

Using Molecular Quantum Similarity Measures under Stochastic Transformation To Describe Physical Properties of Molecular Systems

Xavier Gironés and Ramon Carbó-Dorca*

Institute of Computational Chemistry, University of Girona, Campus Montilivi, 17071 Girona, Catalonia, Spain

Received September 26, 2001

The application of molecular quantum similarity measures (MQSM) to correlate physicochemical properties is reported. Satisfactory quantitative structure–property relationship (QSPR) models are obtained for three molecular sets, where boiling points and chromatographic retention times and indices are studied. In this work, MQSM are scaled using a stochastic transformation and related to molecular properties using the partial least-squares technique.

INTRODUCTION

One of the most promising subjects in present day computational chemistry is the characterization and quantification of molecular properties by means of theoretical descriptors. Given a set of molecules and a related molecular property, a quantitative structure–property or –activity relationship (QSPR/QSAR) model can be constructed by deriving descriptors from molecular structure, which will be related to the molecular feature by some statistical technique.

These disciplines were first started by Cross in 1863,¹ when he observed that toxicity of alcohols to mammals increased as water solubility of such alcohols decreased. Since then, QSPR/QSAR techniques have been largely developed using a variety of parameters as molecular descriptors of diverse structural characteristics. For example, Hammett σ values² are often used as electronic parameters, and other parameters have been devised to account for the shape, size, lipophilicity, polarizability, and many other molecular structural features. In addition, the recent and fast development of computer architectures has increased the computational power to stages that allow the application of the quantum theory to regular organic molecules at fairly accurate computational levels (semiempirical or *ab initio* methods with appropriate basis sets) within reasonable time limits and affordable costs. A number of reviews have been published^{3–5} concerning the historical development, generation of descriptors, and different methodologies in the QSPR/QSAR fields.

Arising from the conception of a molecule provided by quantum chemistry, and according to the postulate that its electronic density function contains all accessible information contained within this molecule, molecular quantum similarity measures (MQSM)^{6–12} stand as a general efficient tool to solve actual chemical problems. MQSM methodology has been successfully applied within pharmacological^{13–18} and toxicological^{19–21} problems. MQSM are based on the concept of molecular similarity^{22,23} and are related to a self-evident

molecular similarity principle: “the more similar two molecules are, the more similar properties they will possess”. This last statement requires a procedure to compare the molecules, and MQSM easily establish the degree of similarity on the basis of the electronic distribution, namely first-order density functions (DF), of the molecular structures.

In this work, it is intended to use the MQSM methodology to correlate boiling points and chromatographic parameters of simple organic molecules. For this purpose, this paper is structured as follows: a survey around the employed protocol is given first, next the presentation of the molecular sets, as well as the results achieved, are described, and the concluding remarks are finally given.

MATERIAL AND METHODS

MQSM and Stochastic Transformation. The practical application of MQSM relies on the use of first-order density functions. These functions are quantum-mechanical observable elements producing information on the molecular electron distribution. Within this framework, two molecular structures are considered to be similar if their electron distributions are similar. Thus, among other possibilities,⁸ quantitative measures of the similarity between two molecules can be defined as the direct volume integral between their density functions weighted by the *Coulomb* operator:

$$Z_{AB} = \int \rho_A(\mathbf{r}_1) |\mathbf{r}_1 - \mathbf{r}_2|^{-1} \rho_B(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2 \quad (1)$$

Here $\rho_A(\mathbf{r})$ and $\rho_B(\mathbf{r})$ are the first-order electron density functions of molecules A and B, respectively, and Z_{AB} is the resulting quantum similarity measure. As MQSM depend on the relative position of both compared objects in space, a molecular alignment procedure is required. When studied molecular sets share common structural features, the topological geometrical superposition algorithm (TGSA)²⁴ is used, as it performs pairwise superpositions according to the molecular backbones. Once superposed, the overall set of MQSM for a series of N molecules is computed and collected in the way of an $(N \times N)$ similarity matrix form, which can be

* Corresponding author phone: (+34)972418367; fax: (+34)972418356; e-mail: quantum@iqc.udg.es.

used to extract the appropriate correlation information and build QSPR/QSAR descriptors. Once the similarity matrix has been computed, a possible scaling of the MQSM can be done by means of eq 2, providing the recently described *stochastic transformation*²⁵ of MQSM. The construction of this stochastic MQSM matrix conceives that the sum of every row elements can be used as scale factors to obtain a new row. This procedure originates a new nonsymmetric square matrix, which will act as a set of input descriptors in any subsequent statistical procedure.

$$S_{ij} = Z_{ij} \left(\sum_j^N Z_{ij} \right)^{-1} \quad (2)$$

Promolecular Atomic Shell Approximation (ASA). To avoid expensive theoretical calculations, the promolecular atomic shell approximation (ASA)^{26–28} has been used here to construct molecular first-order electron density. Within ASA the molecular density is expressed as a sum of discrete atomic density contributions, which are taken as a superposition of 1S Gaussian functions and fitted to atomic ab initio ones. Once the overall atomic densities are built, each molecular density function can be constructed by adding as appropriate these elementary atomic building blocks. Since it has been proved that the MQSM from fitted densities built in this way differ by up to a 2% from the ab initio ones,²⁷ their use is clearly justified.

Molecular Alignment. MQSM are dependent on the relative position of the molecules under comparison. In this way, the topo-geometrical approach (TGSA)²⁴ has been used to perform the needed pairwise molecular alignments. This molecular superposition method overlays the involved molecules according to the maximal common substructure shared by the analyzed compounds.

Molecular Modeling. All molecules have been constructed using the Ampac 6.55²⁹ software package. The geometry optimization of the resulting structures has been carried out in gas phase with the same program at the AM1³⁰ computational level.

Treatment of Quantum Similarity Matrixes and Model Building. Common chemometric tools may be applied to deal with similarity matrices. Particularly, partial least squares (PLS)^{31,32} stands as an ideal technique to obtain a generalized regression to model the association between the matrices **X** (descriptors) and **Y** (responses). In computational chemistry, its main use is to model the relationship between computed variables, which together characterize the structural variation of a set of *N* compounds and a property of interest measured on those *N* substances.^{33–35} This variation of the molecular skeleton is condensed into the matrix **X**, whereas the analyzed properties are recorded into **Y**. In PLS, the matrix **X** is commonly built up from nonindependent data, as it uses more columns than rows; hence it is not called the *independent* matrix but a predictor or descriptor matrix. A good review, as well as its practical application in QSAR, is found in ref 36, and a detailed tutorial in ref 37.

Unlike regression, PLS is not based on the assumption of independent and precise **X** variables, but it is rather based on the more realistic assumption that **X** contains more or less collinear and noisy parameters. PLS summarizes these **X** variables by means of a few orthogonal score vectors (**t**_{*a*}

∈ **T**), and the matrix **Y** is also resumed in few score vectors (**u**_{*a*} ∈ **U**) that are not orthogonal. Plots of columns from **T** and **U** provide visual representation of the configuration of the observations in the **X** or **Y** spaces, respectively. To speed up the PLS routine, the dimensionality of the similarity matrix is a priori reduced in two ways:

(a) removal of columns with a variation coefficient below 2%, thus eliminating those descriptors that are roughly a vector of constants;

(b) extraction of all columns that, after pairwise correlation, score a goodness-on-fit coefficient larger than 0.85, thus eliminating redundant descriptors.

The PLS procedure allows one to derive a number of *factors* and *weights*, which are used to describe the desired properties. QSPR models are built up from these factors and weights.

In this work, all obtained models are evaluated by commonly used statistical parameters: goodness-on fit (*r*²);³⁸ root-mean-square error between experimental and predicted values (*s*)³⁸ by *leave-one-out*;³⁹ predictive capacity (*q*²).³⁸ In addition, and in order to avoid chance correlations and excess of parameters, models are submitted to random tests, where the properties are randomly permuted in their positions and the entire modeling procedure is repeated a number of times. If satisfactory correlations are found within the random test, the model obtained should not be trusted, as the methodology used may be potentially capable to correlate any kind of data.

RESULTS AND DISCUSSION

In this section, it is intended to prove the usefulness of the proposed procedure with three different molecular sets. The molecular properties analyzed consist of boiling points, gas chromatographic retention times, and gas–liquid chromatographic retention indices, respectively. The boiling points are analyzed for a molecular set composed of 529 saturated hydrocarbons⁴⁰ (linear, branched, and cyclic). The retention times are contemplated for a set of 152 diverse chemical compounds.⁴¹ Finally, the retention indices are evaluated over a set of 50 phenol derivatives.⁴²

Boiling Point for 529 Saturated Hydrocarbons. Boiling point at normal pressure (bp) has become a benchmark property for the evaluation of novel QSPR descriptors.^{43,44} In this way, it is interesting to test the present proposed MQSM methodology for such a homologous molecular series. For this purpose, a set of 529 saturated hydrocarbons, whose data are presented in Table 1, was tested for correlation search with their bp using the previously described methodology. The initial 529 × 529 matrix was reduced to 529 × 513 after the redundant variable reduction. QSPR models have been constructed up to the usage of 15 PLS factors, in which both *r*² and *q*² have been inspected. The optimal number of parameters to be used in the final model, which in this case has been fixed to 6 PLS factors, has been chosen from the evolution of these magnitudes, as presented in Figure 1.

As observed in Figure 1, a sharp increase in the predictive capacity occurs up to the sixth PLS factor and a slow increase until the eighth one. A further increase in the number of parameters results, in despite of the expected increase of *r*², in a decrease of the predictive power of the generated model. In this way, the chosen number of PLS factors used was 6,

Table 1. SMILES Notation^{47,48} and Boiling Points (°C) of 529 Saturated Hydrocarbons⁴⁰

compd	bp	compd	bp	compd	bp
CCC	32.8	CCCCCCCCC	99.2	CCCC(C)(C)CCCC	172.5
C ₁ CC ₁	42.1	CC(C)CCCCC	118.2	CCC(CC)(C)CC	166
CCCC	8	CCC(C)CCCC	109.8	CC(C)(C)CC(C)CCC	164.5
CC(C)C	0.7	CCCC(C)CCCC	114.8	CC(C)C(C)C(C)CCC	179.5
C ₁ (C)CC ₁	12.6	CC(C)CCCC(C)C	106.5	CC(C)CCC(C)(C)CC	171.5
C ₁ CCC ₁	0.5	CCCCC(CC)CC	140.5	CC(C)C(CCC)C(C)C	175.8
C ₁₂ CC ₁ C ₂	11.7	CCCC(CC)CCC	168.3	CCC(C)C(C)C(C)CC	160.5
CCCCC	36	CC(C)CCC(C)CC	164.6	CC(C)(C)CC(C)CC	165
CC(C)CC	46	CC(C)C(C)CCCC	162.5	CCC(C)(CC)CCCC	172.5
CC(C)(C)C	39	CCC(C)C(C)CCC	166	CC(C)C(C)C(CC)CC	186.2
C ₁ (CC ₁)CC	33.5	CCC(C)CC(C)CC	161.4	CC(C)C(CC)C(C)CC	174.5
C ₁ (C)CC ₁ C	35.9	CCC(CC)CC(C)C	142.5	CC(C)CC(C)(C)CCC	164
C ₁ (C)(C)CC ₁	32.6	CC(C)(C)CCCCC	147.8	CC(C)(C)C(C)CCCC	160.5
C ₁ (C)CCC ₁	20.6	CCC(C)C(C)CCC	140.7	CCC(C)(C)CC(C)CC	159.5
C ₁ CCCC ₁	36.3	CCCC(CC)C(C)C	161	CCCC(C)(CC)CCC	162
C ₁ C ₂ CC ₁ C ₂	49.3	CC(C)CC(C)CCC	146	CC(C)(C)CC(C)C(C)C	172.5
C ₁ CC ₂ CC ₁₂	36	CCC(C)CC(C)CC	137.5	CC(C)(C)CCC(C)(C)C	165.5
C ₁ CC ₁₂ CC ₂	27.8	CC(C)(C)CCC(C)C	137.5	CC(C)C(C)C(C)C(C)C	159.5
C ₁₂ (C)CC ₁ C ₂	9.5	CCC(C)(C)CCCC	131	CC(C)C(C)(C)CCCC	163
CCCCC	71	CCCC(C)(C)CCC	128	CC(C)CC(C)(CC)CC	150
CC(C)CCC	76	CC(C)C(C)C(C)CC	152.5	CC(C)(C)C(CC)CCC	168
CCC(C)CC	83	CC(C)(C)CC(C)CC	169	CC(C)(C)C(C)CC(C)C	153
CC(C)C(C)C	81	CC(C)C(CC)C(C)C	161.8	CCC(CC)(CC)CCC	168.5
CC(C)(C)CC	69.5	CCC(C)(CC)CCC	163	CC(C)C(C)C(C)C(C)C	160
C ₁ (CC ₁)CCC	60.5	CC(C)CC(C)(C)CC	163	CCC(C)(C)C(C)CCC	152
C ₁ (CC ₁)C(C)C	55	CC(C)(C)C(C)CCC	164	CCC(C)C(C)C(C)CCC	157.5
C ₁ (CC)CC ₁ C	69	CCC(CC)(CC)CC	157.8	CC(C)C(C)(C)CC(C)C	158.5
C ₁ (C)(CC ₁)CC	58.3	CC(C)C(C)(C)CCC	168	CCC(C)(C)C(CC)CC	159.5
C ₁ (C)C(C)C ₁ C	63	CC(C)(C)C(CC)CC	155	CC(C)C(C)(CC)CCC	152
C ₁ (C)(C)CC ₁ C	57	CCC(C)(C)C(C)C	154	CC(C)(C)CC(C)(C)CC	155.1
C ₁ (CC)CCC ₁	63	CC(C)(C)C(C)(CC)CC	152.5	CC(C)(C)C(C)C(C)CC	148.3
C ₁ (C)CC(C ₁)C	52.6	CC(C)(C)CC(C)(C)C	157	CCC(C)(CC)C(C)CC	153.5
C ₁ (C)CCC ₁ C	70.7	CC(C)(C)C(C)C(C)C	159	CC(C)(C)C(CC)C(C)C	146.7
C ₁ (C)(C)CCC ₁	59	CC(C)C(C)(C)C(C)C	149	CC(C)C(C)C(C)(C)CC	146
C ₁ CCC(C ₁)C	62	CC(C)(C)C(C)(C)CC	153	CC(C)C(CC)(CC)CC	158
C ₁ CCCCC ₁	53.6	C ₁ (CC ₁)CCCCC	145	CC(C)C(C)(C)C(C)CC	155.5
C ₁₂ CCC(C ₁)C ₂	71.8	C ₁ (CC ₁)C(C)CCCC	152	CC(C)C(C)(CC)C(C)C	154
C ₁ CC ₁ C ₂ CC ₂	80.7	C ₁ (CC ₁ C)CCCCC	148.2	CC(C)(C)C(C)(C)CC	145.5
C ₁ CC ₂ C ₁ CC ₂	68.7	C ₁ (CC)CC ₁ CCC	151	CC(C)(C)C(C)C(C)(C)C	180
C ₁₂ CCCC ₁ C ₂	60.3	C ₁ CC ₁ (CC)CCCC	146.7	CCC(C)(C)C(C)(C)CC	172
C ₁ CCC ₁₂ CC ₂	63.3	CC(C)(C)CC ₁ (C)CC ₁	151.5	CC(C)(C)C(C)(CC)CC	170.7
C ₁₂ (C)CCC ₁ C ₂	58	CC(C)CC ₁ CC ₁ (C)C	149	CC(C)(C)C(C)(C)C(C)C	176.5
C ₁₂ (C)CC ₁ (C ₂)C	49.7	C ₁ (C)(C)CC ₁ C(C)(C)C	151	C ₁ (C)(CC ₁)CCCC(C)C	172
CCCCCCC	108.5	C ₁ (C)(CC)CC ₁ (C)CC	143	CC(C)CC ₁ CC ₁ C(C)C	153
CC(C)CCCC	104	C ₁ (C)(C)C(C)(C)C ₁ (C)C	138.8	C(C)CCC ₁ CC ₁ (C)C	175.7
CCC(C)CCC	105	C ₁ CCC ₁ C(C)CC	146.5	CCCC ₁ CC ₁ (C)CCC	173.9
CCC(CC)CC	107	C ₁ CCCC ₁ CCCC	142	CC ₁ (C)CC ₁ C(C)(C)CC	173.5
CC(C)CC(C)C	110	C ₁ CCCC ₁ CC(C)C	148.2	C ₁ (C)C(C ₁)C(C)C(C)(C)C	160
CC(C)C(C)CC	107.5	C ₁ C(C)CCC ₁ CCC	150	C ₁ (CCC ₁ C(C)C)C(C)C	172
CC(C)(C)CCC	105.5	C ₁ CCCC ₁ C(C)CC	150.5	C ₁ (CC)C(C)C(C)C ₁ CC	173.9
CCC(C)(C)CC	110	C ₁ C(C)CCC ₁ CC	132.2	C ₁ (C)(C)C(CC)CC ₁ CC	160.5
CC(C)(C)C(C)C	110.5	C ₁ (C)CCCC ₁ CCC	154.5	C ₁ (C)C(CC ₁ (C)C)C(C)C	148.5
C ₁ CC ₁ CCCC	116	C ₁ (CC)CCCC ₁ CC	141	C ₁ CCC(C ₁)CCCC	172
C ₁ CC ₁ C(C)CC	96.5	C ₁ C(C)CCC ₁ C(C)C	131	C ₁ CCC(C ₁)CCC(C)C	171
C ₁ (C)CC ₁ CCC	98.5	C ₁ (C)CCCC ₁ C(C)C	135	C ₁ (C)CCC(C ₁)CCCC	163
C ₁ (CC)CC ₁ CC	100	C ₁ CCCC ₁ (C)CCC	127.2	C ₁ CCC(C ₁)C(C)CCC	164
C ₁ (C)(CC ₁)CCC	103	C ₁ C(C)C(C)CC ₁ CC	149	C ₁ CCC(C ₁ C)CCCC	141.5
C ₁ (CC ₁ C)C(C)C	81.5	C ₁ C(C)CC(C)C ₁ CC	143	C ₁ CC(CC ₁ C)CC(C)C	180.9
C ₁ (CC ₁)C(C)(C)C	92	C ₁ CC(C)C(C)C ₁ CC	153	C ₁ CCCC ₁ C(CC)CC	171.3
C ₁ (CC ₁)(CC)CC	71.5	C ₁ CC(C)C(C)CC	140.2	C ₁ CCCC ₁ CC(C)CC	173.4
C ₁ (CC)C(C)C ₁ C	74	C ₁ (CCCC ₁)C(C)C	125	C ₁ CCCC ₁ C(C)C(C)C	169
C ₁ (C)(CC ₁)C(C)C	78	C ₁ (CC)(CC)CCCC ₁	125.5	CC(C)C ₁ CCCC ₁ CC	179.3
C ₁ (C)(C)CC ₁ CC	84	C ₁ (CC)CCC(C ₁)(C)C	121	C ₁ (CC)CCC(C ₁ C)CC	175.5
C ₁ (C)(CC)CC ₁ C	98	C ₁ (C)CCC(C ₁)(C)CC	130.8	CC(C)(CC)C ₁ CCCC ₁	172
C ₁ (C)(C)C(C)C ₁ C	90.3	C ₁ CCC(C ₁)(C)C(C)C	124.5	CC(C)C ₁ CCC(C)C(C)C ₁	170
C ₁ (C)(C)CC ₁ (C)C	93	C ₁ CC(CC)C(C ₁)(C)C	148.7	CC(C)C ₁ CCC(C)(C)C ₁	176
C ₁ CCC ₁ CCC	90	C ₁ C(C)C(C)C(C)C ₁ C	156.6	C ₁ (C)C(C)CC(C ₁ C)CC	167
C ₁ CCC ₁ C(C)C	84.9	C ₁ CCC(C)C ₁ (C)CC	148	C ₁ (CC)C(C)CC(C)C ₁ C	168.5
C ₁ C(C)CC ₁ CC	81.1	C ₁ C(C)C(C)CC ₁ (C)C	148.3	C ₁ (CCCC ₁ (C)C)C(C)C	171
C ₁ CC(C)C ₁ CC	80.5	C ₁ (C)CC(C)CC ₁ (C)C	154.3	C ₁ (C)C(C)C(C)C(C)C ₁ C	174.3
C ₁ CCCC ₁ CC	88.6	C ₁ (C)C(C)CCC ₁ (C)C	148.2	C ₁ C(C)(C)CC(C)C ₁ (C)C	175
C ₁ C(C)CCC ₁ C	91	C ₁ C(C)(C)CCC ₁ (C)C	149.5	C ₁ CCCC(C ₁)CCCC	171.5
C ₁ (C)CCCC ₁ C	81.5	C ₁ (C)CCC(C)C ₁ (C)C	150.5	C ₁ CCCC(C ₁)CC(C)C	179.5
C ₁ CCCC ₁ (C)C	79.1	C ₁ (C)(C)CCCC ₁ (C)C	141	C ₁ CC(CCC)CCC ₁ C	168

Table 1. (Continued)

compd	bp	compd	bp	compd	bp
C ₁ CCCC(C ₁)C	85.2	C ₁ CCCCC ₁ CCC	145	C ₁ (C)CCCC(C ₁)CCC	166.6
C ₁ CCCCC ₁	78	C ₁ CCCCC ₁ C(C)C	145	C ₁ CCCC(C ₁)C(C)CC	166.5
C ₁ CC ₁ CCC ₂ CC ₂	76	C ₁ CC(C)CCC ₁ CC	142.5	C ₁ CC(C)CCC ₁ CC	177.5
C ₁ CC ₂ CCC ₁ C ₂	100.7	C ₁ C(C)CCCC ₁ CC	138	C ₁ CC(C)CC(C ₁)CC	172.5
C ₁ CC ₂ CC(C ₁)C ₂	92.7	C ₁ (C)CCCCC ₁ CC	147	C ₁ (CC)CCCCC ₁ C	153
C ₁₂ CCCC ₁ CC ₂	89.5	C ₁ (C)CC(C)CC(C ₁)C	151	C ₁ CC(CCC ₁ C)C(C)C	160.3
C ₁₂ CCCCC ₁ C ₂	94	C ₁ (C)CCC(C)C(C ₁)C	145	C ₁ CCCC(C)C ₁ CC	158
C ₁ CCC ₁₂ CCC ₂	103.5	C ₁ CCC(C)C(C)C ₁ C	151	C ₁ C(C)CCCC ₁ C(C)C	153
C ₁ CCCC ₁₂ CC ₂	91.3	C ₁ CCC(C)C(C ₁)CC	133	C ₁ C(C)CC(C)CC ₁ CC	155
C ₁₂ CCC(C)C ₁ C ₂	95.6	C ₁ (C)CCC(C)(C)CC ₁	135.5	C ₁ (C)CCCCC ₁ C(C)C	167
C ₁₂ CCCC ₁ C ₂ C	97.9	C ₁ CCC(C)(C)CC ₁ C	143	C ₁ CCCCC ₁ (C)CCC	161.5
C ₁₂ CCC(C ₁)(C ₂)C	101	C ₁ CCC(C)(C)C(C ₁)C	138	C ₁ (C)C(C)CCCC ₁ CC	182.8
C ₁₂ CCCC ₁ (C ₂)C	118.4	C ₁ CCCCC(C ₁)CC	135.9	C ₁ (CCCCC ₁)C(C)C	164.2
C ₁ C ₂ (C)CC ₁ (C)C ₂	98.5	C ₁ CC(C)CCC(C ₁)C	142.5	C ₁ (CC)(CC)CCCC ₁	173
C ₁₂ (C)CCC ₁ (C ₂)C	90	C ₁ CCC(C)CC(C ₁)C	127	C ₁ CC(C)CCC ₁ (C)CC	170.7
C ₁ CC ₁₂ CC ₂ (C)C	92	C ₁ CCCC(C)C(C ₁)C	129.5	C ₁ C(C)C(C)CC(C ₁)C	177.5
C ₁₂ (C)CC ₁ C ₂ (C)C	93.5	C ₁ CCCCC(C ₁)C	134	C ₁ C(C)CCCC ₁ (C)CC	161
C ₁₂ CC ₂ CC ₃ C ₁ C ₃	80.5	C ₁ CCCC(C)CCC ₁	118.2	C ₁ (C)C(C)CC(C)CC ₁ C	174
C ₁₂ CC ₁ C ₃ CC ₂ C ₃	89.8	C ₁ CCCCCCCC ₁	138	C ₁ CCCCC ₁ (C)C(C)C	191.4
C ₁₂ CCC ₃ C(C ₁)C ₃ C ₂	79.2	C ₁ CCC ₂ CC ₁ CCC ₂	135	C ₁ (C)C(C)CCC(C)C ₁ C	183.5
C ₁₂ CCCC ₃ C ₁ C ₂₃	86.1	C ₁ CCC ₁ CC ₂ CCC ₂	156.7	C ₁ (C)(C)CC(C)CC(C ₁)C	181.5
C ₁₂₃ CC ₁ CCC ₂ C ₃	80.9	C ₁ CCCCC ₂ CC ₁ CC ₂	154.8	C ₁ (C)(C)CCC(C)CC(C ₁)C	185.5
C ₁₂ CC ₃ C ₄ C ₁ C ₂ C ₃₄	137.5	C ₁ CCCC ₂ CCCC ₁₂	150.8	C ₁ (C)CCC(C)CC ₁ (C)C	170.5
C ₁₂ CC ₁ C ₃ C ₄ C ₂ C ₃₄	153	C ₁ CCCC ₂ CCCC ₂ C ₁	150	C ₁ (C)(C)CCC(C)(C)CC ₁	193.6
CCCCCCCC	106	C ₁ (CCCCC ₁)C ₂ CC ₂	154.3	C ₁ (C)(C)CCCC(C ₁)(C)C	142
CC(C)CCCC	142	C ₁ CCCCC ₂ CC ₂ C ₁	139.5	C ₁ (C)(C)CCCC(C)C ₁ C	174.1
CCC(C)CCCC	149	C ₁ C(C)C ₂ CCC ₁ C ₂	144.8	C ₁ (C)(C)CCCCC ₁ (C)C	166.8
CCCC(C)CCC	136	C ₁ CC ₂ CCC(C ₁)C ₂	149.4	C ₁ CCCCC(C ₁)CCC	167.8
C(C)(C)CCC(C)C	125	C ₁ CC ₂ CC(C ₁)C(C)C ₂	152	C ₁ C(C)CCC(C)CC ₁ C	165.7
C(C)CC(C)CC	120.5	C ₁ C(C)CC ₂ C ₁ CCC ₂	136	C ₁ CC(C)C(C)CCC ₁ C	160
C(C)(C)CC(C)CC	103	C ₁₂₃ CCCC ₁ C ₂ CCC ₃	136.6	C ₁ CC(C)CC(C)C(C ₁)C	165.1
C(C)(C)C(C)CCC	111	C ₁₂₃ CCC ₁ C ₂ CCC ₃	145.1	C ₁ CCC(C)C(C)C(C ₁)C	168
C(C)C(C)C(C)CC	115	C ₁ CCC ₂ CCC(C)C ₁₂	163.7	C ₁ CCC(C)(C)CC(C ₁)C	158.5
CC(C)(C)CCCC	102	C ₁₂₃ CC ₁ C ₂ CCCC ₃	153.8	C ₁ CCC(C)C(C)(C)CC ₁	164
CCC(C)C(C)C	129	C ₁ CC ₂ (C)CCC ₁ CC ₂	151	C ₁ CCCCC(C ₁)CC	157
CC(C)C(C)C(C)C	137	C ₁ CC ₂ CC ₁ CC ₂ CC	157	C ₁ CC(C)CCCCC ₁ C	161.8
CCC(C)(C)CCC	136	C ₁₂ CCCC(C ₂)(C)CC ₁	150	C ₁ (C)CCCCCC(C ₁)C	153
CC(C)(C)CC(C)C	133	C ₁ CC ₂ CC ₁ C(C)C ₂ C	168.2	C ₁ CCCCCC(C)C ₁ C	159.5
CCC(C)(C)CC	141	C ₁ CC ₂ (CC)CC ₁ CC ₂	175	C ₁ CCCCCC(C ₁)(C)C	164
CC(C)(C)C(C)CC	125	C ₁ CC ₂ CC(C ₂)C ₁ (C)C	157.8	C ₁ CCCCC(C)CCC ₁	159.7
CC(C)C(C)(C)CC	128	C ₁ CCCC(C)C ₁₂ CC ₂	174.5	C ₁ CCCCCCC ₁	159
C(C)(C)(C)C(C)(C)C	128	C ₁ CCCC ₂ (C)CC ₂ C ₁	150.8	C ₁₂ CCCCC ₂ CCCC ₁	155
C(CCCC)C ₁ CC ₁	130.5	C ₁ C ₂ CCC(C ₂)C ₁ (C)C	142.8	C ₁ CCC(C ₁)C ₂ CCCC ₂	160
CC(CCC)C ₁ CC ₁	125	C ₁ C(C)C ₂ (C)CC ₁ CC ₂	144	C ₁ CCCCC ₂ CCCC ₁₂	145
C ₁ (CC ₁)CCCC	117	C ₁ CCCCC ₁₂ CC ₂ C	142.4	C ₁₂ CCCC(C ₂)CC(C ₁)C	166
C ₁ CC ₁ CC(C)(C)C	138	C ₁ CC ₂ (C)CCC ₁ C ₂ C	134	C ₁₂ CCCC(C ₂)C(C ₁)C	160
C ₁ CC ₁ C(C)C(C)C	125	C ₁ CC ₂ CCC ₁ C ₂ (C)C	143	C ₁₂ CCCC(C ₁)CCCC ₂ C	155.7
C ₁ (CCC)CC ₁ CC	115	C ₁ (C)(C)C ₂ CCCC ₁ C ₂	142.1	C ₁ C(C)CCC ₂ C ₁ CC ₂	166
CC ₁ CC ₁ CC(C)C	91	C ₁ CCCCC ₂ CC ₁₂ CC	136	C ₁ CCCCC ₂ C ₁ CC(C)C ₂	162.4
C ₁ (C)(C)CC ₁ CCC	126.1	C ₁ (C)CCCC ₁₂ CCC ₂ C	133.5	C ₁ CCC ₂ CCCC ₁ C ₂ C	159.5
C ₁ (C)(CC)CC ₁ CC	104	C ₁ CCCCC ₂ C ₁ C ₂ (C)C	140.5	C ₁ CCCCC ₂ (C ₁)CCCC ₂	148.2
C ₁ (C)(C)CC ₁ C(C)C	105	C ₁ (C)CCCC ₂ C ₁ C ₂ (C)C	136	C ₁ (C)CCCCC ₂ C ₁ CC ₂	167
C ₁ (C)(C)CC ₁ (C)CC	159	C ₁₂ (C)CCCC(C ₂)C ₁ (C)C	133.8	C ₁ CCCCC ₂ (C ₁)CCC ₂	159.7
C ₁ (C)(C)C(C)C ₁ (C)C	128	C ₁ CC ₂ (C)CC ₂ (C ₁)CC	132.7	C ₁₂ CCC(C ₂)CC ₁ CC	167
C ₁ (CCC ₁)CC(C)C	117.7	C ₁ (C)CC ₂ (C)CC ₂ (C ₁)C	140.6	C ₁ CCCCC ₂ CCC(C)C ₁₂	161.2
C ₁ (CCC)CC(C ₁)C	124	C ₁₂ CC ₃ CC(C ₁)C(C ₂)C ₃	138	C ₁₂ (C)CCCC(C ₁)CC ₂	157
C ₁ (CCC ₁)C(C)CC	106	C ₁ C ₂ CC ₃ CCC ₁ C ₃ C ₂	131.3	C ₁ CCC ₂ CCC(C)C ₁₂	147
C ₁ (CC)CCC ₁ CC	115.5	C ₁₂ CCC(C ₂)CC ₃ C ₁ C ₃	140.4	C ₁₂ CCC(C ₂)CC(C)C ₁ C	157
C ₁ (C)C(C)C(C)C ₁ C	108	C ₁ CC ₂ C ₃ C ₁ CC ₃ CC ₃	124	C ₁₂ CCC(C ₂)CC ₁ C	160
C ₁ (C)(C)CC(C ₁)(C)C	110	C ₁ CC ₂ (CC ₂)CC ₁₃ CC ₃	137.3	CC ₁ CC ₂ CC(C)CC ₂ C ₁	162
C ₁ CCCC ₁ CCC	105.9	C ₁ (CC ₁)(C ₂ CC ₂)C ₃ CC ₃	135.2	C ₁₂ CC(C)C(C ₂)CCC ₁ C	147.7
C ₁ CCCC ₁ C(C)C	108.9	C ₁ CC ₂ CCC ₁ C ₂₃ CC ₃	139	C ₁₂ (CC)CCC(C ₁)CC ₂	163
C ₁ C(C)CCC ₁ CC	94.4	C ₁₂ C(C)C ₂ C ₃ CC ₁ CC ₃	126.5	C ₁₂ CCC(C ₂)C(C)C ₁ C	152.8
C ₁ (C)CCCC ₁ CC	104.5	C ₁ CCC ₂ (CC ₂)C ₁₃ CC ₃	136.7	C ₁ CCCCC ₂ CCCC ₁₂ C	163
C ₁ (C)CC(C)CC ₁ C	100.5	C ₁ C ₂ CC ₃ C ₁ C ₃ (CC)C ₂	140.6	C ₁ CC(C)C ₂ CCCC(C)C ₁₂	164
C ₁ CCCC ₁ (C)CC	120.1	C ₁ C ₂ C ₃ CC ₁ C(C)(C)C ₂₃	130.7	C ₁ CCC ₂ C ₁ C(C)CC ₂ C	147
C ₁ (C)C(C)CCC ₁ C	117.4	C ₁ C ₂ CC ₃ C ₁ C ₃ (C)C ₂ C	133.6	C ₁₂ CCC(C ₂)CCC ₁ (C)C	163.8
C ₁ C(C)CCC ₁ (C)C	123	C ₁ C ₂ CC ₃ C ₁ (C)C ₃ (C)C ₂	145	C ₁₂₃ CCCC ₁ C ₃ CCC ₂ C	164
C ₁ (C)CCCC ₁ (C)C	119	C ₁₂ CC(C ₃ C ₁ C ₃)C ₄ C ₂ C ₄	137.7	C ₁₂ CCC(C ₂)CC ₁ (C)C	164
C ₁ CCCC(C ₁)CC	114.5	C ₁₂ CCC ₃ C ₂ C ₄ C ₁ C ₃₄	133.8	C ₁₂ (C)CCC(C)C(C ₂)CC ₁	152
C ₁ CC(C)CCC ₁ C	86	CCCCCCCC	122.3	C ₁ CCC ₂ CC(C)CC ₁₂ C	158
C ₁ CC(C)CC(C ₁)C	131	CC(C)CCCC	133	C ₁₂ (C)CCCC(C ₂)(C)CC ₁	155.7
C ₁ CCC(C)C(C ₁)C	126.4	CCC(C)CCCC	141.5	C ₁₂ CC(C)C(C ₂)CC ₁ (C)C	166
C ₁ CCCC(C ₁)(C)C	121	CCCC(C)CCCC	140.2	C ₁ CCCC ₂ CC ₁₂ CCC	147.9

Table 1. (Continued)

compd	bp	compd	bp	compd	bp
C ₁ CCCC(C)CC ₁	124.7	CC(C)CCCC(C)C	160	C ₁₂ CCC(C ₂)C(C)C ₁ (C)C	137
C ₁ CCCCCCC ₁	115	CCCCC(C)CCCC	174	C ₁₂ (C)CCC(C ₂)C(C)C ₁ C	158
C ₁ (CC ₁)CCC ₂ CC ₂	121.5	CCC(CC)CCCCC	155	C ₁₂ CCC(CC ₁ C)C ₂ (C)C	160
C ₁₂ CCCC ₁ CCC ₂	117	CC(C)CCCC(C)CC	175.5	C ₁₂ (C)CCC(C ₂)C(C ₁)(C)C	158
C ₁ (CCC ₁)C ₂ CCC ₂	104.9	CCCC(C)CCCC	219	C ₁₂ CC(CCC ₁ C)C ₂ (C)C	159
C ₁ CCC ₂ CCC ₂ C ₁	114	CC(C)CCC(C)CCC	189.5	C ₁₂ (C)CCC(C)(CC ₁)C ₂ C	148.7
C ₁ CCCCC ₂ C ₁ C ₂	131.8	CCCC(CCC)CCC	158.3	C ₁ C ₁ (C)CCC ₂ C ₁ C ₂ (C)C	166.3
C ₁ CC ₂ CC ₁ C(C)C ₂	121.8	CC(C)CC(C)CCCC	192.7	C ₁₂ (C)CCC(CC ₁)C ₂ (C)C	157
C ₁ CCCC ₁₂ CCC ₂	122.3	CCC(C)CCC(C)CC	160	C ₁ CCCCC ₁₂ CC ₂ (C)C	164
C ₁₂ CCC(CC ₁)C ₂ C	126.6	CC(C)C(C)CCCCC	151	C ₁ CC(C)C ₂ CC ₁₂ C(C)C	164
C ₁ (C)CCC ₂ C ₁ CC ₂	119.5	CCC(CC)CCC(C)C	143.5	C ₁ CCC ₂ (C)CC ₁₂ CCC	153
C ₁₂₃ CC ₁ C ₂ CCCC ₃	134	CCCC(C)CC(C)CC	153	C ₁ C(CC)CC ₂ (C)CC ₁₂ C	165
C ₁₂ (C)CCC(C ₂)CC ₁	149	CCCCCCC(C)(C)C	191.5	C ₁ CCC ₂ (C)CC ₁₂ C(C)C	169
C ₁ CCCC ₂ C ₁ C ₂ C	142	CC(C)CC(C)CCC	190	C ₁₂ CC ₁ CCC(C ₂)C ₃ CC ₃	153
C ₁ CCCC ₂ CC ₁₂ C	202	CC(C)CC(C)CC(C)C	193	C ₁ CC ₂ C ₁ C ₃ CCC ₂ CC ₃	155
C ₁₂ CC ₁ CC(C ₂)(C)C	171	CCC(C)C(C)CCCC	182	C ₁₂ CCC(CC ₂)C ₃ C ₁ CCC ₃	170
C ₁ CC ₂ (C)CC ₁ (C)C ₂	125.7	CCC(C)CC(C)CC	190	C ₁ (C)CC ₂ CCC ₁ C ₃ CC ₂₃	155.3
C ₁₂ CCCC ₁ C ₂ (C)C	117.6	CC(C)CCC(C)C(C)C	187	C ₁ C(C ₁)(C)C ₂ CC ₃ C ₃ CC ₃	162
C ₁₂ C(C ₂ (C)C)C ₁ (C)C	118.9	CC(C)C(CC)CCCC	178	C ₁ CCCC ₁₂ C ₃ CCCC ₂₃	174
C ₁₂ (C)CC ₁ (C)C ₂ (C)C	117.7	CCCC(C)C(C)CCC	173.5	C ₁ CCCC ₂ (CC ₂)C ₁₃ CC ₃	164
C ₁₂ CC ₁ CC ₃ CC ₃ C ₂	109.1	CCCC(CCC)C(C)C	189.9	C ₁₂ CCC(C ₂)C ₃ C ₁ C ₃	169.4
C ₁₂ CC ₁ CCC ₃ C ₂ C ₃	118.5	CC(C)(C)CCCC(C)C	185	C ₁ (C)(C)C ₂ C ₃ CC ₁ (C)CC ₂₃	159.3
C ₁₂ CCC(C ₂)C ₃ C ₁ C ₃	109.4	CCC(CC)C(C)CCC	182	C ₁ C ₂ CC ₃ C ₁ (C)C ₂ C ₂ (C)C	170.5
C ₁ CC ₂ C ₃ CCC ₂ C ₁₃	115.6	CC(C)C(C)CC(C)CC	183	C ₁₂ (C)C ₃ C ₂ CC(C ₃)C ₁ (C)C	168
C ₁ (C)C ₂ CC ₃ C ₁ C ₃ C ₂	117.7	CCC(C)C(CC)CCC	183	C ₁₂ CC ₂ C ₃ CC ₁ C ₄ (CC ₄)C ₃	166
C ₁₂₃ CC ₁ C ₂ C ₄ (C ₃)CC ₄	106.8	CCC(C)(C)CCCC	182	C ₁ CC ₂ C ₃ CC ₄ C(C ₃)C ₂₄ C ₁	188
C ₁ C ₂ CC ₃ C ₁ (C)C ₃ C ₂	115.6	CC(C)CC(C)C(C)CC	178	C ₁₂ CC ₁ C ₃ C ₄ C ₂ C ₃ C ₃ C ₄₅	167
C ₁ CC ₂ (CC ₂)C ₁₃ CC ₃	113.5	CC(C)(C)CCC(C)CC	174		
C ₁₂ CC ₃ C ₄ CC ₁ C ₂ C ₃₄	112	CC(C)CC(CC)C(C)C	174.5		

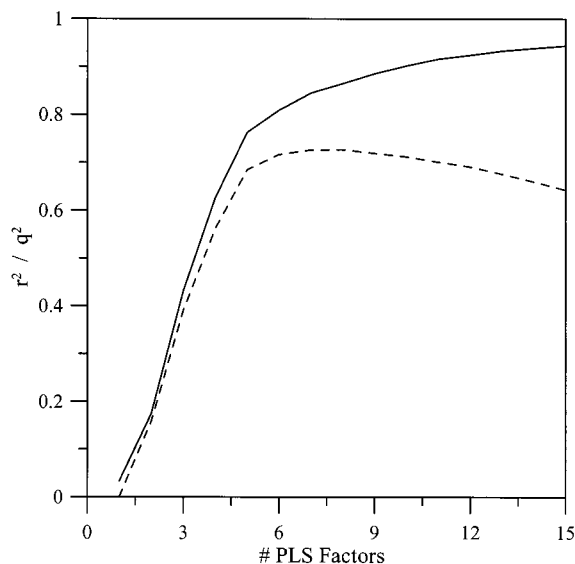


Figure 1. Evolution of r^2 (continuous line) and q^2 (dashed line) versus the number of PLS factors involved in the construction of the QSPR model for a set of 529 saturated hydrocarbons.

as the results provide the optimal balance between predictivity and descriptors applied, and a further increase does not lead to noticeable improved results and could lead to overfitted models. Thus, the final equation and statistical results for this molecular set are

$$\text{bp} = 3147.79\mathbf{f}_1 + 15269.10\mathbf{f}_2 + 34944.14\mathbf{f}_3 + 43295.57\mathbf{f}_4 + 47121.10\mathbf{f}_5 + 28552.70\mathbf{f}_6$$

$$r^2 = 0.808 \quad q^2 = 0.716 \quad s = 20.575$$

The results are also graphically presented by plotting the experimental bp versus the predicted ones from the PLS

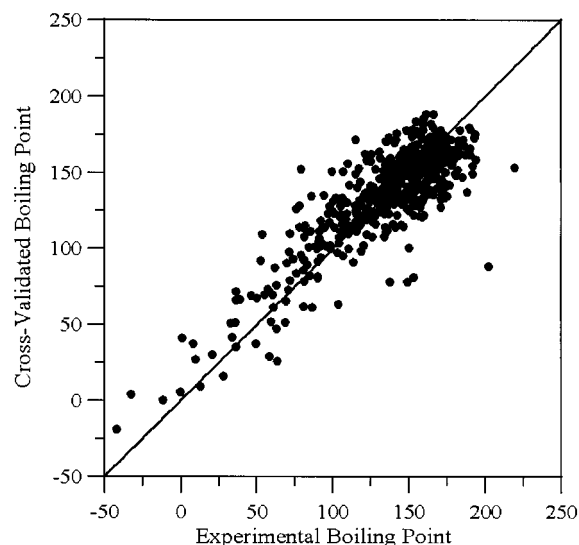


Figure 2. Experimental versus predicted boiling points for a set of 529 saturated hydrocarbons.

procedure, as shown in Figure 2. In addition, the random test was carried out permuting 100 times the bp values and reconstructing the QSPR model each time. The results are graphically compiled in Figure 3, where each point represents the values of r^2 and q^2 achieved in each permutation.

As seen from the statistical and graphical results, satisfactory results are obtained for this molecular set. Most of the compounds are correctly predicted, within a close margin, even few deviations are present, as displayed in Figure 2. The random test shows a clear separation between the original point (+) and the permuted ones (●). Among the permuted models, the highest values of r^2 and q^2 achieved are 0.475 and 0.029, respectively, whereas the average values

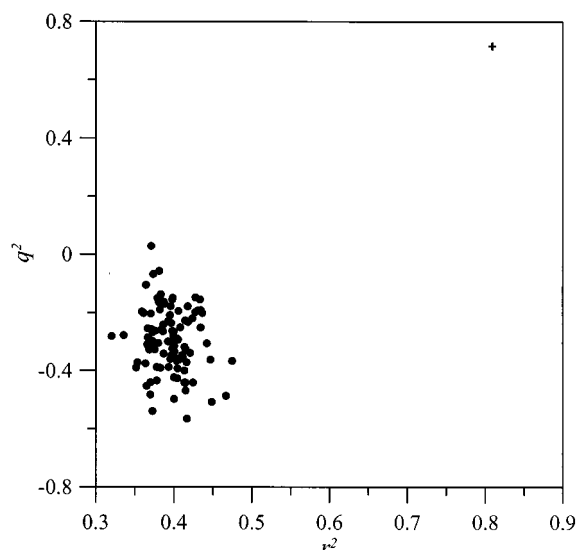


Figure 3. Random test results after permuting the boiling point values for a set of 529 saturated hydrocarbons 100 times. The real value is labeled with a cross.

are 0.395 and -0.293 . The results of the previously published work⁴⁰ are slightly better, yielding to $r^2 = 0.961$ using 6 topological indices but providing no information about the predictive capacity of the proposed models.

Gas Chromatographic Retention Times for 152 Diverse Chemical Compounds. The identification of single organic compounds in chemical mixtures can be made by means of comparing chromatographic peaks with samples of each compound. Due to the costs, or simply by unavailability, of pure samples, sometimes this comparison cannot be done. In this way, it is interesting to develop QSPR models to describe and to predict retention parameters from molecular structures. The present example deals with retention times (t_R), namely the mean time needed for a molecule to elute from a chromatographic column. A molecular set composed by 152 structurally different compounds has been studied together with their retention times, presented in Table 2, with the expressed procedure. The original 152×152 descriptor matrix was reduced to 152×135 and submitted to the PLS routine. The optimal model was found to be built up from five PLS factors, and the obtained results are

$$t_R = 75.83f_1 + 345.17f_2 + 869.53f_3 + 345.55f_4 + 492.72f_5$$

$$r^2 = 0.753 \quad q^2 = 0.541 \quad s = 1.657$$

As seen from the statistical results, an acceptable relationship with valuable predictive capacity is obtained. The predictive power is also graphically represented in Figure 4, where experimental versus predicted t_R values are plotted.

As seen in Figure 4, all compounds are accurately predicted in a narrow margin with the exception of a few molecules, which are mostly overestimated. The results of the random test for this molecular set are, similarly to the previous examples, very satisfactory, yielding to average r^2 and q^2 values of 0.229 and -0.234 , whereas the maximum values achieved were 0.373 and 0.180, respectively. The previous work⁴¹ obtained better results using a large set of 296 classical, topological, and quantum parameters computed

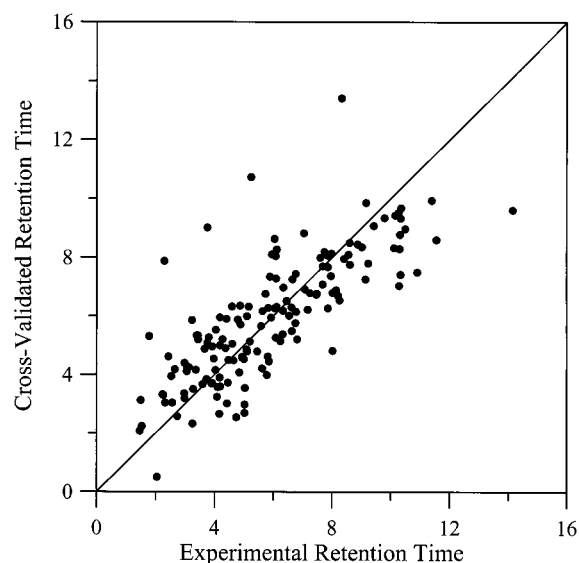


Figure 4. Experimental versus predicted retention times for a set of 152 compounds.

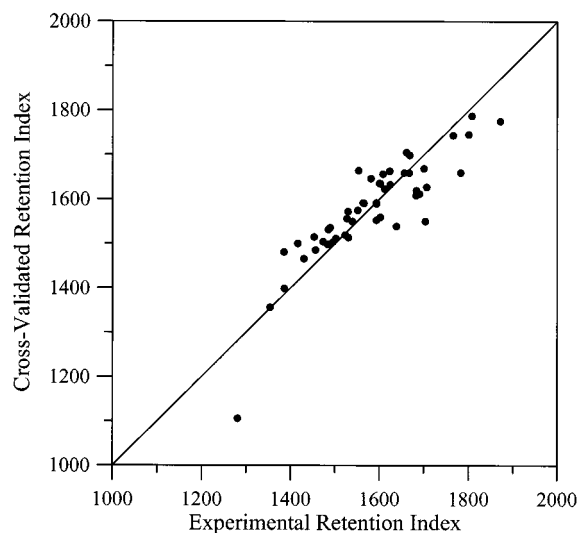


Figure 5. Experimental versus predicted retention indices for a set of 50 phenol derivatives.

with the CODESSA⁴⁵ software. Their models were calculated using linear and nonlinear multilinear regressions yielding to cross-validation regression coefficients (r^2_{cv}) of 0.961 and 0.967, respectively. Excellent results in the random tests were also reported.

Gas-Liquid Chromatography Retention Indices for 50 Phenol Derivatives. This example, similarly to the previous one, deals with chromatography. However, instead of t_R , Kováts retention indices (RI)⁴⁶ are studied. The analyzed molecular set consists of 50 phenol derivatives. The molecular structures, as well as the RI, are presented in Table 3. The original 50×50 descriptor matrix was reduced to 50×38 after the redundancies elimination, and the optimal QSPR model, built up from 6 PLS factors, results in

$$RI = 12938.48f_1 + 31183.31f_2 + 54224.34f_3 + 60530.94f_4 + 25264.86f_5 + 45970.82f_6$$

$$r^2 = 0.897 \quad q^2 = 0.762 \quad s = 59.294$$

Table 2. SMILES Notation^{47,48} and Retention Times of 152 Diverse Chemical Compounds⁴¹

compd	retention time	compd	retention time	compd	retention time
CCCCCCCCCCCC	6.34	c ₁ cm ₃ cc(c ₁)Oc ₂ cm ₃ ccc ₂	7.74	c ₁ cm ₃ cc(O)c ₁ OC(C)C	5.95
CCCCCCCCCCCC CCC	8.56	c ₁ cm ₃ cc(c ₁)C#N	4.1	c ₁ cnccc ₁ C#N	3.89
CCCCCCC=C	2.53	c ₁ cm ₃ cc(c ₁)Sc ₂ cm ₃ ccc ₂	9.02	c ₁ (CC)ccnc ₁	3.92
CCCCCCCCC=C	4.45	c ₁₂ CCCCc ₁ c(O)ccc ₂	8.02	c ₁ c(C)cc(C)nc ₁ C	4.37
C ₁ CCC(O)CC ₁	3.27	c ₁ cm ₃ ccc ₁ C(=O)C ₂ CCCCC ₂	9.14	c ₁ c(C)cc(C)nc ₁ N	5.85
C ₁ CCC(=O)CC ₁	3.25	c ₁ cm ₃ cc(c ₁)NC	4.94	c ₁ cm ₃ cc(c ₁)C=CC=O	6.63
CCCCCCCCC	4.88	c ₁₂ cm ₃ cc(O)c ₂ cm ₃ c(c ₁)O	10.33	c ₁₂ cm ₃ ccc ₁ OC(=O)CC ₂	7.46
CCCCCCCCCCN	6.62	C ₁ =CC=CS ₁	1.49	c ₁ cm ₃ cc(c ₁)C=CC(=O)O	7.86
CCCCCC#N	3.06	C ₁ =CC=CN ₁	2.04	c ₁₂ cm ₃ ccc ₁ OC(=O)C=C ₂	7.87
C ₁ CCC(C ₁)=O	2.26	c ₁ (CN)ccccc ₁	4.47	c ₁ (CC)ccncc ₁	3.97
CCCCCCCCSCCCCCC	10.9	c ₁ (C)ccc(cc ₁)n(o)o	6.24	c ₁ (CCN)ccncc ₁	6.08
CCOC(=O)OCC	2.29	c ₁ cm ₃ c(cc ₁)C(=O)Oc ₂ cm ₃ ccc ₂	9.42	c ₁ CCN(CC ₁)C(=O)OCC	6.11
C ₁ CCC(C ₁)O	2.32	c ₁ (C)ccc(N)cc ₁	5.02	c ₁ cm ₃ nc(c ₁)CC	3.44
CCCCCCCCCO	5.1	c ₁₂ cm ₃ ccc ₁ Cc ₃ c ₂ cm ₃ cc ₃	10.28	c ₁ cm ₃ nc(c ₁)CCO	5.63
CCCCCCCCCCCO	6.77	c ₁ (O)ccc(cc ₁)CC	5.81	c ₁ cm ₃ c(c ₁)C(=O)C(=O)c ₂ cm ₃ ccc ₂	10.31
CCCCCCCCCCCCO	8.26	c ₁ (N)ccccc ₁	4.08	c ₁ (cnccc ₁)C(=O)OCC	6.04
C ₁ CCCC(N)CC ₁	4.4	c ₁ (N)ccccc ₁ C	5	c ₁ cm ₃ nc(c ₁)C(C)=O	4.6
C ₁ CCC(N)CC ₁	3.05	c ₁ cm ₃ ccc ₁ OCCCC	6.08	c ₁ cncc(c ₁)C(C)=O	5.2
CCCCNCCCC	4.18	c ₁ cm ₃ ccc ₁ Br	3.75	c ₁ cm ₃ c ₂ CCOCc ₂ c ₁	6.04
CCCCCCCCCCCCCN	8.14	c ₁ cm ₃ ccc ₁ I	4.74	C ₁ (CN)=CC=CO ₁	2.73
CCCCCCCCCCBr	7.42	c ₁ cm ₃ ccc ₁ C=O	3.92	C ₁ =CC=CN(C ₁)C=O	5.78
N ₁ (C)CCNCC ₁	2.98	c ₁ cm ₃ ccc ₁ CO	4.64	c ₁ cm ₃ c ₂ c(c ₁)cnc ₃ cm ₃ ccc ₂₃	10.49
C(O)CCO	2.65	c ₁ cm ₃ (C)ccc ₁ C=O	5.1	c ₁ cm ₃ c(CC)c(c ₁)O	5.58
C(O)CCCO	3.72	C ₁ (=O)C=CC(=O)c ₂ c ₁ cm ₃ cc ₂	7.68	c ₁ cm ₃ c(cnc ₁)C=O	4.17
C ₁ CCC(C ₁)CC ₁	3.36	c ₁₂ cm ₃ ccc ₁ C=CN ₂	6.81	c ₁ cm ₃ c ₂ cm ₃₃ cm ₃ ccc ₃ nc ₂ c ₁	10.35
C ₁ CCCC(=O)CC ₁	4.42	c ₁ (ccccc ₁)n(o)o	5.09	c ₁ cm ₃ c(cc ₁)Nc ₂ cm ₃ ccc ₂	9.16
C ₁ CCC(CC ₁)C(=O)O	5.62	c ₁₂ cm ₃ ccc ₁ N=CS ₂	6.33	c ₁ (C)ccc(O)cc ₁	5.02
c ₁ cm ₃ c ₂ cm ₃ ccc ₁₂	6.09	c ₁₂ cm ₃ ccc ₁ NC=N ₂	8.02	c ₁ (C)ccccc ₁	2.56
c ₁ cm ₃ cc(c ₁)CCc ₂ cm ₃ ccc ₂	7.6	c ₁₂ cm ₃ ccc ₁ cm ₃ c(n ₂)C	7.04	c ₁ cm ₃ cc(c ₁)C ₂ cm ₃ ccc ₂	7.96
c ₁ cm ₃ cc(c ₁)CCC	4.04	c ₁ (C)ncccc ₁ C	3.8	c ₁ cm ₃ cc(c ₁)CCc ₂ cm ₃ ccc ₂	8.61
c ₁ cm ₃ cc(c ₁)CC=C	3.39	c ₁ (C)cccc(n ₁)C	3.23	c ₁ cm ₃ c(cc ₁)CC(C)(C)C ₂ cm ₃ ccc ₂	9.79
c ₁ cm ₃ c ₂ CCCc ₁₂	4.79	c ₁ (ccccc ₁)SSc ₂ cm ₃ ccc ₂	10.15	c ₁ cm ₃ (C)ccc ₁ Cl	4.03
c ₁ cm ₃ cc(c ₁)C(C)C	3.75	c ₁ (C)ncccc(c ₁)C	3.66	c ₁ cm ₃ cc(c ₁)CBr	5.24
c ₁ cm ₃ c ₂ c ₁ cm ₃ c ₃ cm ₃ ccc ₂₃	10.27	c ₁ (C)ccccc ₁ O	4.84	CC(C)C(=O)C(C)C	2.43
c ₁ cm ₃ c ₂ CCCCc ₂ c ₁	5.92	c ₁ (OC)cc(OC)cc(c ₁)OC	7.68	c ₁ cm ₃ ccc ₁	1.46
c ₁ cm ₃ c ₂ CCCNc ₂ c ₁	7.17	c ₁ (O)ccc(cc ₁)C(C)C	6.32	c ₁ cm ₃ ccc ₁ C	2.23
c ₁ cm ₃ c ₂ cm ₃ nc ₂ c ₁	6.45	c ₁ cm ₃ cc(c ₁)C(=O)c ₂ cm ₃ ccc ₂	9.23	c ₁ cm ₃ cc(c ₁)C ₁ C	3.43
c ₁ c(C)cc(C)ccc ₁ C	4.18	CC(=O)c ₁ cm ₃ ccc ₁	4.93	C(CC)CCCCCCC	4.59
c ₁ cm ₃ c ₂ CNCCc ₂ c ₁	6.75	CCCCCCCCC(=O)O	6.76	C(CC)CCCCCCCCCCCCCCCC	11.55
c ₁ cm ₃ c ₂ cm ₃ cm ₃ ccc ₃ cm ₃₂ c ₁	10.34	c ₁ cm ₃ ccc ₁ OC(=O)C	5.73	C ₁ CCCCC ₁	1.53
c ₁ cm ₃ c(cc ₁)C(=C)C	4.27	c ₁ (O)cccc ₁ C(C)C	6.11	CCCCCCC	1.78
c ₁ cm ₃ c ₂ c(c ₁)Oc ₃ cm ₃ ccc ₂₃	8.6	C(C#N)c ₁ cm ₃ ccc ₁	5.45	c ₁ cm ₃ ccc ₁ CC	3.12
c ₁ cm ₃ c ₂ c(c ₁)Sc ₃ cm ₃ ccc ₂₃	10.11	C(C(=O)O)c ₁ cm ₃ ccc ₁	6.54	c ₁ cm ₃ ccc ₁ C(=O)CC	5.82
c ₁ (ccccc ₂ cm ₃ ccc ₁₂)c ₃ cm ₃ cc ₄ cm ₃ ccc ₃₄	14.14	c ₁ cm ₃ cc(c ₁)CCCC(=O)O	7.86	c ₁ cm ₃ ccc ₁ OC	3.59
c ₁ cm ₃ c ₂ CCCC(=O)c ₁₂	7.48	c ₁ cm ₃ cc(c ₁)SC	5.16	c ₁ (O)ccc ₂ cm ₃ cc(O)c ₂ c ₁	10.3
c ₁ cm ₃ c ₂ cm ₃ cc(Br)c ₁₂	8.33	c ₁ cm ₃ cc(c ₁)CSC	5.88	c ₁ cm ₃ c ₂ cm ₃ nc ₂ c ₁	6.65
c ₁ cm ₃ c ₂ cm ₃ cc(O)c ₁₂	8.41	c ₁ ncccc(c ₁)C	2.98	c ₁ cm ₃ cn ₁ c ₂ ncccc ₂	7.98
c ₁ cm ₃ ccc ₁ C ₂ CCCCC ₂	7.24	c ₁ ncccc ₁ C	2.97	C ₁ (C)cc(C)cc(C)c ₁ C(O)=O	7.82
c ₁ cm ₃ c ₂ cm ₃ cc(Cc ₃ cm ₃ ccc ₃)c ₂ c ₁	11.39	c ₁ c(C)cccc ₁ O	5.03	c ₁ cm ₃ ccc ₁ S(C)=O	7.07
c ₁ cm ₃ cc(c ₁)O	4.16	c ₁ (C)c(C)c(C)c(C)c(C)c ₁ C	8.21	c ₁ cm ₃ ccc ₁ CCC=C	4.87
c ₁ cm ₃ cc(c ₁)OCc ₂ cm ₃ ccc ₂	8.87	n ₁ cm ₃ cc(c ₁)C#N	4.19		

As seen from the previous statistical results, a satisfactory QSPR model is found with a remarkable predictive capacity, as judged from the q^2 value achieved. A clearer picture of the obtained results can be seen by browsing Figure 5, where the experimental RI are plotted versus the cross-validated ones. As seen in Figure 5, all RI are correctly predicted except the one referred to the unsubstituted benzene, which is underestimated. In addition, the random test is satisfactorily surpassed, yielding to averaged values of 0.437 and -0.352 for r^2 and q^2 , respectively, whereas the maximum values were 0.675 and 0.377.

The previous work⁴⁶ achieved valuable results with a few number of molecular descriptors, consisting of vertex- and edge-weighted molecular graphs. Results ranging from 0.978 to 0.982 for r^2_{cv} values are obtained with up to 4 topological descriptors.

CONCLUSIONS

A number of molecular sets have been tested for correlation search with molecular properties, such as boiling points and chromatographic parameters. Satisfactory correlations are obtained using a small number of descriptors for molecular sets making use of quantum similarity measures.

Even if other methodologies may provide better results, it must be emphasized that the procedure used along this work, as well as the generation of molecular descriptors, has been kept unmanipulated. In other words, the exposed QSPR protocol consists of a general methodology made of unbiased and universal MQSM descriptors able to characterize different molecular properties without introducing further information. Additional refinements may be applied to the procedure to improve the results according to each

Table 3. Smiles Notation^{47,48} and Retention Indices of 50 Phenol Derivatives⁴²

compd	retention index	compd	retention index
c ₁ (O)cccc ₁ C	1281	c ₁ (O)c(C)ccc(C)c ₁ C	1551
c ₁ (O)cccc ₁ C	1354	c ₁ (O)cc(C)c(C)cc ₁ C	1593
c ₁ (O)cccc(c ₁)C	1386	c ₁ (O)cc(C)c(C)c(c ₁)C	1667
c ₁ (O)ccc(C)cc ₁	1385	c ₁ (O)ccc(cc ₁)CC(C)C	1612
c ₁ (O)cccc ₁ CC	1430	c ₁ (O)cccc ₁ CCCC	1600
c ₁ (O)cccc(c ₁)CC	1483	c ₁ (O)cccc(c ₁)CCCC	1668
c ₁ (O)ccc(cc ₁)CC	1473	c ₁ (O)ccc(cc ₁)CCCC	1661
c ₁ (O)cccc(C)c ₁ C	1495	c ₁ (O)ccc(cc ₁)CCC	1623
c ₁ (O)ccc(C)cc ₁ C	1456	c ₁ (O)c(C)cccc ₁ CCC	1553
c ₁ (O)cc(C)ccc ₁ C	1453	c ₁ (O)cc(C)ccc ₁ CCC	1602
c ₁ (O)c(C)cccc ₁ C	1416	c ₁ (O)ccc(C)cc ₁ CCC	1593
c ₁ (O)cc(C)cc(c ₁)C	1489	c ₁ (O)ccc(CC)cc ₁ CC	1602
c ₁ (O)ccc(C)c(c ₁)C	1530	c ₁ (O)cc(CC)ccc ₁ CC	1624
c ₁ (O)ccc(cc ₁)C(C)C	1527	c ₁ (O)ccc(CC)c(c ₁)CC	1682
c ₁ (O)cccc ₁ CCC	1502	c ₁ (O)cc(C)c(C)c(c ₁)C	1782
c ₁ (O)cccc(c ₁)CCC	1565	c ₁ (O)c(C)cc(C)c(c ₁)C	1690
c ₁ (O)ccc(cc ₁)CCC	1563	c ₁ (O)c(C)c(C)cc(C)c ₁ C	1683
c ₁ (O)ccc(C)cc ₁ CC	1523	c ₁ (O)cc(C)c(C)cc ₁ CC	1656
c ₁ (O)cc(C)ccc ₁ CC	1529	c ₁ (O)cccc ₁ CCCC	1700
c ₁ (O)c(C)cccc ₁ CC	1485	c ₁ (O)ccc(cc ₁)CCCC	1765
c ₁ (O)cc(C)cc(c ₁)CC	1581	c ₁ (O)ccc(cc ₁)CC(C)(C)C	1703
c ₁ (O)ccc(CC)cc ₁ C	1539	c ₁ (O)c(CC)cccc ₁ CCC	1706
c ₁ (O)ccc(CC)c(c ₁)C	1608	c ₁ (O)cccc ₁ CCCCC	1800
c ₁ (O)ccc(C)c(C)c ₁ C	1638	c ₁ (O)ccc(cc ₁)CCCCC	1871
c ₁ (O)cc(C)cc(C)c ₁ C	1593	c ₁ (O)cc(CCCC)ccc ₁ CC	1807

molecular set under study; however, the exposed protocol may constitute an excellent starting point for subsequent research.

ACKNOWLEDGMENT

This research has been partially supported by Project No. SAF2000-0223-C03-01 from the Spanish Ministerio de Ciencia y Tecnología. X.G. benefits from a predoctoral fellowship from the University of Girona. Financial support from the Fundació Maria Francisca de Roviralta is also acknowledged.

REFERENCES AND NOTES

- Quoted in: Borman, S. New QSAR Techniques Eyed for Environmental Assessments. *Chem. Eng. News* **1990**, 68, 20–23.
- Hammett, L. P. The Effect of Structures upon the Reactions of Organic Compounds. Benzene Derivatives. *J. Am. Chem. Soc.* **1937**, 59, 96–103.
- Jurs, P. C. Quantitative Structure–Property Relationships (QSPR). In *Encyclopedia of Computational Chemistry*; Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III; Schreiner, P. R., Eds.; John Wiley and Sons Ltd.: Chichester, U.K., 1998; Vol. 4, pp 2320–2330.
- Waterbeemd, H. v. d., Ed. *Structure–Property Correlations in Drug Research*; Academic R. G. Landes Co.: Austin, TX, 1996.
- Kubinyi, H. Quantitative Structure–Activity Relationships in Drug Design. In *Encyclopedia of Computational Chemistry*; Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III; Schreiner, P. R., Eds.; John Wiley and Sons Ltd.: Chichester, U.K., 1998; Vol. 4, pp 2309–2319.
- Carbó, R.; Arnau, J.; Leyda, L. How similar is a molecule to another? An electron density measure of similarity between two molecular structures. *Int. J. Quantum Chem.* **1980**, 17, 1185–1189.
- Carbó-Dorca, R.; Besalú, E. A general survey of molecular quantum similarity. *THEOCHEM* **1998**, 451, 11–23.
- Carbó-Dorca, R.; Amat, L.; Besalú, E.; Lobato, M. Quantum Similarity. In *Advances in Molecular Similarity*; Carbó-Dorca, R., Mezey, P. G., Eds.; JAI Press: Greenwich, CT, 1998; Vol. 2, pp 1–42.
- Carbó-Dorca, R. Fuzzy sets and Boolean tagged sets; vector semispaces and convex sets; quantum similarity measures and ASA density functions; diagonal vector spaces and quantum chemistry. In *Advances in molecular similarity*; Carbó-Dorca, R., Mezey, P. G., Eds.; JAI Press: Greenwich, CT, 1998; Vol. 2, pp 43–72.
- Besalú, E.; Carbó, R.; Mestres, J.; Solà, M. Foundations and recent developments on molecular quantum similarity. *Top. Curr. Chem.* **1995**, 173, 31–62.
- Carbó, R.; Besalú, E. Theoretical foundations of quantum molecular similarity. In *Molecular similarity and reactivity: from quantum chemical to phenomenological approaches*; Carbo, R., Ed.; Kluwer: Amsterdam, 1995; pp 3–30.
- Carbó-Dorca R.; Besalú, E.; Amat, L.; Fradera, X. Quantum molecular similarity measures: concepts, definitions and applications to quantitative structure-properties relationships. In *Advances in molecular similarity*; Carbó-Dorca, R., Mezey, P. G., Eds.; JAI Press: Greenwich, CT, 1996; Vol. 1, pp 1–42.
- Carbó, R.; Besalú, E.; Amat, L.; Fradera, X. Quantum molecular similarity measures (QMSM) as a natural way leading towards a theoretical foundation of quantitative structure-properties relationships. *J. Math. Chem.* **1995**, 18, 237–246.
- Fradera, X.; Amat, L.; Besalú, E.; Carbó-Dorca, R. Application of molecular quantum similarity to QSAR. *Quant. Struct.-Act. Relat.* **1997**, 16, 25–32.
- Lobato, M.; Amat, L.; Carbó-Dorca, R. Structure–activity relationships of a steroid family using quantum similarity measures and topological quantum similarity indices. *Quant. Struct.-Act. Relat.* **1997**, 16, 465–472.
- Amat, L.; Robert, D.; Besalú, E.; Carbó-Dorca, R. Molecular quantum similarity measures tuned 3D QSAR: an antitumoral family validation study. *J. Chem. Inf. Comput. Sci.* **1998**, 38, 624–631.
- Robert, D.; Amat, L.; Carbó-Dorca, R. Three-dimensional quantitative structure–activity relationships from tuned molecular quantum similarity measures. Prediction of the corticosteroid-binding globulin binding affinity for a steroid family. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 333–344.
- Robert, D.; Gironés, X.; Carbó-Dorca, R. Facet diagrams for quantum similarity data. *J. Comput.-Aided Mol. Des.* **1999**, 13, 597–610.
- Gironés, X.; Amat, L.; Robert, D.; Carbó-Dorca, R. Use of electron–electron repulsion energy as a molecular descriptor in QSAR and QSPR studies. *J. Comput.-Aided Mol. Des.* **2000**, 14, 477–485.
- Robert, D.; Carbó-Dorca, R., Aromatic compounds aquatic toxicity QSAR using quantum similarity measures. *SAR QSAR Environ. Res.* **1999**, 10, 401–422.
- Gironés, X.; Amat, L.; Carbó-Dorca, R. Using molecular quantum similarity measures as descriptors in quantitative structure-toxicity relationships. *SAR QSAR Environ. Res.* **1999**, 10, 545–556.
- Dean, P. M. Molecular Similarity. In *3D QSAR in Drug Design: Theory, Methods and Applications*; Kubinyi, H., Ed.; ESCOM Science Publishers BV.: Leiden, The Netherlands, 1993; pp 150–172.
- Richards, W. G., Molecular Similarity and Dissimilarity. In *Modeling of Biomolecular Structures and Mechanisms*; Pullman, A., Jortner, J., Pullman, B., Eds.; Kluwer Academic Publishers: Dordrecht, The Netherlands, 1995; pp 365–369.
- Gironés, X.; Robert, D.; Carbó-Dorca, R. TGSA: a molecular superposition program based on Topo-Geometrical Considerations. *J. Comput. Chem.* **2001**, 22, 255–263.
- Carbó-Dorca, R. Stochastic Transformation of Quantum Similarity Matrixes and Their Use in Quantum QSAR (QQSAR) Models. *Int. J. Quantum Chem.* **2000**, 79, 163–177.
- Constants, P.; Carbó, R. Atomic shell approximation: electron density fitting algorithm restricting coefficients to positive values. *J. Chem. Inf. Comput. Sci.* **1995**, 35, 1046–1053.
- Amat, L.; Carbó-Dorca, R. Quantum similarity measures under atomic shell approximation: first-order density fitting using elementary Jacobi rotations. *J. Comput. Chem.* **1997**, 18, 2023–2039.
- Amat, L.; Carbó-Dorca, R. Fitted electronic density functions from H to Rn for use in quantum similarity measures: cis-diamminedichloro-platinum(II) complex as an application example. *J. Comput. Chem.* **1999**, 20, 911–920.
- AMPAC 6.55; Semichem: Shawnee, KS, 1999.
- Dewar, M. J. S.; Zoebisch, E. G.; Healy, E. F.; Stewart, J. J. P. AM1: A New General Purpose Quantum Chemical Molecular Model. *J. Am. Chem. Soc.* **1985**, 107, 3902–3909.
- Höskuldsson A. *Prediction Methods in Science and Technology*; Thor: Copenhagen, 1996.
- Tenenhaus, M. *Regression de PLS*; Editions Technip: Paris, 1997.
- Wold, S.; Johansson, E.; Cocchi, M. PLS–Partial Least-Squares Projections to Latent Structures. In *3D QSAR in Drug Design, Theory, Methods and Applications*; Kubinyi, H., Ed.; ESCOM Science: Leiden, The Netherlands, 1993; pp 523–550.
- Wold, S. PLS for multivariate linear modelling. In *Methods and Principles in Medicinal Chemistry. Vol. 2: Chemometric Methods in Molecular Design*; van de Waterbeemd, H., Ed.; VCH: Weinheim, Germany, 1995; pp 195–218.

- (35) Cramer, R. D., III; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect on Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (36) Wold, S.; Sjöström, M.; Eriksson, L. Partial Least Squares Projections to Latent Structures (PLS) in Chemistry. In *Encyclopedia of Computational Chemistry*; Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer, H. F., III., Sreiner, P. R., Eds.; John Wiley and Sons Ltd.: Chichester, U.K., 1994; Vol. 4, pp 2006–2021.
- (37) Geladi, P.; Kowalski, B. R. Partial Least-Squares Regression: A Tutorial. *Anal. Chim. Acta* **1986**, *185*, 1–17.
- (38) Montgomery, D. C.; Peck, E. A. *Introduction to linear regression analysis*; Wiley: New York, 1992.
- (39) Wold, S. Cross-validatory estimation of a number of components in factor and principal component models. *Technometrics* **1978**, *20*, 397–405.
- (40) Rücker, G.; Rücker, C. On Topological Indices, Boiling Points, and Cycloalkanes. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 788–802.
- (41) Lucic, B.; Trinajstić, N.; Sild, S.; Karelson, M.; Katritzky, A. R. A New Efficient Approach for Variable Selection Based on Multi-regression: Prediction of Gas Chromatographic Retention Times and Response Factors. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 610–621.
- (42) Ivanciuc, O.; Ivanciuc, T.; Cabrol-Bass, D.; Balaban, A. T. Comparison of Weighting Schemes for Molecular Graph Descriptors: Application in Quantitative Structure-Retention Relationship Models for Alkylphenols in Gas–Liquid Chromatography. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 732–743.
- (43) Katritzky, A. R.; Lobanov, V. S.; Karelson, M. Normal Boiling Points for Organic Compounds: Correlation and Prediction by a Quantitative Structure–Property Relationship. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 28–41.
- (44) Ivanciuc, T.; Balaban, A. T. Quantitative Structure–Property Relationship Study of Normal Boiling Points for Halogen-/Oxygen-/Sulphur-Containing Organic Compounds Using the CODESSA Program. *Tetrahedron* **1998**, *54*, 9129–9142.
- (45) Katritzky, A. R.; Ignatchenko, E. S.; Barcock, R. A.; Lobanov, V. S.; Karelson, M. Prediction of Gas Chromatographic Retention Times and Response Factors Using a General Quantitative Structure–Property Relationship Treatment. *Anal. Chem.* **1994**, *66*, 1799–1807.
- (46) Buryan, P.; Nabivach, V. M.; Dimitrikov, V. P. Structure-Retention Correlations of Isomeric Alkylphenols in Gas–Liquid Chromatography. *J. Chromatogr.* **1990**, *509*, 3–14.
- (47) 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.
- (48) SMILES stands for simplified molecular input line entry system. More information can be found at the following: <http://esc.syrres.com/interkow/docsmile.html>.

CI0103370