the representational error (vide supra) since the set of features that spans the chemical space (i.e., the complete basis set of features) cannot be determined.

Section 2 provides a detailed account of the mathematical methodology behind the approach, which draws heavily from abstract vector space theory³ and the inner-product kernel functions used in many modern machine-learning methods.^{7–9} The first three subsections provide a simple two-dimensional example based on a skewed (i.e., nonorthogonal) coordinate

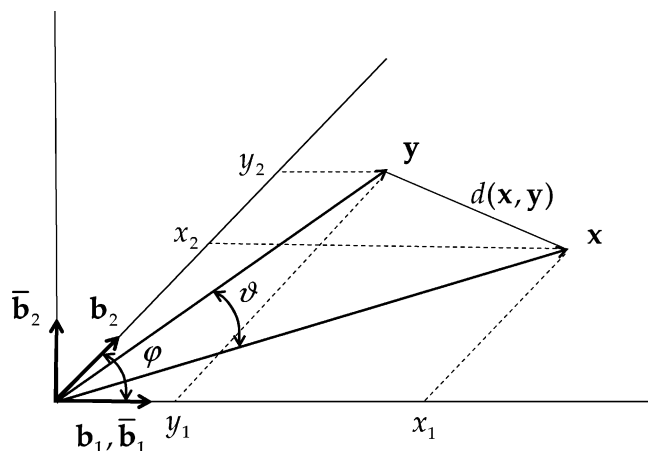


Figure 1. Skewed coordinate system with basis vectors \mathbf{b}_1 and \mathbf{b}_2 , oriented at an angle φ with respect to each other, and the corresponding set of orthonormal basis vectors $\bar{\mathbf{b}}_1$ and $\bar{\mathbf{b}}_2$. Two vectors, \mathbf{x} and \mathbf{y} , oriented at an angle ϑ with respect to each other, and the distance between them, $d(\mathbf{x}, \mathbf{y})$, are also depicted.

system that illustrates the consequences that nonorthogonality has on inner products and distances. The remaining subsections describe how to select molecular basis sets, how to construct an orthonormal molecular basis, how such an orthonormal molecular basis is used to represent molecules in chemical space, and how to determine the representational error. Section 3 describes the set of compounds obtained from the NCI AIDS compound collection,²⁸ which were used to characterize the proposed method. Section 4 provides an example of how the method performs on test data as a function of basis-set size and examines four statistics that assess the suitability of the high-dimensional spaces upon which the current method is based. Section 5 concludes with a summary of the salient features of the method and draws some conclusions regarding its potential strengths and weaknesses. A preliminary report of the current approach was recently published that illustrates many of its features.²

2. THEORY AND METHODOLOGY

2.1 Nonorthogonal (“Skewed”) Coordinate Systems.

Although orthogonal coordinate systems are employed routinely in many exploratory data analyses, the “true” coordinate systems may be inherently nonorthogonal or skewed. As will be seen in the sequel, this can have a major impact on properties such as distances and angles in these spaces. Figure 1 provides an illustration of this for a simple two-dimensional skewed coordinate system. As is seen in the figure, the skewed axes are represented by the basis vectors, \mathbf{b}_1 and \mathbf{b}_2 , respectively, which each have unit *norm*

$$\|\mathbf{b}_i\| = \sqrt{\langle \mathbf{b}_i, \mathbf{b}_i \rangle} = 1 \text{ for } i = 1, 2 \quad (1)$$

where $\langle \mathbf{b}_i, \mathbf{b}_i \rangle$ is an inner product that is given by

$$\langle \mathbf{b}_i, \mathbf{b}_j \rangle = \|\mathbf{b}_i\| \|\mathbf{b}_j\| \cos(\mathbf{b}_i, \mathbf{b}_j) = \|\mathbf{b}_i\| \|\mathbf{b}_j\| \cos \varphi = \cos \varphi \quad (2)$$

Generally, the vectors \mathbf{x} and \mathbf{y} in the $\{\mathbf{b}_1, \mathbf{b}_2\}$ basis are represented as

$$\mathbf{x} = x_1 \mathbf{b}_1 + x_2 \mathbf{b}_2 = (\mathbf{b}_1 \ \mathbf{b}_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \mathbf{b} \mathbf{x}$$

$$\mathbf{y} = y_1 \mathbf{b}_1 + y_2 \mathbf{b}_2 = (\mathbf{b}_1 \ \mathbf{b}_2) \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \mathbf{b} \mathbf{y} \quad (3)$$

where $(x_1 x_2)^T$ and $(y_1 y_2)^T$ are the coordinate (“column”) vectors, \mathbf{x} and \mathbf{y} , of \mathbf{x} and \mathbf{y} in the $\{\mathbf{b}_1, \mathbf{b}_2\}$ basis, and the superscript “T” represents the transpose of the column vector. Similarly to eq 2, the inner product between the two vectors \mathbf{x} and \mathbf{y} is

$$\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x}\| \|\mathbf{y}\| \cos(\mathbf{x}, \mathbf{y}) = \|\mathbf{x}\| \|\mathbf{y}\| \cos \vartheta \quad (4)$$

Equation 4 can also be written in expanded form:

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} \rangle &= \langle \mathbf{b} \mathbf{x}, \mathbf{b} \mathbf{y} \rangle \\ &= \mathbf{x}^T \langle \mathbf{b}, \mathbf{b} \rangle \mathbf{y} \\ &= \mathbf{x}^T \mathbf{G} \mathbf{y} \\ &= (x_1 \ x_2)^T \begin{pmatrix} G_{1,1} & G_{1,2} \\ G_{2,1} & G_{2,2} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \end{aligned} \quad (5)$$

where \mathbf{G} is the *metric matrix*²⁹ whose elements are given explicitly by

$$G_{ij} = \langle \mathbf{b}_i, \mathbf{b}_j \rangle \quad \text{for } i, j = 1, 2 \quad (6)$$

From Figure 1 and the unit norms of the basis vectors, the metric matrix is

$$\mathbf{G} = \begin{pmatrix} 1 & \cos \varphi \\ \cos \varphi & 1 \end{pmatrix} \quad (7)$$

In expanded form, eq 5 becomes

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 + x_2 y_2 + x_1 y_2 \cos \varphi + x_2 y_1 \cos \varphi \quad (8)$$

Thus, the skew of the axes enters directly into and influences the magnitude of the inner product. If the axes are orthogonal, $\varphi = \pi/2$ rad and $\cos \varphi = 0$ and \mathbf{G} becomes the unit matrix \mathbf{I} . In this case, the “off-diagonal” terms equal zero and the inner product reduces to its usual Euclidean form in an orthonormal basis

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 + x_2 y_2 \quad (9)$$

The distance between the two vectors is given by

$$\begin{aligned} d(\mathbf{x}, \mathbf{y}) &= \sqrt{\langle (\mathbf{y} - \mathbf{x}), (\mathbf{y} - \mathbf{x}) \rangle} \\ &= \sqrt{(y_1 - x_1 \ y_2 - x_2) \begin{pmatrix} \langle \mathbf{b}_1, \mathbf{b}_1 \rangle & \langle \mathbf{b}_1, \mathbf{b}_2 \rangle \\ \langle \mathbf{b}_2, \mathbf{b}_1 \rangle & \langle \mathbf{b}_2, \mathbf{b}_2 \rangle \end{pmatrix} \begin{pmatrix} y_1 - x_1 \\ y_2 - x_2 \end{pmatrix}} \\ &= \sqrt{(\mathbf{y} - \mathbf{x})^T \mathbf{G} (\mathbf{y} - \mathbf{x})} \\ &= \sqrt{\sum_{i=1}^2 \sum_{j=1}^2 (y_i - x_i)(y_j - x_j) G_{ij}} \end{aligned} \quad (10)$$

In an orthonormal basis, the expression for the distance reduces to the commonly used expression for the Euclidean distance

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^2 (y_i - x_i)^2} \quad (11)$$

Although eqs 9 and 11 are used routinely in data analysis applications, the assumption of an orthonormal basis is rarely explicitly stated and in many instances is unwarranted by the data.³⁰

2.2. Orthonormalization. Nonorthogonal basis vectors can be orthonormalized in many ways. One of the most common is by *Gram–Schmidt orthonormalization*,³¹ which provides a very intuitive procedure for obtaining orthonormal basis sets. When constructing our orthonormal molecular basis vectors, however, another method called *symmetric orthonormalization*^{3,4} will be employed because it possesses a number of properties that make it more suitable for the current application.³²

Gram–Schmidt orthonormalization is a sequential process. Only two iterations are necessary here since the basis set is two-dimensional. In the initial step, an arbitrary basis vector, say \mathbf{b}_1 , is chosen

$$\bar{\mathbf{b}}_1 = N_1 \mathbf{b}_1, \quad (12)$$

where the normalization constant $N_1 = 1$ since \mathbf{b}_1 is a unit vector. In the second step, a new, normalized basis vector is constructed that is orthonormal to $\bar{\mathbf{b}}_1$, namely,

$$\bar{\mathbf{b}}_2 = N_2(\mathbf{b}_2 - a\mathbf{b}_1) \quad (13)$$

Since $\bar{\mathbf{b}}_2$ is orthonormal to $\bar{\mathbf{b}}_1$ by construction, it satisfies

$$\langle \bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2 \rangle = 0 = N_2(\langle \mathbf{b}_1, \mathbf{b}_2 \rangle - a\langle \mathbf{b}_1, \mathbf{b}_1 \rangle) = N_2(\langle \mathbf{b}_1, \mathbf{b}_2 \rangle - a) \quad (14)$$

and consequently

$$\begin{aligned} a &= \langle \mathbf{b}_1, \mathbf{b}_2 \rangle \\ &= \cos \varphi \end{aligned} \quad (15)$$

Substituting back into eq 12 and using the normalization condition yields

$$\begin{aligned} 1 &= \|\bar{\mathbf{b}}_2\| = \sqrt{\langle \bar{\mathbf{b}}_2, \bar{\mathbf{b}}_2 \rangle} \\ &= N_2 \sqrt{1 - \cos^2 \varphi} \\ &= N_2 \sqrt{1 - \cos^2 \varphi} \\ &= N_2 \sin \varphi \end{aligned} \quad (16)$$

so that $N_2 = 1/\sin \varphi$, and thus, $\bar{\mathbf{b}}_2$ becomes

$$\bar{\mathbf{b}}_2 = \frac{1}{\sin \varphi}(\mathbf{b}_2 - \cos \varphi \mathbf{b}_1) \quad (17)$$

Substituting $\bar{\mathbf{b}}_1$ and $\bar{\mathbf{b}}_2$ into eq 6 for \mathbf{b}_1 and \mathbf{b}_2 , respectively, yields the unit matrix \mathbf{I} , demonstrating that the new basis is, indeed, orthonormal. For larger sets of basis vectors, the Gram–Schmidt process is continued until a complete orthonormal set is obtained.

A little algebraic manipulation shows that the new orthonormal basis is obtained from the original skewed basis by a nonsingular, linear transformation. Specifically,

$$(\bar{\mathbf{b}}_1 \ \bar{\mathbf{b}}_2) = (\mathbf{b}_1 \ \mathbf{b}_2) \begin{pmatrix} 1 & -\frac{\cos \varphi}{\sin \varphi} \\ 0 & \frac{1}{\sin \varphi} \end{pmatrix} \quad (18)$$

or, more succinctly,

$$\bar{\mathbf{b}} = \mathbf{b}\mathbf{T} \quad (19)$$

Clearly, \mathbf{T} is a nonorthogonal transformation since

$$\mathbf{T}^T \mathbf{T} \neq \mathbf{T} \mathbf{T}^T \neq \mathbf{I} \Rightarrow \mathbf{T}^{-1} \neq \mathbf{T}^T \quad (20)$$

This is not surprising since a nonorthogonal set of basis vectors cannot be transformed into an orthogonal set by an orthogonal transformation, which preserves angles and distances.

2.3. Representing Vectors in Orthonormal Bases. Consider the vector \mathbf{z} shown in Figure 2. As is seen in the figure, \mathbf{z} does not lie completely within the two-dimensional plane (subspace) represented by the basis vectors $\bar{\mathbf{b}}_1$ and $\bar{\mathbf{b}}_2$. Rather, \mathbf{z} has two components, one that lies within the plane, \mathbf{z}^{\parallel} , and one, \mathbf{z}^{\perp} , that lies in the orthogonal subspace “perpendicular” to the plane, so that

$$\mathbf{z} = \mathbf{z}^{\parallel} + \mathbf{z}^{\perp} \quad (21)$$

and

$$\langle \mathbf{z}^{\parallel}, \mathbf{z}^{\perp} \rangle = 0 \quad (22)$$

From Figure 2 it follows that \mathbf{z}^{\parallel} is given by

$$\mathbf{z}^{\parallel} = z_1 \bar{\mathbf{b}}_1 + z_2 \bar{\mathbf{b}}_2 = (\bar{\mathbf{b}}_1 \ \bar{\mathbf{b}}_2) \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \quad (23)$$

The components of \mathbf{z} in the subspace spanned by $\bar{\mathbf{b}}_1$ and $\bar{\mathbf{b}}_2$ are obtained by computing the inner product of \mathbf{z} with each of the orthonormal basis vectors

$$\begin{aligned} \langle \bar{\mathbf{b}}_1, \mathbf{z} \rangle &= \langle \bar{\mathbf{b}}_1, \mathbf{z}^{\parallel} \rangle + \langle \bar{\mathbf{b}}_1, \mathbf{z}^{\perp} \rangle \\ \langle \bar{\mathbf{b}}_2, \mathbf{z} \rangle &= \langle \bar{\mathbf{b}}_2, \mathbf{z}^{\parallel} \rangle + \langle \bar{\mathbf{b}}_2, \mathbf{z}^{\perp} \rangle \end{aligned} \quad (24)$$

Substituting \mathbf{z}^{\parallel} from eq 23 and using the fact that $\langle \bar{\mathbf{b}}_1, \mathbf{z}^{\perp} \rangle = \langle \bar{\mathbf{b}}_2, \mathbf{z}^{\perp} \rangle = 0$ yields the desired result

$$\begin{aligned} \langle \bar{\mathbf{b}}_1, \mathbf{z} \rangle &= z_1 \overbrace{\langle \bar{\mathbf{b}}_1, \bar{\mathbf{b}}_1 \rangle}^1 + z_2 \overbrace{\langle \bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2 \rangle}^0 = z_1 \\ \langle \bar{\mathbf{b}}_2, \mathbf{z} \rangle &= z_1 \overbrace{\langle \bar{\mathbf{b}}_2, \bar{\mathbf{b}}_1 \rangle}^0 + z_2 \overbrace{\langle \bar{\mathbf{b}}_2, \bar{\mathbf{b}}_2 \rangle}^1 = z_2 \end{aligned} \quad (25)$$

As is seen from Figure 2, z_1 and z_2 are the *projections* of \mathbf{z} onto the two orthonormal basis vectors, $\bar{\mathbf{b}}_1$ and $\bar{\mathbf{b}}_2$. An important aspect of this approach, which is reminiscent of that typically used in multivariate linear regression methods,⁴ is the *squared error*, \mathcal{E}^2 , due to the lower dimensional approximation and is given by

$$\begin{aligned} \|\mathbf{z}\|^2 &= \|\mathbf{z}^{\parallel}\|^2 + \|\mathbf{z}^{\perp}\|^2 \\ &= (z_1^2 + z_2^2) + \mathcal{E}^2 \end{aligned} \quad (26)$$

The situation in nonorthogonal basis sets is considerably more complicated and is not discussed here.

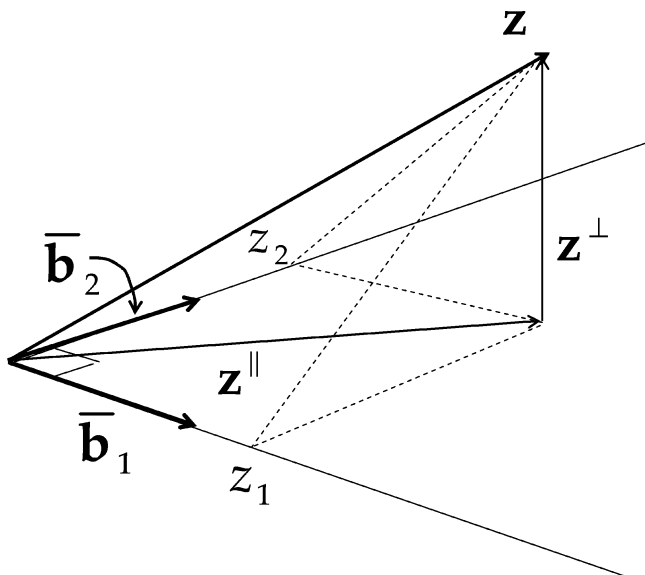


Figure 2. Representation of 3-D vector, \mathbf{z} , with respect to a 2-D orthonormal subspace $\{\mathbf{b}_1, \mathbf{b}_2\}$. The vector \mathbf{z}^{\parallel} is the projection of \mathbf{z} onto the 2-D subspace, while the vector \mathbf{z}^{\perp} lies in the orthogonal complement of $\{\mathbf{b}_1, \mathbf{b}_2\}$, and thus, $\langle \mathbf{z}^{\parallel}, \mathbf{z}^{\perp} \rangle = 0$.

2.4. Constructing a Molecular Basis Set. Consider a set \mathcal{B} of p molecules that constitute a *molecular basis* for a chemical space

$$\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_p\} \quad (27)$$

where each molecule can be thought of as an *abstract molecular basis vector* \mathbf{b}_i . The set of molecular basis vectors, selected such that they *approximately* span chemical space, are grouped together into a row matrix of basis vectors

$$\mathbf{b} = (\mathbf{b}_1 \ \mathbf{b}_2 \ \dots \ \mathbf{b}_p) \quad (28)$$

Forming the inner product of \mathbf{b} with itself yields the $p \times p$ dimensional matrix of inner products which is equivalent to

$$\langle \mathbf{b}, \mathbf{b} \rangle = \begin{pmatrix} \langle \mathbf{b}_1, \mathbf{b}_1 \rangle & \dots & \langle \mathbf{b}_1, \mathbf{b}_p \rangle \\ \vdots & \ddots & \vdots \\ \langle \mathbf{b}_p, \mathbf{b}_1 \rangle & \dots & \langle \mathbf{b}_p, \mathbf{b}_p \rangle \end{pmatrix}, \quad (29)$$

the metric matrix shown in eq 6 for the case of two basis vectors, that is, $\langle \mathbf{b}, \mathbf{b} \rangle \equiv \mathbf{G}$, where each element of \mathbf{G} is given by

$$G_{ij} = \langle \mathbf{b}_i, \mathbf{b}_j \rangle \quad \text{for } i, j = 1, 2, \dots, p \quad (30)$$

The following *ansatz* is used to evaluate the elements of metric matrix:

$$\langle \mathbf{b}_i, \mathbf{b}_j \rangle = G_{ij} \Leftrightarrow S_{ij} \quad (31)$$

where S_{ij} is an element of a molecular similarity matrix, thus

$$\begin{pmatrix} \langle \mathbf{b}_1, \mathbf{b}_1 \rangle & \dots & \langle \mathbf{b}_1, \mathbf{b}_p \rangle \\ \vdots & \ddots & \vdots \\ \langle \mathbf{b}_p, \mathbf{b}_1 \rangle & \dots & \langle \mathbf{b}_p, \mathbf{b}_p \rangle \end{pmatrix} \triangleq \begin{pmatrix} 1 & \dots & S_{1,p} \\ \vdots & \ddots & \vdots \\ S_{p,1} & \dots & 1 \end{pmatrix} = \mathbf{S}. \quad (32)$$

Each element of the similarity matrix satisfies the inequality $0 \leq S_{ij} \leq 1$.³³ The “1’s” along the diagonal correspond

to self-similarities, that is, the similarity of a molecule to itself.³⁴ As long as $S_{ij} < 1$ for all $i \neq j$, \mathbf{S} will be a *positive definite* matrix³⁵ and all of its eigenvalues will be positive, although numerical difficulties can arise as $S_{ij} \rightarrow 1$ if the level of precision used in the computations is not high enough. Importantly, $S_{ii} = \langle \mathbf{b}_i, \mathbf{b}_i \rangle = 1$, for $i = 1, 2, \dots, p$; thus, the molecular basis vectors are all of unit norm, $\|\mathbf{b}_i\| = \sqrt{\langle \mathbf{b}_i, \mathbf{b}_i \rangle} = 1$.

Since molecular similarity is usually taken to be symmetric,³⁶ $\mathbf{S} = \mathbf{S}^T$, and thus \mathbf{S} can be diagonalized by an orthogonal similarity transformation

$$\mathbf{V}^T \mathbf{S} \mathbf{V} = \Lambda, \quad (33)$$

where \mathbf{V} is the orthonormal matrix of eigenvectors, $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}$, and Λ is the diagonal matrix of eigenvalues, $\lambda_i > 0$, $i = 1, 2, \dots, p$, which follows since \mathbf{S} is a positive definite matrix unless there are dependencies within the set of molecular basis vectors.

As the similarity matrix is not, in general, diagonal and has no zero eigenvalues, the molecular basis set is linearly independent but not orthonormal. It is beneficial to orthonormalize the molecular basis set as it facilitates many of the operations performed in constructing a suitable abstract vector-based representation of chemical space.³⁷ There are many ways to construct an orthonormal set of vectors. One that is particularly appropriate in this work is *symmetric orthonormalization*.^{3,4}

$$\bar{\mathbf{b}} = \mathbf{b} \mathbf{S}^{-1/2} \quad (34)$$

The inner product obtained from eq 34 is given by

$$\begin{aligned} \langle \bar{\mathbf{b}}, \bar{\mathbf{b}} \rangle &= \langle \mathbf{S}^{-1/2} \mathbf{b}, \mathbf{S}^{-1/2} \mathbf{b} \rangle \\ &= (\mathbf{S}^{-1/2})^T \langle \mathbf{b}, \mathbf{b} \rangle \mathbf{S}^{-1/2} \\ &= \mathbf{S}^{-1/2} \mathbf{S} \mathbf{S}^{-1/2} \\ &= \mathbf{I} \end{aligned} \quad (35)$$

where $(\mathbf{S}^{-1/2})^T = \mathbf{S}^{-1/2}$, since $\mathbf{S}^{-1/2}$ is symmetric; \mathbf{I} is the identity matrix, which confirms that the transformation given in eq 34 does generate an orthonormal set of basis vectors, $\{\bar{\mathbf{b}}_1, \dots, \bar{\mathbf{b}}_i, \dots, \bar{\mathbf{b}}_p\}$. A desirable property of symmetric orthonormalization is that the basis set obtained is as close as possible, in a least-squares sense, to the original basis set.⁵ In addition, the similarity of the transformed, orthonormal molecular basis vectors to the “original” basis vectors depends inversely on the magnitude of the off-diagonal elements of the similarity matrix. This follows from the expansion of a symmetrically orthonormalized molecular basis vector in powers of the similarities of the original molecular basis set^{3,4}

$$\bar{\mathbf{b}}_i = \mathbf{b}_i - \frac{1}{2} \sum_k \mathbf{b}_k S_{k,i} + \frac{3}{8} \sum_k \sum_l \mathbf{b}_k S_{k,i} S_{l,k} + O(S^3) \quad (36)$$

From eq 36 it is clear that if the similarities are small, an orthonormal basis vector, $\bar{\mathbf{b}}_i$, is nearly identical to the original basis vector, \mathbf{b}_i . However, when the opposite holds, namely,

that the similarities tend to be large, an orthonormal basis vector contains linear combinations of basis vectors from the original nonorthogonal basis set with large similarities to the basis vector in question, as seen in eq 36. Thus, it is desirable in selecting an initial molecular basis set to ensure that the intermolecular similarities of its members be as small as possible so that the description of the orthonormalized molecular basis set derived from it is, in fact, close to it. This constraint also ensures that the chemical space of interest is reasonably well-covered by as small a set of molecular basis vectors as possible. A potentially undesirable feature of symmetric orthonormalization is that linear combinations of all of the original, nonorthogonal molecular basis vectors are effectively included to some degree in each of the orthonormal basis vectors produced by the procedure. Thus, unlike the example of Gram–Schmidt orthogonalization described in section 2.2, the addition of new vectors to a symmetrically orthonormalized basis can influence all of the previously orthonormalized basis vectors.

The $\mathbf{S}^{-1/2}$ matrix is computed as follows. Since \mathbf{S} is a positive definite matrix, its eigenvalues are all positive; thus, it is possible to transform the diagonal eigenvalue matrix Λ into a diagonal matrix whose elements are inverse square roots of the eigenvalues, namely, $\Lambda^{-1/2}$. Using the back transformation of eq 33 yields the desired result:

$$\mathbf{S}^{-1/2} = \mathbf{V}\Lambda^{-1/2}\mathbf{V}^T \quad (37)$$

Since $\mathbf{S}^{-1/2}$ is symmetric, but the inverse of $\mathbf{S}^{-1/2}$ is $\mathbf{S}^{1/2}$, it follows that $(\mathbf{S}^{-1/2})^T \neq (\mathbf{S}^{-1/2})^{-1}$, which clearly shows that $\mathbf{S}^{-1/2}$ is *not* an orthogonal transformation. A useful property of square-root matrices is that they can be handled in a similar fashion to exponential functions, namely, $\mathbf{S}^{-1/2} \mathbf{S}^{1/2} = \mathbf{S}^{-1/2} (\mathbf{S}^{-1/2})^{-1} = \mathbf{I}$.

2.5. Representing Molecules Using an Orthonormal Molecular Basis Set. In this subsection, mathematical expressions presented in the two-dimensional example described earlier are extended to cover the more general situation that constitutes the focus of this work, namely, the representation of molecules in a chemical space in terms of a general set of orthonormal molecular basis vectors. As was the case for molecular basis vectors shown in eq 28, the set of n molecules in a chemical space is represented by a row matrix of molecular vectors

$$\mathbf{m} = (\mathbf{m}_1 \quad \cdots \quad \mathbf{m}_r \quad \cdots \quad \mathbf{m}_n) \quad (38)$$

With respect to the orthonormal molecular basis set $\{\bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2, \dots, \bar{\mathbf{b}}_p\}$, a given molecular vector is represented by

$$\mathbf{m}_r = \sum_{i=1}^p \bar{\mathbf{b}}_i \langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle \quad r = 1, 2, \dots, n \quad (39)$$

which is a generalization to p dimensions of the example given in eqs 21–25. As noted earlier, the components $\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle$ of \mathbf{m}_r in the orthonormal basis are the coefficients of the generalized Fourier expansion of \mathbf{m}_r in the $\{\bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2, \dots, \bar{\mathbf{b}}_p\}$ basis. To evaluate the Fourier coefficients using the similarity *ansatz* given in eq 32, the orthonormal basis vectors in the inner product terms in eq 39 must be transformed back into the original nonorthogonal basis. To accomplish this, eq 39 is augmented to include the entire set of n molecular vectors given in eq 38

$$\mathbf{m} = \bar{\mathbf{b}} \langle \bar{\mathbf{b}}, \mathbf{m} \rangle \quad (40)$$

Substituting from eq 34 for $\bar{\mathbf{b}}$ and rearranging and regrouping terms yields

$$\begin{aligned} \mathbf{m} &= \bar{\mathbf{b}} (\mathbf{b} \mathbf{S}^{-1/2})^T \langle \mathbf{b}, \mathbf{m} \rangle \\ &= \bar{\mathbf{b}} (\mathbf{S}^{-1/2})^T \langle \mathbf{b}, \mathbf{m} \rangle \\ &= \bar{\mathbf{b}} (\mathbf{S}^{-1/2} \tilde{\mathbf{S}}) \end{aligned} \quad (41)$$

where the elements of the $p \times n$ dimensional matrix $\tilde{\mathbf{S}}$ correspond to the inner products of each of the molecular vectors, \mathbf{m}_r , with each of the basis vectors, \mathbf{b}_i , in the original nonorthogonal basis $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_p\}$. Thus, these terms can be determined in a completely analogous manner to the *ansatz* used to evaluate inner products among the original set of molecular basis vectors (see eq 31), namely, $\langle \mathbf{b}_k, \mathbf{m}_r \rangle \leftrightarrow \tilde{S}_{k,r}$. The term in parentheses in the last line of eq 41, $\mathbf{S}^{-1/2} \tilde{\mathbf{S}}$, is a $p \times n$ dimensional matrix whose elements are the generalized Fourier coefficients with respect to the orthonormal basis

$$\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle = \left(\mathbf{S}^{-1/2} \tilde{\mathbf{S}} \right)_{i,r} \quad (42)$$

Equation 42 clearly shows that the Fourier coefficients depend on the similarities of both the basis vectors to the vectors describing the molecules in chemical space ($\tilde{\mathbf{S}}$) and to the similarities of the basis vectors to each other (\mathbf{S}).

In analogy to the discussion given in section 2.3, molecular vectors can be decomposed into two terms:

$$\mathbf{m}_r = \mathbf{m}_r^{\parallel} + \mathbf{m}_r^{\perp} \quad (43)$$

where \mathbf{m}_r^{\parallel} is the component of \mathbf{m}_r that lies in the subspace spanned by the molecular basis set and \mathbf{m}_r^{\perp} is its *orthogonal complement*. Since $\langle \mathbf{m}_r^{\parallel}, \mathbf{m}_r^{\perp} \rangle = 0$, and $\|\mathbf{m}_r\| = 1$, as all molecular vectors are similar to themselves ($\|\mathbf{m}_r\| = S_{r,r} = 1$)

$$\|\mathbf{m}_r\|^2 = \|\mathbf{m}_r^{\parallel}\|^2 + \|\mathbf{m}_r^{\perp}\|^2 = 1 \quad (44)$$

It is worth noting that this is analogous to the procedure underlying least-squared-error-based linear regression methods.⁴ Importantly, $\|\mathbf{m}_r^{\perp}\|^2$ is the *squared error* associated with the incompleteness of the molecular basis set, that is, $\|\mathbf{m}_r^{\perp}\|^2 = \mathcal{E}_r^2$. Thus, eq 44 can be rearranged to

$$\begin{aligned} \mathcal{E}_r^2 &= 1 - \|\mathbf{m}_r^{\parallel}\|^2 \\ &= 1 - \sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2 \end{aligned} \quad (45)$$

It can be shown that

$$\sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2 = (\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}})_{r,r} \quad (46)$$

so that eq 45 can also be written as

$$\mathcal{E}_r^2 = 1 - (\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}})_{r,r} \quad (47)$$

which will be of use in section 4.1 for describing the average squared error for a set of molecules with respect to a given chemical space.

Equations 45 and 47 show that all molecules in the chemical space defined by the orthonormal basis set $\{\bar{\mathbf{b}}_1, \bar{\mathbf{b}}_2, \dots, \bar{\mathbf{b}}_p\}$ lie within a p -dimensional solid hypersphere of unit radius³⁸ unless $\mathcal{E}_r^2 = 0$. In that case, a highly unlikely occurrence, the molecular vectors will lie on the p -dimensional hypersphere that bounds it. Thus, arguments concerning the *relative volume* of a solid hypersphere inscribed within a hypercube of similar dimension^{39,40} are not entirely applicable. This follows since all the “molecular points” lay within the hypersphere, and thus the remaining volume of the enveloping hypercube is irrelevant to this chemical space.⁴¹

Combining eq 46 with eq 45 exhibits another useful feature of the current approach, namely,

$$\sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2 = 1 - \mathcal{E}_r^2 \quad r = 1, 2, \dots, n \quad (48)$$

Each of the squared generalized Fourier coefficients $|\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2$ in the summation on the left-hand side of the equation corresponds to the fraction of the molecular vector \mathbf{m}_r that is represented by each of the orthonormal molecular basis vectors $\bar{\mathbf{b}}_k$, $k = 1, 2, \dots, p$ (vide supra). The squared error term is equal to the fraction of \mathbf{m}_r not represented by the orthonormal basis set.

2.6. Computational Methodology. All computations reported in this work were carried out using Matlab version 7.0 and the Matlab statistical analysis module.⁴² Similarities were evaluated using the Tanimoto similarity coefficient^{2,10} based on MACCS key fingerprints, which were obtained from the program MOE⁴³ available from Chemical Computing Group (Montreal, Canada). Molecules were represented using the sd file format originally developed by MDL Information Systems.¹¹

3. MOLECULAR DATA SETS

3.1. Molecular Basis and Test Sets. A set of molecules was obtained from the NCI AIDS database to illustrate the feasibility of the approach proposed in this work. The database contains about 43 000 diverse compounds tested as possible inhibitors of the AIDS virus. Importantly, it is publicly available and can be downloaded from the NCI Web site.²⁸ In this study, a set of 6000 compounds was selected from the first 6000 compounds in the database and divided into two subsets. The first subset, \mathcal{B}_{2000} , which consists of the first 2000 compounds in the original set, is used to construct a family of eight molecular basis sets of increasing size taken sequentially in increments of 250 compounds: $\mathcal{F}_{\mathcal{B}} = \{\mathcal{B}_{250}, \mathcal{B}_{500}, \mathcal{B}_{750}, \mathcal{B}_{1000}, \mathcal{B}_{1250}, \mathcal{B}_{1500}, \mathcal{B}_{1750}, \mathcal{B}_{2000}\}$. The second subset, \mathcal{T}_{4000} , consists of the remaining 4000 compounds in the original set. This subset is used as a test set to evaluate the effect of basis-set size on \mathcal{E}_r^2 .

3.2. Characteristics of the Molecular Sets. It is desirable that the similarities among the molecules constituting a molecular basis set be as small as possible to minimize collinearities in the set of molecular basis vectors (vide supra). This condition is tantamount to ensuring that the

Table 1. Values of Average Similarities and Their Corresponding Standard Deviations Computed for the Test Set of 4000 Molecules with Respect to Molecular Basis Sets of Increasing Size from 250 to 2000 Molecular Basis Vectors

	molecular sets	average similarity	standard deviation
basis sets	1–250	0.2109	0.1246
	1–500	0.2255	0.1304
	1–750	0.2325	0.1254
	1–1000	0.2355	0.1248
	1–1250	0.2371	0.1238
	1–1500	0.2377	0.1230
	1–1750	0.2414	0.1245
	1–2000	0.2423	0.1227
test set	2000–6000	0.2679	0.1215

molecular basis is made up of a diverse set of molecules. Table 1 provides a statistical summary for the eight basis sets and the test set. Not unexpectedly, as can be seen from the data in the table, the sample averages of the eight basis sets, \mathcal{B}_{250} , \mathcal{B}_{500} , \mathcal{B}_{750} , \mathcal{B}_{1000} , \mathcal{B}_{1250} , \mathcal{B}_{1500} , \mathcal{B}_{1750} , and \mathcal{B}_{2000} , increase monotonically with the size of the basis set from 0.2109 to 0.2423. The corresponding standard deviations remain relatively constant around 0.1250 and, thus, do not show the same, clearly monotonic, behavior with increasing basis set size as do the average similarities. As seen in Table 1, the test set values are also comparable to those of the basis sets. The averages and standard deviations for all of the molecular sets in Table 1 clearly show that most of the molecules considered in the basis sets are quite dissimilar as are those considered in the test set. The key issue regarding the molecular basis sets is whether or not they are sufficiently complete to adequately represent the chemical space of interest represented in this work by \mathcal{T}_{4000} . It is here that the value of the error term will clearly be seen (vide infra).

4. RESULTS AND DISCUSSION

4.1. Relationship of Average Squared Error to Molecular Basis Set Size. As noted in the previous section, eight nested molecular basis sets of increasing size were selected from the first 2000 molecules in the NCI AIDS database to test the proposed methodology. The average squared error for a basis set of size p , $\langle \mathcal{E}^2 \rangle_p$, can be computed for each molecule in the test set, \mathcal{T}_{4000} , as

$$\begin{aligned} \langle \mathcal{E}^2 \rangle_p &= \frac{1}{n} \sum_{r=1}^n \mathcal{E}_r^2 \\ &= \frac{1}{n} \sum_{r=1}^n \left(1 - \sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2 \right) \\ &= \frac{n}{n} - \frac{1}{n} \sum_{r=1}^n (\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}})_{r,r} = 1 - \text{Tr}(\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}}) \end{aligned} \quad (49)$$

where n is the size of the test set, which is 4000 in the present case, “Tr” is the trace of the matrix $\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}}$, and \mathcal{E}_r^2 can be computed using either eq 45 or eq 47.

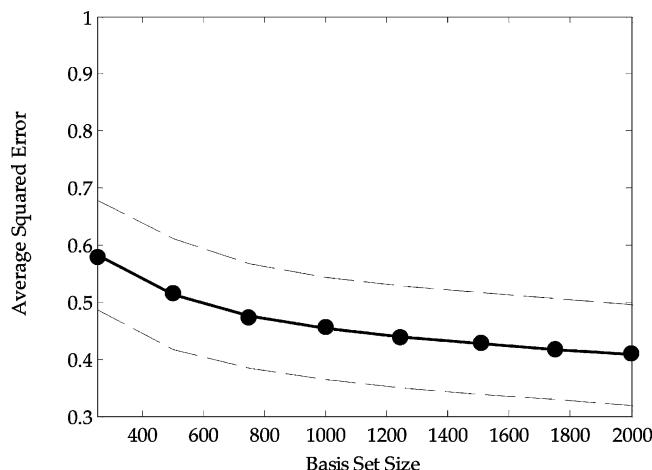


Figure 3. Plot of the average squared basis set error (solid, heavy line) bounded from above and from below by the corresponding standard deviation (dashed line) calculated for the test set of 4000 molecules with respect to the eight basis sets.

The results are plotted in Figure 3 for basis sets of increasing size. From the figure, it is clear that $\langle \mathcal{G}^2 \rangle_p$ decreases monotonically with increased basis-set size, while the associated standard deviations (indicated by the light, dashed lines bounding the squared error) remain relatively constant. This is not surprising since the capability of a basis set to represent a chemical space should improve as its size increases unless linear dependencies arise. Such an occurrence is, however, not likely to occur in the method proposed here since linear dependencies are detected as zero, or approximately zero, eigenvalues of the metric matrix and can be removed.

Alternatively, the average norm of the set of 4000 test-set molecular vectors in the subspaces spanned by molecular basis sets of increasing size is given by

$$\begin{aligned} \langle \|\mathbf{m}^{\parallel}\| \rangle_p &= \frac{1}{n} \sum_{r=1}^n \sqrt{\sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2} \\ &= \frac{1}{n} \sqrt{\text{Tr}(\tilde{\mathbf{S}}^T \mathbf{S}^{-1} \tilde{\mathbf{S}})} \end{aligned} \quad (50)$$

where \mathbf{m}^{\parallel} is the projection of \mathbf{m} onto the given basis and $\|\mathbf{m}^{\parallel}\|$ is its norm. Figure 4 shows that the average norm increases with basis-set size, asymptotically approaching unity as the basis-set size approaches infinity. The standard deviation indicated by the light, dashed lines bounding the norm decreases slightly with an increase in basis-set size. When the basis set is complete, molecular vectors can be represented exactly. In this case, which cannot be realized in practice, all of the molecular vectors will lie on the p -dimensional hypersphere.⁴⁴ In that case, the norm of all of the molecular vectors is unity, as is the average norm, and the standard deviation is zero.

Table 2 provides a summary of the average of the squared generalized Fourier coefficients for each of the basis sets of size p . The average is computed as

$$\langle \text{GFC}^2 \rangle_p = \frac{1}{pn} \sum_{r=1}^n \sum_{i=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2 \quad (51)$$

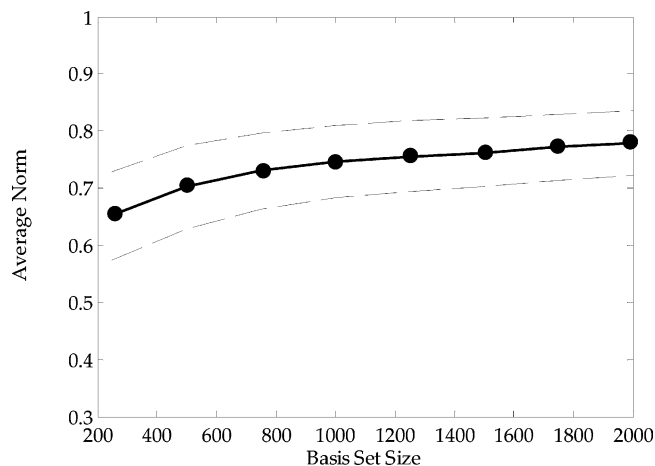


Figure 4. Plot of average norm (solid, heavy line) bounded from above and from below by the corresponding standard deviation (dashed line) calculated for the test set of 4000 molecules with respect to the eight basis sets.

Table 2. Values of Average Squared Generalized Fourier Coefficients, $\langle \text{GFC}^2 \rangle$, Computed for the Test Set of 4000 Molecules with Respect to Molecular Basis Sets of Increasing Size from 250 to 2000 Molecular Basis Vectors

basis set	$\langle \text{GFC}^2 \rangle$
1–250	1.6816×10^{-3}
1–500	9.7803×10^{-4}
1–750	7.0198×10^{-4}
1–1000	5.4850×10^{-4}
1–1250	4.5135×10^{-4}
1–1500	3.8349×10^{-4}
1–1750	3.3467×10^{-4}
1–2000	2.9757×10^{-4}

where GFC^2 stands for the square of the generalized Fourier coefficient and $\sum_{k=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2$ is the square of the magnitude of the projection of \mathbf{m}_r (i.e., \mathbf{m}_r^{\parallel}) onto the appropriate molecular basis.⁴⁵ As shown in Figure 3, $\langle \mathcal{G}^2 \rangle_p$ decreases modestly with basis set size, and thus, $\sum_{k=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2$ must, correspondingly, only increase modestly in magnitude as well. However, since the range of the summation increases rapidly as the basis-set size increases, the magnitude of each $|\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2$ must on average decrease in order to maintain the near constancy of $\sum_{k=1}^p |\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle|^2$. Thus, as shown in Table 2, the average value of the squared generalized Fourier coefficients should become smaller as the basis set size increases.

From Figure 3, it is clear that the size of the basis set required to reduce the error to, say, 10% will be substantial. This raises the question as to how small the error or, alternatively, how large the basis set needs to be to represent the salient characteristics of chemical space required in a given type of analysis. This question is considered in the following sections.

4.2. Exploring the Characteristics of the High-Dimensional Chemical Spaces. Even though the dimensions of some of the chemical spaces examined in this work are quite high, it may be the case that lower-dimensional representations in the range of 250–500 or fewer basis vectors will be sufficient if certain features of the chemical space are largely preserved. To address this issue, four statistics are considered: (1) the *average Euclidean distance* between two molecules in a chemical space of specified dimension, (2)

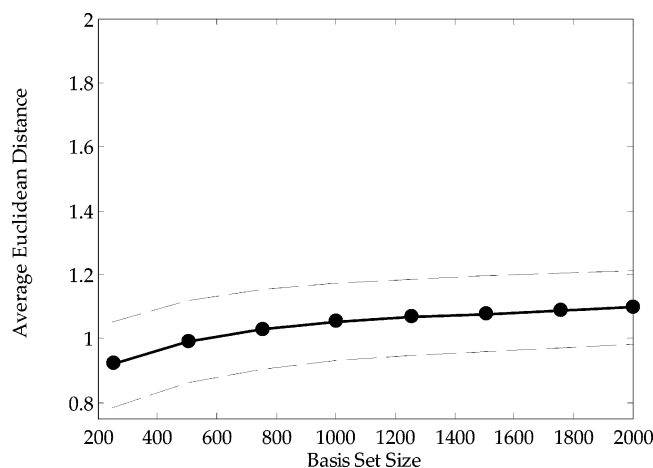


Figure 5. Plot of average Euclidean distance (solid, heavy line) bounded from above and from below by the corresponding standard deviation (dashed line) calculated for the test set of 4000 molecules with respect to the eight basis sets.

the *average difference in Euclidean distance* between two molecules with respect to two spaces of different dimension, (3) the *Spearman rank correlation*⁴⁶ between two spaces of different dimension, and (4) a variant of *Kruskal's stress function*^{2,47} between two spaces of different dimension.

All four statistics make use of the *Euclidean distance* between two molecules, \mathbf{m}_r and \mathbf{m}_s , which is given by

$$d(\mathbf{m}_r, \mathbf{m}_s)_p = \sqrt{\sum_{i=1}^p (\langle \bar{\mathbf{b}}_i, \mathbf{m}_r \rangle - \langle \bar{\mathbf{b}}_i, \mathbf{m}_s \rangle)^2} \quad (52)$$

where p is the dimension of the chemical space. Figure 5 shows a plot of the average Euclidean distance between a pair of molecules in the test set as a function of the basis-set size

$$\langle d \rangle_p = \frac{1}{n(n-1)/2} \sum_{r=1}^n \sum_{s>r}^n d(\mathbf{m}_r, \mathbf{m}_s)_p \quad (53)$$

The plot indicates that the average Euclidean distance slowly increases with the basis-set size and that the corresponding standard deviations increase in magnitude.

Because of the transformation given in eq 42, the generalized Fourier coefficients of molecular vectors in the orthonormalized basis set can take on negative as well as positive values even though the coefficients in the original nonorthogonal basis, $\langle \mathbf{b}_i, \mathbf{m}_r \rangle = (\tilde{\mathbf{S}})_{i,r}$ are all, by construction, positive and bounded by unity. Thus, molecular vectors in the orthonormalized chemical space, in principle, can be distributed over the entire chemical space, which lies within a unit p -dimensional solid hypersphere. In such chemical spaces, a strict upper bound to the distance can be deduced as illustrated in Figure 6, which shows that the *maximum* distance between any two molecular vectors represented in a chemical space is bounded from above by 2 and occurs when two molecular vectors are diametrically opposed to one another. Vectors that lie on the dashed circle in Figure 6 represent vectors whose norms approach but are less than the maximum possible value of unity, which only obtains in the case of a complete basis set $p = \infty$.⁴⁸ As illustrated in the figure for the two-dimensional case, the maximum

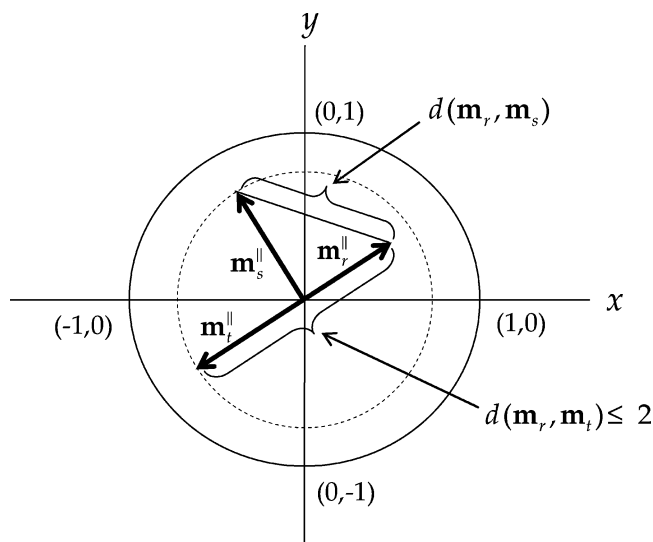


Figure 6. Depiction of Euclidean distance relationships among molecular vectors in a 2-D subspace.

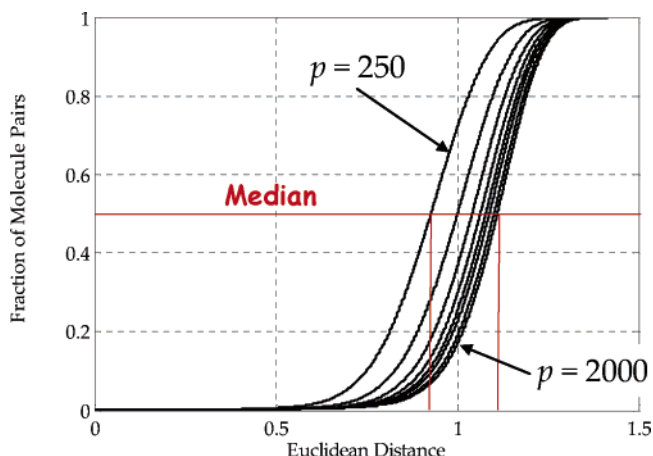


Figure 7. Cumulative distribution function plots of the pairwise Euclidean distance distributions calculated for a random sample of 1000 test set molecules with respect to the eight basis sets. The left most cdf is associated with \mathcal{B}_{250} , the next with \mathcal{B}_{500} , ..., and the right most with \mathcal{B}_{2000} . The median is indicated in red.

distance between two vectors, $\mathbf{m}_s^||$ and $\mathbf{m}_t^||$, occurs when, as noted above, they are collinear but diametrically opposed. In such cases, $d(\mathbf{m}_s, \mathbf{m}_t) = \|\mathbf{m}_s^||\| + \|\mathbf{m}_t^||\|$; as $p \rightarrow \infty$, \mathcal{E}_s^2 and $\mathcal{E}_t^2 \rightarrow 0$ and $\|\mathbf{m}_s^||\|$ and $\|\mathbf{m}_t^||\| \rightarrow 1$, and thus $d(\mathbf{m}_s, \mathbf{m}_t) \rightarrow 2$.⁴⁹ Although this situation cannot be realized practically, it can serve as a rigorous upper bound to the distances between molecular vectors in a given chemical space. Figure 3 shows that a considerable extension in the size of the molecular basis set is needed before that upper bound is closely approached. Moreover, even in the case of zero error, where the norm of all of the molecular vectors is unity, the average distance will lie below the upper bound of 2 in nearly all nontrivial cases.⁵⁰ This follows since most pairs of molecular vectors in a chemical space are not diametrically opposed, and thus the maximum can only be obtained for a small fraction of all possible pairs of molecules.

An alternative, more detailed, view of the distance data is presented in Figure 7 in the form of cumulative distribution function (cdf) plots,⁴⁷ one for each of the eight basis sets. The figure shows that the medians of the cdf's move toward

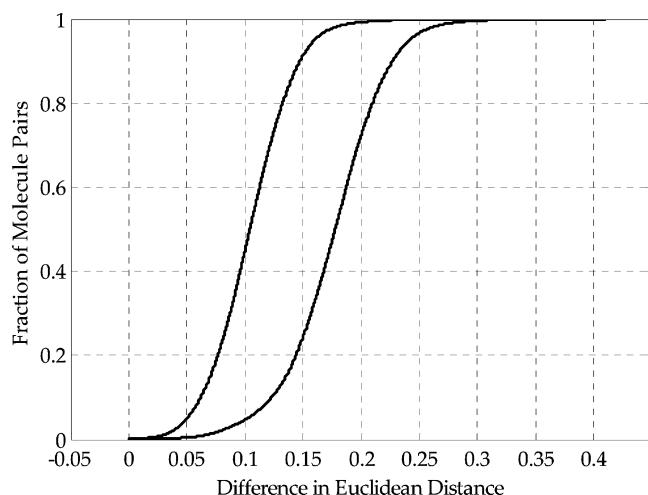


Figure 8. Cumulative distribution function plots of the pairwise differences in Euclidean distance for the \mathcal{B}_{250} and \mathcal{B}_{500} basis sets compared to the “complete” \mathcal{B}_{2000} basis set, calculated for a random sample of 1000 test-set molecules. The left most cdf corresponds to \mathcal{B}_{500} and the right most to \mathcal{B}_{250} .

longer distances as the basis-set size increases. This is entirely expected behavior, as the distance between any two objects in a space increases monotonically as the dimension increases, and thus, both the average and median should also increase.

Figure 8 depicts cdf plots of the *difference* in Euclidean distances for pairs of test set molecules, \mathbf{m}_r and \mathbf{m}_s

$$\Delta(\mathbf{m}_r, \mathbf{m}_s)_{p,q} = d(\mathbf{m}_r, \mathbf{m}_s)_q - d(\mathbf{m}_r, \mathbf{m}_s)_p \quad \begin{cases} q = 2000 \\ p = 250, 500 \\ r < s = 1, \dots, 1000 \end{cases} \quad (54)$$

where q corresponds to the 2000-dimensional “reference chemical space”, which is compared to two lower p -dimensional spaces. The plots are generated from a single random sample of 1000 molecules taken from the test set of 4000 molecules.

The plot clearly shows that differences are more pronounced in the case of the $p = 250$ basis set than the $p = 500$ basis set, since the medians of the cdf’s differ by approximately 0.075. The general trend in the differences is expected since as the basis-set size increases the squared error decreases, and thus differences in distances between the reference basis set and the smaller basis sets should get smaller.

Spearman rank correlation⁴⁶ assesses the correlation between the rank orderings, from largest to smallest, of the $n(n-1)/2$ distances computed for a set of n molecules with respect to chemical spaces of different dimensions and is given by

$$r_{\text{rank}} = 1 - \frac{6 \sum_{i=1}^N \rho_i^2}{N(N^2 - 1)} \quad (55)$$

where $N = n(n-1)/2$ and ρ_i^2 is the square of the *difference in ranking* of the distances between pairs of molecules computed in each of the respective spaces. Thus, it provides a comparative means for assessing interchanges in the order

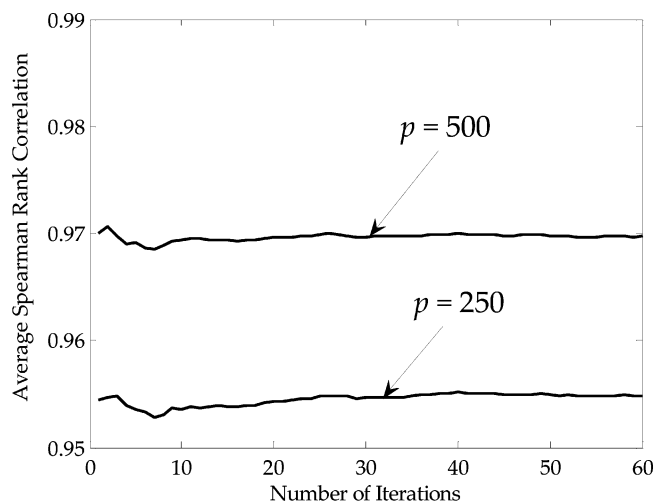


Figure 9. Plot of the average Spearman rank correlation of Euclidean distances in the \mathcal{B}_{250} and \mathcal{B}_{500} basis sets compared to the “complete” \mathcal{B}_{2000} basis set, calculated for multiple random samples of 1000 test-set molecules.

of distances between molecules in a set of molecules computed with respect to chemical spaces of different dimension. The *average* Spearman rank correlation was determined using eq 55 by averaging the average rank correlations on samples of 1000 molecules taken from the test set until convergence was reached. Figure 9 depicts the results for basis-sets of $p = 250$ and $p = 500$ dimensions, where the rank correlations were computed for these basis sets against the reference basis set where $p = 2000$. From the figure, it is clear that the sampling process is well-converged in both cases after about 40 iterations. As expected, rank correlation is improved in the larger basis set ($p = 500$) compared to the smaller basis set ($p = 250$), but the rank correlation is quite good in both cases.

Kruskal’s stress function,^{2,47} which is given by

$$\mathcal{K}_{p,q} = \frac{\sum_{r=1}^n \sum_{s>r}^n [d(r,s)_p - d(r,s)_q]^2}{\sum_{r=1}^n \sum_{s>r}^n d(r,s)_q^2} \quad (56)$$

provides a comparative measure of the aggregate effect of basis-set size on distances calculated in two chemical spaces of dimensions p and q , respectively. The term in the denominator defines the q -dimensional “reference chemical space” that in this work is 2000 while $p = 250$ and 500. Stress is approximately zero in the case where $d(r,s)_p \approx d(r,s)_q$ and approaches unity as $d(r,s)_q \gg d(r,s)_p$, which occurs as the p -dimensional “test chemical space” decreases in dimension, that is, as $p \rightarrow 0$. A sampling procedure identical to that used to determine the average Spearman rank correlation is also used to determine the average Kruskal stress function. The results are depicted in Figure 10, which clearly shows the very rapid convergence of the sampling procedure. As expected, the stress was considerably less for the $p = 500$ basis set than for the $p = 250$ basis set.

4.3. Are These Chemical Spaces of Use Practically?

Essentially all of the results given in the previous two subsections support the notion that the high-dimensional chemical spaces described in the current approach behave

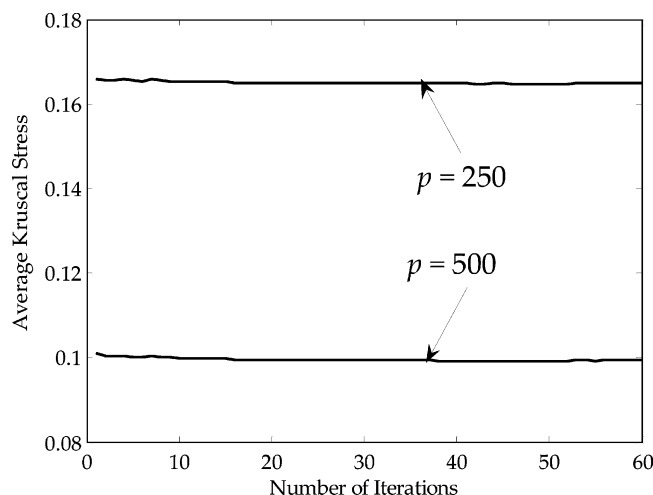


Figure 10. Plot of the average Kruskal stress in the \mathcal{B}_{250} and \mathcal{B}_{500} basis sets compared to the “complete” \mathcal{B}_{2000} basis set, calculated for multiple random samples of 1000 test-set molecules.

in a reasonable and consistent manner. Many of the seemingly paradoxical behaviors associated with such high-dimensional chemical spaces do not provide any substantive impediments in this work. For example, consider the well-known hypersphere—hypercube paradox:³⁹ as dimension increases, the ratio of the volume of a unit solid hypersphere embedded within a unit hypercube and tangent to its faces goes to zero, implying that the hypersphere becomes a multidimensional δ function with a vanishingly small volume. Since all of the molecular vectors lie within the unit solid hypersphere,⁵¹ the result seems to imply that all of the intermolecule distances in these spaces will tend toward zero. This is definitely not the case as shown by the various measures of intermolecule distances described in the previous subsection, especially the distance averages and standard deviations depicted in Figure 5 and the cdf's describing the distance distributions depicted in Figure 7. Moreover, as seen in these two figures, the fact that the averages and medians of the distance distributions are close suggests, but certainly does not prove, that the distributions are approximately Gaussian. Interestingly, as described by Agrafiotis,⁴⁷ such is the case in diverse sets of molecules represented in high-dimensional chemical spaces.

Another example is provided by the ratio of the volume of the *outer hyperspherical shell* of a unit solid hypersphere of given dimension with a thickness equal to say 10% of its radius.^{39,52,53} As dimension increases the ratio goes to unity, indicating that all of the volume effectively resides in the outer hyperspherical shell. This fact is also entirely consistent with the behavior of the high-dimensional chemical spaces described in the previous subsection. As shown in Figure 4, the average norm of a set of molecular vectors increases while its corresponding standard deviation decreases, albeit only modestly. This indicates that as basis-set size, and thus the accuracy of the representation, increases, the molecular vectors are on average moving toward the surface of the hypersphere and will lie within the outer hyperspherical shell.⁵⁴

Results for both average Spearman rank correlation and average Kruskal stress shown in Figures 9 and 10 provide further support of consistent behavior of the chemical spaces, even for the smaller basis sets considered here, namely, $p =$

250 and $p = 500$. In the former case, the average rank correlation is greater than 0.95 for both basis sets compared to the reference basis set, where $p = 2000$. The average Kruskal stress also lies below 0.20 for both basis sets, again in comparison to the largest ($p = 2000$) basis set. Both of these statistics clearly show that the lower dimensional basis sets are still capable of capturing many of the distance-related properties of the largest basis set considered in this work. Since distance is an important concept in metric spaces, and since the distance properties of the chemical spaces considered here appear to be reasonably consistent over the range of basis-set sizes examined in this work, it is likely that the lower dimensional basis set ($p = 250$) may be sufficient for most chemical space applications. In fact, chemical spaces of even lower dimension may also be applicable, but this is a subject for future investigation.

5. SUMMARY AND CONCLUSIONS

From the above discussion, it is clear that the high-dimensional spaces described in this work do not present insurmountable difficulties that would render the spaces incapable of providing effective representations of chemical space. The power of the approach is 3-fold. First, it utilizes a well-known mathematical method, namely, generalized Fourier series, for representing chemical spaces in terms of a set of “molecular basis vectors” rather than by typical molecular descriptors many of which lack chemically intuitive appeal.⁵⁵ Because the methodology is based upon generalized Fourier series, all of the mathematical apparatus that underlies that method can be used, an important aspect of which is that it provides a clear description of the convergence of the molecular vectors, in terms of their squared error, as the size of the basis set is increased. Second, it affords a means for directly evaluating the inner products of molecular vectors without the need for an explicit representation of the molecular vectors themselves. This is accomplished using an *ansatz* that equates inner products between these vectors with the similarities of the corresponding molecules. This approach is well-documented in many kernel-based, machine-learning applications that describe the direct construction of inner products without the need first to construct a related feature space.^{7–9} Thus, an added advantage of the approach is that it can be applied to *any* problem domain where inner products can be directly evaluated.⁸ Because of this, it is possible to construct orthonormal basis sets, such as those described in this work, which can facilitate the representation of the relevant spaces significantly. Third, although the dimension of the chemical spaces described in this work can be quite large, their metric behavior is reasonable on the basis of the results of a number of statistics discussed in section 4.

An issue that has not been considered in the current work is how to construct molecular basis sets optimally. Although this can be done in a number of ways, the choice may also depend on other factors such as, for example, whether a very general chemical-space representation is desired or whether a more restricted representation “tuned” to specific classes of compounds is preferred.⁵⁶ In any case, the basis sets used in the current work are sufficiently diverse, as measured by the average similarities and standard deviations reported in Table 1, to provide a reasonable starting point for the study. A closely related issue is what size basis set is needed to

represent the chemical space of interest effectively. In addition, as noted earlier, the choice of similarity method can have an impact. For example, if stereochemical information is required, a 3-D similarity method may be needed. These issues are being addressed as part of our ongoing research and will be dealt with in future publications.

One disadvantage of the high-dimensionality of the spaces constructed by the current method is that the position of objects (molecules) in these spaces cannot be easily visualized. An initial attempt based on principal components analysis was unsuccessful. In that study, it was found that, as the size of the basis set increased, so also did the number of principal components needed to describe a constant percent of the sample variance, say 75%. This can be understood from the discussion on the squares of the generalized Fourier coefficients presented in section 4.1, especially the discussion surrounding eq 51. The argument is based on the fact that the squared error decreases relatively slowly with respect to increases in basis-set size. Hence, many generalized Fourier coefficients, mostly of approximately equal but small magnitude, are needed to adequately represent chemical spaces as their dimension increases. This implies that the dimensions associated with these coefficients are relevant to an overall description of the inherent variance in the data. Moreover, it implies that the data are approximately spherically distributed in the full space, and thus, the principle components will be approximately degenerate. In such cases, the variances along each of the principal components are approximately equal, but small, in magnitude.⁵⁷ Thus, an increasing number of principal components are needed as the basis-set size increases to account for a fixed amount of variance in the sample.

Fortunately, a number of other methods exist for embedding high-dimensional data into lower dimensional spaces. The methods include nonlinear mapping,^{58,59} multidimensional scaling,⁶⁰ isometric mapping,⁶¹ local linear embedding,⁶² exploratory project pursuit,⁶³ stochastic proximity embedding,^{64,65} and eigenvalue-based methods,⁶⁶ to name a few, and we are currently investigating their feasibility. However, even if these methods are unable to produce reasonable lower-dimensional representations that are suitable for visualization, it does not preclude the use of the chemical spaces described here from being of value in many areas of cheminformatics.

ACKNOWLEDGMENT

The authors would like to thank Dr. Jose Medina-Franco for help in manuscript preparation and the BIO5 Institute of the University of Arizona for partial support of this work.

REFERENCES AND NOTES

- Scior, T.; Bernard, P.; Medina-Franco, J. L.; Maggiora, G. M. Large Compound Databases for Structure–Activity Relationships Studies in Drug Discovery. *Mini-Rev. Med. Chem.* In press.
- Maggiora, G. M.; Shanmugasundaram, V. Molecular Similarity Measures. In *Cheminformatics: Concepts, Methods, and Tools for Drug Discovery*; Bajorath, J., Ed.; Humana Press: Totowa, New Jersey, 2004; pp 1–50.
- Löwdin, P. O. *Linear Algebra for Quantum Theory*; John Wiley & Sons: New York, 1998.
- Löwdin, P. O. On Linear Algebra, the Least Square Method, and the Search for Linear Relations by Regression Analysis in Quantum Chemistry and Other Sciences. *Adv. Quantum Chem.* **1992**, *23*, 83–126.
- Carlson, B. C.; Keller, J. M. Orthogonalization Procedures and the Localization of Wannier Functions. *Phys. Rev.* **1957**, *105*, 102–103.
- In kernel-based machine-learning applications, this matrix is usually called the Gram matrix.^{7–9}
- Shawn-Taylor, J.; Cristianini, N. *Kernel Methods for Pattern Analysis*; Cambridge University Press: Cambridge, U. K., 2004.
- Herbrich, R. *Learning Kernel Classifiers*; The MIT Press: Cambridge, MA, 2002.
- Schölkopf, B.; Smola, A. J. *Learning with Kernels*; The MIT Press: Cambridge, MA, 2002.
- Willett, P.; Barnard, J. P.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.
- MACCS Structural Keys; MDL Information Systems Inc.: San Leandro, CA. <http://www.mdli.com> (accessed Mar 21, 2007).
- Weininger, D. SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules. *J. Chem. Inf. Comput. Sci.* **1988**, *28*, 31–36.
- Leach, A. R.; Gillet, V. J. *An Introduction to Chemoinformatics*; Kluwer Academic Publishers: Dordrecht, The Netherlands, 2003.
- Mestres, J.; Rohrer, D. C.; Maggiora, G. M. MIMIC: A Molecular-Field Matching Program. Exploiting Applicability of Molecular Similarity Approaches. *J. Comput. Chem.* **1997**, *18*, 934–954.
- Lemmen, C.; Lengauer, T.; Klebe, G. FLEXS: A Method for Fast Flexible Ligand Superposition. *J. Med. Chem.* **1998**, *41*, 4502–4520.
- OpenEye Scientific Software, Santa Fe, NM. <http://www.eyesopen.com> (accessed Mar 21, 2007).
- Lemmen, C.; Lengauer, T. Computational Methods for the Structural Alignment of Molecules. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 215–232.
- Good, A. C.; Richards, W. G. Explicit Calculation of 3D Molecular Similarity. *Perspect. Drug Discovery* **1998**, *9/10/11*, 321–338.
- This is identical to the use of a basis of orthogonal polynomials in least squares approximations of functions. See, for example: Johnson, L. W.; Riess, R. D. *Numerical Analysis*; Addison-Wesley Publishing Company: Reading, MA, 1982.
- Note that this error will tend to be much smaller in the case that one of the molecular vectors also happens to be equal or very close to one of the original *non-orthogonal* basis vectors. However, because the original basis set has been normalized, the value of the squared Fourier coefficient can deviate significantly from unity.
- Jolliffe, I. T. *Principal Component Analysis*, 2nd ed.; Springer: New York, 2002.
- Gower, J. C. Some Distance Properties of Latent Roots and Vector Methods Used in Multivariate Analysis. *Biometrika* **1966**, *53*, 325–338.
- Klein, C.; Kaiser, D.; Kopp, S.; Chiba, P.; Ecker, G. F. Similarity Based SAR (SIBAR) as [a] Tool for Early ADME Profiling. *J. Comput.-Aided Mol. Des.* **2003**, *16*, 785–793.
- Oprea, T. I.; Gottfries, J. Chemography: The Art of Navigating in Chemical Space. **2001**, *3*, 157–166.
- cf. the discussion in section 2.4, in particular refs 33 and 35.
- Randic, M. Resolution of Ambiguities in Structure–Property Studies by Use of Orthogonalized Descriptors. *J. Chem. Inf. Comput. Sci.* **1991**, *31*, 311–320.
- Typical scaling procedures include transformation of a set of values to zero mean and unit variance, $z_i = (x_i - x_{\text{ave}})/\sigma$, and transformation of a set of values to lie within the unit interval, $z_i = (x_i - x_{\text{min}})/(x_{\text{max}} - x_{\text{min}})$.
- National Institutes of Health. Screening Services. http://dtp.nci.nih.gov/docs/aids/aids_screen.html and http://dtp.nci.nih.gov/docs/aids/aids_data.html (accessed Mar 21, 2007).
- Also called the Gram matrix in the machine learning literature.^{7–9}
- The variance–covariance matrix is an appropriate metric matrix for many statistical applications.⁴
- Lay, D. C. *Linear Algebra and Its Applications*, 2nd ed.; Addison-Wesley: Reading, MA, 1997.
- See section 2. 4 for further discussion.
- It is important to note that *electrostatic similarity* does not necessarily satisfy the inequality $0 \leq S_{ij} \leq 1$ as discussed in: Maggiora, G. M.; Petke, J. D.; Mestres. A General Analysis of Field-Based Molecular Similarity Indices. *J. Math. Chem.* **2002**, *31*, 251–270.
- It should be noted that some similarities involve 2-D or 3-D structure matches.² In such cases, multiple solutions can arise, and the optimum solution, corresponding to the maximum similarity, should generally be chosen.
- Since electrostatic similarity can have negative values,³³ it will not lead to a positive definite similarity matrix and, therefore, cannot be used in the current methodology. This, however, is not a serious limitation as electrostatic similarities are usually combined with some type of shape similarity such that the overall similarity generally takes on a positive value.

- (36) Note that asymmetric similarities are sometimes used in specific situations. See, for example: Tversky, A. Features of Similarity. *Psychol. Rev.* **1977**, *84*, 327–352.
- (37) See section 2.5 for further discussion.
- (38) In mathematics, a solid hypersphere embedded in a p -dimensional Euclidean space is also called a *closed* p -ball. The surface (boundary) of the solid hypersphere in that space is a p -dimensional hypersphere. Removal of the hypersphere that is removing the boundary of the p -dimensional solid hypersphere (closed p -ball) generates an *open* p -ball. For consistency, the terms solid hypersphere and hypersphere will be used in this work. For additional discussion on these points see the following: *Encyclopedic Dictionary of Mathematics*; MIT Press: Cambridge, MA, 1980; Vol. 1. Oden, J. T.; Demkowicz, L. F. *Applied Functional Analysis*. Chemical Rubber Publishing Company: Boca Raton, FL, 1996. Lastly, the p -dimensional hypersphere can also be considered as a manifold in a $(p-1)$ -dimensional subspace that is embedded within the p -dimensional Euclidean space, but this approach will not be exploited here. See, for example: Small, C. G. *The Statistical Theory of Shape*; Springer: New York, 1996.
- (39) Matoušek, J. *Lectures on Discrete Geometry*; Springer: New York, 2002.
- (40) Scott, D. W. *Multivariate Density Estimation*; John Wiley & Sons: New York, 1992.
- (41) See section 4.3 for further discussion on this point.
- (42) *Matlab*; MathWorks: Natick, MA. <http://www.mathworks.com/products/matlab/> (accessed Mar 21, 2007).
- (43) *Molecular Operating Environment (MOE)*; Chemical Computing Group Inc.: Montreal, Quebec, Canada, 2005. <http://www.chemcomp.com> (accessed Mar 21, 2007).
- (44) They can also be considered as lying on an closed $(p-1)$ -dimensional hyperspherical manifold.³⁸
- (45) See sections 2.4 and 2.5 for additional mathematical details.
- (46) Altman, D. G. *Practical Statistics for Medical Research*; Chapman & Hall: London, 1991; pp 285–288.
- (47) Agrafiotis, D. K. A Constant Time Algorithm for Estimating the Diversity of Large Chemical Libraries. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 159–167.
- (48) In the case of a complete molecular basis set, all molecular vectors will lie on an p -dimensional hypersphere; that is, they will lie on a closed $(p-1)$ -dimensional hyperspherical manifold.^{38,44}
- (49) Note that, although the diagram in Figure 6 is in two dimensions, the argument holds in higher dimensions as well.
- (50) This can be seen if one considers a uniform random distribution of vectors on the surface of a unit sphere.
- (51) In the case of zero squared error, which cannot be realized in practice, all molecular vectors will lie on the surface of the hypersphere.⁴⁸
- (52) Wegman, E. J. Hyperdimensional Data Analysis Using Parallel Coordinates. *J. Am. Stat. Assoc.* **1990**, *85*, 664–675.
- (53) We are concerned primarily with the “outer” hyperspherical shell that lies in the last small percent of the radius (r), say between $0.9r$ and $1.0r$.
- (54) In the case of a complete set of basis vectors, where $p = \infty$, all molecular vectors are “perfectly” represented by the basis set and have unit norms. Thus, they all lie on the surface of an infinite-dimensional hypersphere.^{38,44}
- (55) Although not considered here, molecular fragments, in principle, can also be used as suitable basis-set elements.
- (56) cf. the methods presented in refs 23 and 24.
- (57) Approximately degenerate principal components have approximately equal eigenvalues.
- (58) Domine, D.; Devillers, J.; Chastrette, M.; Karcher, W. Non-Linear Mapping for Structure–Activity and Structure–Property Modeling. *J. Chemom.* **1993**, *7*, 227–242.
- (59) Rassokhin, D. N.; Lobanov, V. S.; Agrafiotis, D. K. Non-Linear Mapping of Massive Data Sets by Fuzzy Clustering and Neural Networks. *J. Comput. Chem.* **2001**, *22*, 373–386.
- (60) Borg, I.; Groenen, P. *Modern Multidimensional Scaling—Theory and Applications*; Springer: New York, 1997.
- (61) Tenenbaum, J. B.; de Silva, V.; Langford, J. C. A Global Geometric Framework for Non-Linear Dimensionality Reduction. *Science* **2000**, *290*, 2319–2323.
- (62) Roweis, S. T.; Saul, L. K. Non-Linear Dimensionality Reduction by Local Linear Embedding. *Science* **2000**, *290*, 2323–2326.
- (63) Friedman, J.; Tukey, J. A Projection Pursuit Algorithm for Exploratory Data Analysis. *IEEE Trans. Comput.* **1974**, *C 23*, 881–889.
- (64) Agrafiotis, D. K. Stochastic Proximity Embedding. *J. Comput. Chem.* **2003**, *24*, 1215–1221.
- (65) Agrafiotis, D. K.; Xu, H. A Geodesic Framework for Analyzing Molecular Similarities. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 475–484.
- (66) Donoho, D. L.; Grimes, C. Hessian Eigenmaps: Local Linear Embedding Techniques for High-Dimensional Data. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100*, 5591–5596.

CI600552N