# General topological patterns of known drugs

J. Gálvez [a,*], J.V. de Julián-Ortiz [b], R. García-Domenech [a]

[a] *Unidad de Investigación de Diseño de Fármaos y Conectividad Molecular, Dept. Química Fisica, Facultad de Farmacia,
Universitat de València, Valencia, Spain*
[b] *JVDJO Investigación de nuevos productos, Valencia, Spain*

## Abstract

Discriminating "drug-like" from "non-drug-like" compounds is a relatively emerging topic within the drug research. The basic assumption is that it is possible to obtain relevant information from structural features common to the known drugs, in order to discard a huge number of candidate chemical structures with low probability of becoming drugs. A graph-theoretical contribution to this subject is reported in this paper, by making exclusive use of linear relationships. The results suggest that it is possible to achieve a pattern of general pharmacological activity based on molecular topology. Conclusions are tentative pending verification of the results with larger compound libraries. © 2001 Elsevier Science Inc. All rights reserved.

*Keywords:* Topology; Compounds; Drugs

## 1. Introduction

The philosophy of the applicability of non-mechanistical SAR approaches has experienced an evolution. The success reached in each step has, unevenly, prompted the following challenges, QSAR models have been widely attempted using *homologous* and *congeneric* series of active compounds for the correlation of biological properties [1]. On the other hand, studies on structurally *heterogeneous* sets of compounds, having similar mechanism but different parent structure, are not very common [2]. Some examples of activity modeling into mechanistically *diverse* groups of molecules, mainly correlations of toxicity, can also be found [3]. Finally, the most *general* approach, the attempt to discriminate between drug and non-drug molecules irrespective to the activity they present, has arisen. This methodology has been proposed for the design of libraries in combinatorial chemistry as well as molecular structure databases for virtual screening [4].

More and more complex mathematical tools are used in the modeling of the activity patterns. But these techniques are not essential features. In all these extra-mechanistical approaches, the key is structural similarity. The level of subtlety used in the molecular structural patterns allows one to model the biological activities within series of homologous, congeneric, heterogeneous, diverse or general drugs. In these

five kinds of sets, arranged from lesser to greater chemical diversity, the accuracy of the results that the different mathematical models are able to reach, generally decreases in the same sequence also.

The interest in the identification of drug-like structures comes from the necessity of designing drugs that have desirable characteristics common with known drugs, such as synthetic accessibility at low cost, solubility, oral availability, useful pharmacokinetic and pharmacodynamic properties, chemical and metabolical stability, low toxicity, carcinogenesis and teratogenesis, minimum addictive potential, etc.

An important consideration must be regarded here, the definitions of drug and non-drug and their variations with the historical context. The concept of drug, and particularly the notion of non-drug, is not currently defined in a strict manner in the studies on general drugs discrimination. Probably this is not possible due to the evolution of the concept of drug. Indeed, the molecular structures of compounds considered as drugs have gradually changed. Two analyses, one with the drugs known today and another with the drugs used 50 years ago would provide different conclusions on what "drug-like" means, structurally speaking. The question is whether drugs used throughout the last decades would be recognized as such, according to their structural similarity with the more recently marketed substances. In other words, the question is if valuable candidates could be neglected because their structures are not *in fashion*. On assuming this risk, the effort to focus the drug discovery on the most likely candidates

* Corresponding author.
*E-mail address:* jorge.galvez@uv.es (J. Gálvez).

must continue, confident that theoretical improvements will take place in the future and they will be able to *rescue* other useful candidates previously disdained.

To sum up, if the compounds to be screened are "like drugs", the probability of success increases. Bemis and Murcko first published this working hypothesis in l996 [5]. They analyzed the frequencies of the molecular scaffolds present in 5120 known drugs from a commercially available database. A substructure was considered as a scaffold if it was constituted by rings or rings joined by molecular covalent bridges. All the other substructures were considered side chains. At most, a given molecule could have one scaffold. If the chemical nature of atoms and bonds was not considered, only 32 scaffolds were present in one-half of the database drugs. Following the same criterion, the total number of scaffolds was one-fourth of the number of drugs approximately, 783 (66%) scaffolds were found in only one molecule, and 306 (6%) compounds had no scaffold because, they did not present cycles. The most common structure was the six-membered ring with a frequency of 606 molecules. This analysis was repeated considering the chemical identity of atoms and bonds. Anyway, a low molecular diversity was found in the set of known drugs. No estimable differences between drugs and non-drugs were revealed in this preliminary work by using this analysis, but some strategies were proposed to apply it to the design of new drugs. For example, a pharmacological promiscuity index was defined for a given scaffold as the ratio between the number of pharmaceutical targets and the number of drugs. A value of this parameter close to one means that the basic structure could be used in the design of virtual libraries with probable success. Another strategy would be the assembly of "dumb" substructures with high frequency into a receptor considering only the steric factors, and then assign atom and bond types according to electrostatic complementarity. The same authors have released more recently a second paper with a detailed analysis of the common drug side chains [6].

Gillet et al. have introduced biological activity profiles and a genetic algorithm based scoring scheme to obtain rankings [7]. Two databases were analyzed, one of about 15,000 drugs and another containing a representative sample of ca. 17000 extracted from more than 1,500,000 presumed non-drugs. The descriptors used were, molecular weight (MW), Kier $^2\kappa_\alpha$ shape index [8], numbers of aromatic rings, rotatable bonds, hydrogen bond donors, and hydrogen bond acceptors. It is noteworthy that differences in the distribution diagrams for each descriptor between the two databases were found. A formalism similar to the substructural analysis [9] was used. Weights were assigned to the different values of the discrete descriptors as well as to arbitrary intervals of MW and $^2\kappa_\alpha$. Each given value of every descriptor was characterized through a weight which was calculated from the relative frequency of occurrences in known active and known inactive compounds, and it gave a measure of the likelihood of activity. A molecule is scored by summing the different weights that exhibits for each descriptor considered. The use

of more than a single descriptor often resulted in no increase of the predictive ability. A genetic algorithm showed better ability to combine discriminatory information from several descriptors. In the best result obtained, 360 drug molecules were found in the top 1000 ranked positions, within a subset of 1000 drugs and 16,807 non-drugs. About 50% of active compounds were identified within the top 17% of the total set. This methodology was also applied to specific therapeutic classes with diverse results, being the best discrimination the one reached with antibiotics.

Ajay et al. have used two descriptor sets, both, in a Bayesian neural network (BNN) and in a machine learning algorithm [10]. They have analyzed database subsets of 3500 compounds each from the drug and non-drug databases. The first group of descriptors was very similar to the one used by Gillet et al., these were, $\log P$, molecular weight, aromatic density, $^2\kappa_\alpha$ shape index, numbers of hydrogen bond donors, hydrogen bond acceptors, and rotatable bonds. The second group of descriptors was constituted by indicator variables, the ISIS fingerprint [11]. The two groups afforded reasonable accuracy for the classification into drug-like and non-drug-like molecules using BNN models, but combinations of both descriptor sets lowered the classification errors. One of the best models was constituted by a BNN with the seven descriptors from the first set and the 71 more relevant from the second one as inputs, having five hidden nodes and one output. The classification error percentages ranged from 10 to 12 for drugs and from 12 to 13 for non-drugs. A success of 80% was obtained in the extrapolation to a test database constituted by known drugs.

Sadowski and Kubinyi have developed a scoring model based on atom type descriptors and neural networks [12]. The Ghose–Crippen system for encoding organic molecules is based in counts of 120 atom types in a molecule [13]. Only 92 of such counts were considered relevant enough in this work, and were included as inputs for a neural network, with five hidden nodes and one output. Subsets of 5000 compounds were extracted from drug and non-drug databases, respectively. About 77% of the drugs and 83% of the non-drugs set were correctly classified. These results are comparable to those of Ajay et al. [10]. The extrapolating ability of the approach was tested by removing several entire therapeutic categories from the training set. It was shown that the model trained without these sets was only slightly worse in predicting such molecular structures as drug-like.

Wang and Ramnarayan have used the concept of the multilevel chemical compatibility between a compound and a drug library as a measure of the drug-like character of such compound [14]. This is based on the assumption that the local chemical environment of each atom or group of atoms in a compound contributes to its drug-likeness. A systematic comparison of the local environments within a compound and those within the existing drugs provides a similarity basis for determining such character. The method was applied to four test sets, top selling drugs, compounds under biological testing prior to the preclinical test, anticancer drugs, and

compounds known to have poor drug-like character. Among the conclusions obtained, it stood out that known drugs contain about 80% of all the viable types of local structure types; therefore, discovery of a drug with new local structures is relatively unlikely. The method was selective in discerning drug-like compounds, most of the top drugs, about one-quarter of the biological testing compounds, and about one-fifth of the anticancer drugs were drug-like, following the outlined criterion.

Frimurer et al. have developed a carefully designed neural network score approach using atom-type descriptors within non-redundant huge data sets, with a critical selection of training series [15]. With a model having 80 inputs and 200 hidden nodes, the best discrimination reported until now was obtained: 88% of accurate assignation both, in drug and non-drug sets. In order to avoid false negatives, the cut-off threshold can be moved and the model correctly predicts 98% of non-drugs and 63% of drugs. The model was also applied successfully to the drug-likelihood estimation of GABA uptake inhibitors.

Another useful approach has been the development of the "leadlike" concept by Oprea and coworkers as a tool for the optimization of combinatorial chemistry libraries [16]. The aim was the discrimination of potential leads more than drugs, since combinatorial synthesis projects are devoted to the identification of new lead drugs and their subsequent optimization. This is not essentially different from the discrimination between drugs and non-drugs. Every new drug that exhibits an unusual structure can be converted in a new *lead*. The distribution histograms of some descriptors were studied [17] for the same databases used by Ajay et al. [10] and Sadowski and Kubinyi [12]. About 70% of the drug-like compounds were found between the following limits, number of hydrogen bonds donors between 0 and 2, acceptors between 2 and 9, number of rotatable bonds between 2 and 8, number of rings between 1 and 4. About 60% of non-drugs and 30% of drugs showed two rings or less, and $<17$ rigid bonds whereas $<30\%$ of non-drugs and 60% of drugs were out of these figures. These findings suggest that the probability of identifying drug-like structures is not entirely independent of molecular *complexity* in the databases used. Thus, it would be valuable to discriminate drugs within sets of molecules with similar sizes. These sets would require a much more careful molecular selection than an extensive number of compounds.

Once it has been pointed out that the discrimination of general drugs is possible, it may be interesting to test linear models. In previous literature, there are models with very simple descriptors although with complex mathematical treatment. A graph-theoretical approach to the problem of discriminating drugs from non-drugs is attempted in this work, trying to avoid the bias due to different molecular size and complexity between the two groups. Once a molecular structure is recognized as a drug, more insight could be gained with the help of models able to discriminate a particular activity.

## 2. Descriptors and method

In Chemical Graph Theory, molecular structures are normally represented as hydrogen-depleted graphs, whose vertices and edges act as atoms and covalent bonds, respectively. Graph-theoretical indices, also known as topological indices, are descriptors that characterize a molecular graph and are able to give account of their structural properties. They have shown their usefulness in classification analysis [18–20], PLS [21,22] and, in general, in the modeling of biological activities [23–26].

Well-known descriptors were used in this analysis, subgraph Randić–Kier–Hall indices until fourth-order ($^m\chi_t$, $^m\chi_t^v$ [27,28], topological charge indices until fifth-order ($J_m$, $G_m$, $J_m^v$, $G_m^v$) [29,30], quotients of connectivity indices ($^mC_t = {}^m\chi_t/{}^m\chi_t^v$) Wiener path number (W) [31], $PR_n$ (number of pairs of ramifications at topological distance $n$, with $n$ ranging from 0 to 3), and other graph-theoretical descriptors which were not selected in the final models.

Stepwise linear discriminant analysis (LDA) was applied in order to obtain classification functions through linear combination of variables by using the BMDP software [32]. The maximum Snedecor F was the criterion for selecting the descriptors that appear in each equation. The quality of the classification functions obtained was evaluated through the Wilks parameter ($\lambda$) and the percentage of correct classification.

The models were obtained by using small but structurally heterogeneous sets of compounds, and it is assumed that the molecular diversity is high enough to have representative sets. Sets of known drugs and others of non-drugs were used in the obtaining of dicotomic models. As a quantification of the molecular diversity, the standard deviations of $N$ (number of non-hydrogen atoms) and $^1\chi$ (Randić index, that can be interpreted as a measure of the molecular branching) in the drugs set were used. In order to guarantee homogeneity in size and molecular complexity the mean values of $N$ and $^1\chi$ between the pairs of drug and non-drug sets were maintained close. As an initial criterion, only xenobiotic substances with known therapeutic utility were considered as drugs. All other molecules were considered as non-drugs. This criterion allowed a homogeneous complexity to be reached between drug and non-drug sets more easily.

LDA was also used to classify the drugs in five therapeutic groups.

## 3. Results and discussion

A first study was performed with 89 drugs and 73 presumed non-drugs (Table 1). The average values $\pm$ standard deviations obtained for the homogeneity measures were: $\langle N \rangle_d = 22.1 \pm 9.9$, $\langle {}^1\chi \rangle_d = 10.5 \pm 4.7$, and $\langle N \rangle_{nd} = 23.0 \pm 5.8$, $\langle {}^1\chi \rangle_{nd} = 10.8 \pm 5.8$ for drugs and non-drugs, respectively. Using all the molecules as a training set, the

Table 1
Classification results obtained with Eq. (2) by LDA

| Compound | Probability[a] | Classification[b] | Compound | Probability[a] | Classification[b] |
|---|---|---|---|---|---|
| Training group (drugs) | | | | | |
| Acifran. | 0.761 | D | Mechlorethamine | 0.342 | ND |
| Adrenaline | 0.744 | D | Mefenamic acid | 0.586 | D |
| Altretamine | 0.993 | D | Mercaptoacetamide | 0.539 | D |
| Amiloride | 0.326 | ND | Metofoline | 0.943 | D |
| Amsacrine | 0.431 | ND | Metrotexate | 0.526 | D |
| Atropine | 0.852 | D | Mevastatin | 0.747 | D |
| Azetazolamide | 0.376 | ND | Naproxen | 0.552 | D |
| Benzalbutyramide | 0.364 | ND | Nicoclonate | 0.679 | D |
| Bezofibrate | 0.655 | D | Nicomol | 0.914 | D |
| Bumetanide | 0.656 | D | Orciprenaline | 0.787 | D |
| Busulfan | 0.122 | ND | Oxiniacic acid | 0.429 | ND |
| Carboquone | 0.940 | D | Oxitropium | 0.931 | D |
| Carmustine | 0.681 | D | Paracetamol | 0.477 | ND |
| Chlorambucil | 0.725 | D | Pipobroman | 0.901 | D |
| Chlorthalidone | 0.550 | D | Piretanide | 0.772 | D |
| Citarabine | 0.755 | D | Pravastatin | 0.653 | D |
| Dacarbazine | 0.858 | D | Procarbazole | 0.542 | D |
| Ddticarbazine | 0.421 | ND | Quinethazone | 0.795 | D |
| Ephedrine | 0.534 | D | Salsalate | 0.525 | D |
| Esplacto | 0.943 | D | Sinvastatin | 0.714 | D |
| Etersalate | 0.882 | D | Sitosterol | 0.105 | ND |
| Ethacrynate | 0.750 | D | Soterenol | 0.590 | D |
| Etoposide | 0.485 | ND | Sulindac | 0.873 | D |
| Etozolin | 0.961 | D | Tenoxican | 0.863 | D |
| Fenofibrate | 0.610 | D | Theobromine | 0.886 | D |
| Fenoprofen | 0.541 | D | Thyrotropic acid | 0.628 | D |
| Fenoterole | 0.774 | D | Thyroxine | 0.681 | D |
| Fenspiride | 0.633 | D | Tiadenol | 0.268 | ND |
| Fluorouracil | 0.316 | ND | Tolmetin | 0.736 | D |
| Furosemide | 0.806 | D | Torasemide | 0.762 | D |
| Halofenate | 0.683 | D | Tretoquinol | 0.980 | D |
| Hidrochlorothiazide | 0.869 | D | Triamterene | 0.489 | ND |
| Indomethacin | 0.653 | D | Vinblastine | 0.908 | D |
| Ipratropium | 0.945 | D | Vincristine | 0.866 | D |
| Ketorolac | 0.586 | D | | | |
| Training group (non-drugs) | | | | | |
| 3,3′-Thiodipropionic acid | 0.440 | D | Indigo Carmine | 0.674 | ND |
| Acesulfame | 0.562 | ND | Isibutyl Isovalerate | 0.600 | ND |
| Acetic acid | 0.957 | ND | Isoamyl Isovalerate | 0.498 | D |
| Acid Green | 0.980 | ND | Isophytol | 0.842 | ND |
| Ad1066b | 0.810 | ND | Lactic acid | 0.816 | ND |
| Ad1320b | 0.597 | ND | Lactose | 0.583 | ND |
| Adipic acid | 0.663 | ND | Maltol | 0.495 | D |
| Aspartame | 0.330 | D | Naphthoresorcinol | 0.637 | ND |
| Benzalkonium | 0.600 | ND | Neohesperidin | 0.511 | ND |
| Benzoic acid | 0.870 | ND | Phytol | 0.803 | ND |
| Brilliant Blue | 0.980 | ND | Ponceau 3R | 0.728 | ND |
| Canthaxanthin | 0.937 | ND | Quinoline | 0.741 | ND |
| Capsathin | 0.971 | ND | Riboflavine | 0.408 | D |
| Carminic acid | 0.947 | ND | Sanguinarine | 0.303 | D |
| Citric acid | 0.510 | ND | Sorbic acid | 0.786 | ND |
| Clorophene | 0.654 | ND | Sorbitol | 0.703 | ND |
| Curcumin | 0.151 | D | Succinic acid | 0.828 | ND |
| Dianiside | 0.204 | D | Sucrose | 0.152 | D |
| Dibenzalacetone | 0.856 | ND | Sunset yellow | 0.795 | ND |
| Erythrocentaurin | 0.299 | D | Tartrazine | 0.735 | ND |
| Erythrosine | 0.714 | ND | Ursolic acid | 0.412 | D |
| Ethoxyquin | 0.066 | D | Vit. E Acet. | 0.538 | ND |
| Ethyl benzoate | 0.627 | ND | Vitamin E | 0.560 | ND |
| Flavoxanthin | 0.921 | ND | α-Carotene | 0.977 | ND |

Table 1 (*Continued*)

| Compound | Probability[a] | Classification[b] | Compound | Probability[a] | Classification[b] |
|---|---|---|---|---|---|
| Fumaric acid | 0.771 | ND | β-Carotene | 0.973 | ND |
| Glucosamine | 0.684 | ND | δ-Tocopherol | 0.535 | ND |
| Glutamic acid | 0.690 | ND | γ-Carotene | 0.972 | ND |
| Glycerol | 0.893 | ND | γ-Tocopherol | 0.555 | ND |
| *Test group (drugs)* | | | | | |
| AAS | 0.597 | D | Floctafenine | 0.573 | D |
| Bamifylline | 0.952 | D | Flutropium | 0.933 | D |
| Benfluzi | 0.572 | D | Hidroure | 0.064 | ND |
| Bevonium | 0.755 | D | Isonixin | 0.870 | D |
| Clofibrate | 0.840 | D | Metoxfen | 0.853 | D |
| Clomestrone | 0.838 | D | Piroxican | 0.809 | D |
| Clonitazene | 0.941 | D | Pirozadil | 0.990 | D |
| Diclofenac | 0.515 | D | Probucol | 0.868 | D |
| Doxofylline | 0.958 | D | Rimiterol | 0.664 | D |
| Etiroxate | 0.799 | D | Vindesin | 0.841 | D |
| *Test group (non-drugs)* | | | | | |
| Amaranth | 0.917 | ND | Malic acid | 0.693 | ND |
| Ascorbic acid | 0.327 | D | Ponceau SX | 0.615 | ND |
| Bixin | 0.858 | ND | Quinoline Y. | 0.741 | ND |
| Ciclamaet | 0.876 | ND | Rhodoxanthin | 0.960 | ND |
| Citrus Red | 0.160 | D | Ribixanthin | 0.969 | ND |
| Cryptoxanthin | 0.971 | ND | Saccharin | 0.589 | ND |
| EDTA | 0.151 | D | Tartaric acid | 0.634 | ND |
| Erythrosine | 0.912 | ND | β-Tocopherol | 0.535 | ND |
| Eucalyptol | 0.678 | ND | | | |

[a] Of classification.
[b] D: drug; ND: non-drug.

following classification function was obtained:

$$D_1 = -0.736\,{}^3\chi_c + 0.239G_1^v + 14.3J_3 + 29.4J_5^v$$
$$-0.00044W - 0.259\mathrm{PR}_3 - 3.58$$
$$N = 162, \ F = 13.46, \ \lambda = 0.657 \tag{1}$$

where ${}^3\chi_c$ and $\mathrm{PR}_3$ are the three-order cluster connectivity index and the number of pairs of vertices with topological valence = 3 (i.e. tertiary ramifications). Both encode pure structural information. The Wiener number, $W$, is the sum of the overall topological distances between every pair of vertices of the graphs. It is also structural and is related to the different possible ways of arranging atoms inside the molecule. Finally, the $G_i$ and $J_i$ indices are the topological charge indices, which evaluate intramolecular charge transfers between atoms placed at a topological distance $i = 1$, 3, 5. They take into account molecular properties, such as dipole moments and electronic polarizability.

Its classification success was 83.1 and 80.8% for drug and non-drug compounds, respectively.

The stability of the Eq. (1) was tested taking random test sets having 20% of the compounds, using the same six variables of Eq. (1) and performing new classification analyses. As it can be seen in Table 2, the classification results are stable. As an example, one of the functions obtained is given by the equation:

Table 2
Classification matrix obtained in the random test using the variables of the Eq. (1)

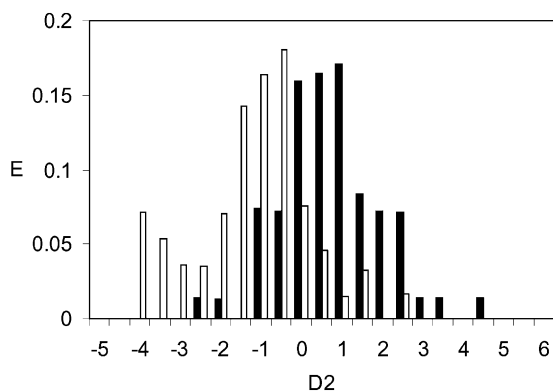| | Correct (%) | Drugs | Non-drugs |
|---|---|---|---|
| Run no. 1 | | | |
| Drugs | 80 | 55 | 14 |
| Non-drugs | 79 | 12 | 44 |
| Test drugs | 95 | 19 | 1 |
| Test non-drugs | 82 | 3 | 14 |
| Total | 84 | 89 | 73 |
| Run no. 2 | | | |
| Drugs | 85 | 58 | 10 |
| Non-drugs | 82 | 10 | 46 |
| Test drugs | 86 | 18 | 3 |
| Test non-drugs | 77 | 4 | 13 |
| Total | 82 | 90 | 72 |
| Run no. 3 | | | |
| Drugs | 82 | 61 | 13 |
| Non-drugs | 82 | 11 | 49 |
| Test drugs | 80 | 12 | 3 |
| Test non-drugs | 85 | 2 | 11 |
| Total | 82 | 86 | 76 |
| Run no. 4 | | | |
| Drugs | 81 | 61 | 14 |
| Non-drugs | 82 | 10 | 46 |
| Test drugs | 93 | 13 | 1 |
| Test non-drugs | 71 | 5 | 12 |
| Total | 82 | 89 | 73 |

Fig. 1. Pharmacological distribution diagram obtained for the test group by using the discriminant function $D_2$, Eq. (2). $E$, expectancy of belonging to a given set. Black bars, drugs set white bars, non-drugs set.

$$D_2 = -0.713^3\chi_c + 0.264G_1^v + 13.8J_3 + 23.3J_5^v$$
$$-0.00045W - 0.262PR_3 - 3.42$$
$$N = 125, \ F = 9.16, \ \lambda = 0.682 \tag{2}$$

The coefficients of the variables were very close to that obtained with Eq. (1). The probabilities of classification, as well as the categories assigned for each compound by Eq. (2) are displayed for the training and test groups in Table 1. This test set was also used to obtain the optimal discriminant interval of $D_2$ (and, by extrapolation, of $D_1$). Fig. 1 illustrates the pharmacological distribution diagram [33] with the expectancy profile for both drug and non-drug compounds for each interval of $D_2$. In general, expectancy for a set A along each interval $x$ is defined as, $E$: percentage of A in $x$/(percentage of non-A in $x + 100$). Despite the presence of an overlapping region, it is remarkable that these results are better than the usually achieved by other authors using more sophisticated models with more simple descriptors. Values of $D_2$ between 1 and 3.5 were taken as a threshold for the classification as drug. In this interval, the classification was considered optimal.

The standard deviations of $N$ and $^1\chi$, that were about one-half of the average value, were considered high enough to assess the molecular diversity of the drugs set and so, to guarantee its molecular heterogeneity. On the other hand, it is not important for these values to be high in the non-drugs set. The results obtained point out that Eq. (1) is stable and reproducible enough within the sets studied. In order to validate Eq. (1), it was applied to a set which included about 2000 compounds from the Merck index database. About 82% of accuracy in the classification of drugs is achieved for the $D_1$ range between 1 and 3.5. Table 3 illustrates the results obtained for a representative set of molecules.

It is noteworthy that despite the fact that some compounds were erroneously classified as non-drugs, their active–inactive gap is usually very short (for instance, triamterene shows probabilities for activity–inactivity of 48.9 versus 51.1%, respectively).

On the other hand, some of the compounds included into the non-drugs set, but possessing drug-desirable properties, such as vitamin E acetate or lactose, generally showed a non-negligible probability to be drugs (44 and 46.2%, respectively, Table 1). Other vitamins, such as riboflavin, included within the non-drug set, were classified as drugs (59.2% of probability). That also occurred with sucrose, which showed a probability of *activity* of about 85%.

Although all these results show the discriminant ability of the actual model, strictly speaking its significance must be still considered as tentative, pending verification of the model with a larger database.

As pointed out in the introduction, one arising question here is the notion of drug. After the restricted definition followed here, a drug is an exobiotic compound that shows a significant therapeutical utility, that is to say, an acceptable therapeutical index (low toxicity) as well as good bioavailability, etc. But what is a non-drug? By contrast, we could define it as a substance not showing the conditions expressed above. However, while we can be reasonably sure that a given compound is a drug, it is possible that another one not considered as such yet, could be demonstrated to be in the future. Thus, the absence of an entry in the Pharmacopeia is not evidence of absence of pharmacological action, The question about the drug nature of some *endobiotic* compounds included by us here in non-drugs sets is open to further discussion. Thus, in principle, any vitamin is considered as a natural product playing a physiological (not a pharmacological) role, despite the fact that some of them, such as vitamins E and C, may be used as antioxidant agents at high doses. The main reason to include them into the non-drug set was to obtain the desired size and structural complexity.

Once the active–inactive character has been discriminated, it is possible to classify it into a specific pharmacological activity. Along this goal, five different pharmacological activities, namely non-narcotic analgesics, bronchodilators, diuretics, hypolipidemics and antineoplastics were input into a new discriminant analysis. The results are displayed in Tables 4 and 5.

As may be realized, a significant level of success is achieved just including three topological indices ($J_1^v$, $^3C_c$, $^4C_{pc}$). More than 50% of average accuracy is obtained for all groups. This quality of the discrimination could be interpreted as the existence of a certain uniqueness in the mechanisms of action within every group. However, it does not work for antineoplastics, whose diverse mechanisms of action leads to a poor level of correct classification (just over 35%). As pointed out in the introduction, *hetrogeneous* sets give better models than *diverse* ones. Again, these results are pending of verification by using larger sets of compounds.

Certainly, an improvement in these results is to be expected in the future, for instance using more sophisticated techniques, such as neural networks, but the objective of this article was just to reveal the efficacy of the graph-theoretical approaches in the solutions of this kind of problems. It must

Table 3
Validation test of discriminant function $D_1$ for a representative set of compounds in a database comprised of 2000 compounds from the Merck index

| Compound | $D_1$[a] | Therapeutic category[b] | Compound | $D_1$[a] | Therapeutic category[b] |
|---|---|---|---|---|---|
| Acebutolol | 2.13 | Antihypertensive | Eugenol | 1.39 | Analgesic |
| Acecainide | 2.44 | Antiarrhythmic | Fenalcomine | 1.52 | Anesthesic |
| Aceclofenac | 1.51 | Analgesic | Fenarimol | 1.53 | Non-drug |
| Acenocoumarol | 1.12 | Anticoagulant | Fenbutrazate | 2.99 | Anorexic |
| Aceperone | 2.22 | Non-drug | Fentanyl | 2.45 | Analgesic |
| Acetazolamide | 1.85 | Diuretic | Fenthion | 1.99 | Anticholinergic |
| Acifran. | 1.68 | Hypolipidemic | Florfenicol | 2.26 | Antibacterial |
| Acipimox | 1.38 | Hypolipidemic | Flubendazole | 3.02 | Anthelmintic |
| Acitretin | 1.71 | Antipsoriatic | Flunisolide | 2.75 | Antiasmathic |
| Acrivastine | 1.24 | Antihistaminic | Formebolone | 1.1 | Anabolic |
| Actinoquinol | 1.84 | Non-drug | Fospirate | 2.46 | Anthelmintic |
| Afloqualone | 2.51 | Muscle relaxant | Furazolidone | 2.58 | Anti-infective |
| A-furildioxime | 1.49 | Non-drug | Gabexate | 2.44 | Non-drug |
| Alclofenac | 1.52 | Analgesic | Ganglefene | 1.46 | Vasodilatador |
| Amphotalide | 2.05 | Anthelmintic | Glucametacin | 2.26 | Anti-inflamatory |
| Aniracetam | 1.77 | Nootropic | Glutethimide | 2.05 | Sedative |
| Azaperone | 2.8 | Sedative | Glyburide | 2.1 | Antidiabetic |
| Azure c | 1.77 | Non-drug | Glyhexamide | 1.16 | Antidiabetic |
| Bamipine | 1.31 | Antihistaminic | Glymidine | 2.94 | Antidiabetic |
| Barban | 2.02 | Non-drug | Guanacline | 1.39 | Antihypertensive |
| Befunolol | 1.44 | Antiglaucoma | Halazepam | 2.02 | Anxiolytic |
| Bendazac | 2.81 | Anti-inflamatory | Halcinonide | 3.46 | Anti-inflamatory |
| Bendroflumethiazide | 1.77 | Diuretic | Halethazole | 2.08 | Antifungal |
| Benperidol | 3.28 | Antipsychotic | Heptachlor | 2.51 | Insecticide |
| Benzetimide | 2.87 | Anticholinergic | Heptenophos | 1.82 | Insecticide |
| Benzimidazole | 1.39 | Non-drug | Hexachlorobenzene | 1.18 | Non-drug |
| Benzimidazole | 1.39 | Non-drug | Hydracarbazine | 1.51 | Antihypertensive |
| Benznidazole | 2.55 | Antiprotozoal | Isoflupredone | 2.66 | Anti-inflamatory |
| Benzoctamine | 1.18 | Muscle relaxant | Isomethadone | 3.19 | Analgesic |
| Benzodepa | 2.27 | Antineoplastic | Lidocaine | 3.28 | Anesthetic |
| Betamethasone | 3.36 | Glucocorticoide | Linuron | 2.26 | Non-drug |
| Betazole | 1.75 | Non-drug | Lumazine | 2.65 | Non-drug |
| Bifluranol | 1.53 | Antiandrogen | Mebendazole | 2.65 | Anthelmintic |
| Binedaline | 1.51 | Antidepressant | Naloxone | 2.45 | Narcotic |
| Carbetapentane | 2.76 | Antitussive | Napropamide | 1.77 | Herbicide |
| Carbinoxamine | 2.49 | Antihistaminic | Nicotine | 2.06 | Non-drug |
| Carpipramine | 3.29 | Antipsychotic | Nifuradene | 2.44 | Antibacterial |
| Carsalam | 1.79 | Analgesic | Nimidane | 1.52 | Acaricide |
| Carvedilol | 3.27 | Antihypertensive | Norhyoscyamine | 1.39 | Non-drug |
| Cefaclor | 2.67 | Antibacterial | N-phenylmaleimide | 1.4 | Non-drug |
| Chlordantoin | 2.64 | Antifungal | Oxazolam | 2.82 | Anxiolytic |
| Chlordimeform | 2.79 | Acaricide | Oxymetazoline | 3.31 | Adrenergic |
| Cinnarizine | 1.31 | Antihistaminic | Oxymorphone | 2.81 | Analgesic |
| Cinoxacin | 3.06 | Antibacterial | Ozagrel | 1.76 | Antianginal |
| Cinoxate | 2.2 | Non-drug | Perazine | 2.82 | Antipsychotic |
| Clocinizine | 1.39 | Antihistaminic | Pheniramine | 1.76 | Antihistaminic |
| Clothiapine | 2.64 | Antipsychotic | Phenoxybenzamine | 1.52 | Antihypertensive |
| Clotrimazole | 1.39 | Antifungal | Phensuximide | 2.05 | Anticonvulsant |
| Cloxacillin | 2.65 | Antibacterial | Pifarnine | 3.04 | Antiulcerative |
| Coniferyl alcohol | 1.77 | Non-drug | Pirprofen | 2.26 | Anti-inflamatory |
| Coumaphos | 1.89 | Insecticide | Podocarpic acid | 2.65 | Non-drug |
| Coumetarol | 2.63 | Anticoagulant | Propenzolate | 3.02 | Anticholinergic |
| Crotamiton | 2.44 | Fungicide | Propicillin | 2.65 | Antibacterial |
| Crotamiton | 2.44 | Scabicide | Proquazone | 1.53 | Anti-inflamatory |
| Cyclazocine | 1.37 | Narcotic antagonist | Roquinimex | 3.03 | Antineoplastic |
| Cyclexanone | 1.87 | Antitussive | Salsoline | 3.28 | Non-drug |
| Cyclizine | 1.78 | Antiemetic | Sulbentine | 2.44 | Antifungal |
| Delapril | 2.05 | Antihypertensive | Sulfacytine | 2.45 | Antibacterial |
| Diacetyldihydromorph. | 2.65 | Non-drug | Sulfaethidole | 2.07 | Antibacterial |
| Diampromide | 2.44 | Analgesic | Sulfamethoxazole | 1.52 | Antibacterial |
| Diflorasone | 3.32 | Anti-inflamatory | Tetroquinone | 2.06 | Keratolytic |
| Diphenoxylate | 3.18 | Antidiarrheal | Thiazolsulfone | 2.27 | Antibacterial |

Table 3 (*Continued*)

| Compound | $D_1$[a] | Therapeutic category[b] | Compound | $D_1$[a] | Therapeutic category[b] |
|---|---|---|---|---|---|
| Dipivefrin | 2.77 | Antiglaucoma | Thiram | 2.81 | Antiseptic |
| Endralazine | 3.12 | Antihypertensive | Tolazamide | 1.53 | Antidiabetic |
| Endrin | 1.65 | Non-drug | Tribenoside | 2.05 | Non-drug |
| Enprostil | 1.72 | Antiulcerative | Trifluomeprazine | 2.64 | Tranquilizer |
| Enviroxime | 1.33 | Non-drug | Trilostane | 2.27 | Adrenocortical suppres |
| Erdosteine | 2.43 | Mucolytic | Trimipramine | 2.27 | Antidepressant |
| Esaprazole | 2.5 | Antiulcerative | Tuberin | 1.52 | Non-drug |
| Estramustine | 2.06 | Antineoplastic | Velnacrine | 2.05 | Nootropic |
| Etaqualone | 2.27 | Sedative | Xylometazoline | 3.03 | Decongestant |
| Eucatropine | 1.97 | Anticholinergic | Yohimbine | 2.05 | α-adrenergic blocker |

[a] Value from discriminant function $D_1$.
[b] From the Merck index.

Table 4
Discriminant function and classification matrix obtained by LDA study using five different therapeutic groups[a]

| Descriptor | Analgesic | Bronchodilators | Diuretic | Hypolipidemic | Antineoplastics |
|---|---|---|---|---|---|
| $J_1^v$ | 40.918 | 47.228 | 38.018 | 39.503 | 48.771 |
| $^3C_c$ | 6.197 | 3.938 | 7.450 | 3.312 | 4.591 |
| $^4C_{pc}$ | −0.485 | 0.683 | −1.242 | 0.628 | 0.840 |
| Constant | −22.513 | −24.776 | −22.020 | −17.971 | −28.010 |
| Classification matrix | | | | | |
| Group | Correct (%) | Analgesic | Bronchodilators | Diuretic | Hypolipidemic | Antineoplastics |
| Analgesic | 50 | 9 | 2 | 5 | 1 | 1 |
| Bronchodilators | 53 | 1 | 9 | 0 | 4 | 3 |
| Diuretic | 71 | 2 | 1 | 10 | 1 | 0 |
| Hypolipidemic | 65 | 0 | 3 | 2 | 13 | 2 |
| Antineoplastics | 35 | 2 | 7 | 3 | 1 | 7 |
| Total | 54 | 14 | 22 | 20 | 20 | 13 |

[a] $N = 89$, $l = 0.486$, $F = 5.68$.

Table 5
Classification results obtained with LDA by using five different therapeutical groups (equations, in Table 4)

| Drug | Probability of classification | | | | | |
|---|---|---|---|---|---|---|
| | Analgesic | Bronchodilators | Diuretic | Hypolipid | Neoplastic | Classification |
| Therapeutical category: analgesics | | | | | | |
| AAS | 0.335 | 0.034 | 0.519 | 0.013 | 0.099 | Diuretic |
| Clonitazene | 0.147 | 0.366 | 0.052 | 0.15 | 0.286 | Bronchodilators |
| Diclofenac | 0.276 | 0.086 | 0.267 | 0.331 | 0.041 | Hypolipidemic |
| Etersalate | 0.304 | 0.132 | 0.238 | 0.028 | 0.298 | Analgesic |
| Fenoprofen | 0.319 | 0.088 | 0.303 | 0.223 | 0.067 | Analgesic |
| Floctafenine | 0.326 | 0.03 | 0.562 | 0.013 | 0.07 | Diuretic |
| Indomethacin | 0.282 | 0.158 | 0.201 | 0.256 | 0.103 | Analgesic |
| Isonixin | 0.316 | 0.068 | 0.346 | 0.224 | 0.046 | Diuretic |
| Ketorolac | 0.294 | 0.166 | 0.205 | 0.202 | 0.133 | Analgesic |
| Mefenamic acid | 0.294 | 0.067 | 0.315 | 0.286 | 0.038 | Diuretic |
| Metofoline | 0.113 | 0.399 | 0.034 | 0.235 | 0.219 | Bronchodilators |
| Naproxen | 0.285 | 0.132 | 0.215 | 0.281 | 0.088 | Analgesic |
| Paracetamol | 0.298 | 0.137 | 0.223 | 0.03 | 0.312 | Neoplastic |
| Piroxican | 0.335 | 0.115 | 0.293 | 0.149 | 0.109 | Analgesic |
| Salsalate | 0.36 | 0.05 | 0.441 | 0.065 | 0.085 | Diuretic |
| Sulindac | 0.274 | 0.228 | 0.169 | 0.099 | 0.23 | Analgesic |
| Tenoxican | 0.306 | 0.175 | 0.218 | 0.143 | 0.159 | Analgesic |
| Tolmetin | 0.272 | 0.218 | 0.164 | 0.153 | 0.193 | Analgesic |
| Therapeutical category: bronchodilators | | | | | | |
| Adrenaline | 0.174 | 0.273 | 0.065 | 0.126 | 0.362 | Neoplastic |
| Atropine | 0.139 | 0.37 | 0.05 | 0.255 | 0.188 | Bronchodilators |
| Bamifylline | 0.109 | 0.391 | 0.031 | 0.156 | 0.313 | Bronchodilators |

Table 5 (*Continued*)

| Drug | Probability of classification | | | | | Classification |
|------|-----------|----------------|----------|----------|-----------|----------------|
| | Analgesic | Bronchodilators | Diuretic | Hypolipid | Neoplastic | |
| Bevonium | 0.138 | 0.218 | 0.058 | 0.496 | 0.09 | Hypolipidemic |
| Doxofylline | 0.028 | 0.369 | 0.004 | 0.02 | 0.579 | Neoplastic |
| Ephedrinet | 0.156 | 0.132 | 0.085 | 0.572 | 0.054 | Hypolipidemic |
| Fenoterol | 0.244 | 0.245 | 0.13 | 0.165 | 0.215 | Bronchodilators |
| Fenspiride | 0.186 | 0.293 | 0.086 | 0.286 | 0.148 | Bronchodilators |
| Flutropium | 0.155 | 0.294 | 0.063 | 0.348 | 0.139 | Hypolipidemic |
| Ipratropium | 0.089 | 0.469 | 0.024 | 0.188 | 0.23 | Bronchodilators |
| Metoxiphen | 0.102 | 0.411 | 0.029 | 0.221 | 0.237 | Bronchodilators |
| Orciprenaline | 0.113 | 0.382 | 0.033 | 0.109 | 0.363 | Bronchodilators |
| Oxitropium | 0.081 | 0.483 | 0.02 | 0.154 | 0.262 | Bronchodilators |
| Rimiterol | 0.238 | 0.194 | 0.128 | 0.287 | 0.152 | Hypolipidemic |
| Soterenolt | 0.281 | 0.205 | 0.175 | 0.134 | 0.205 | Analgesic |
| Theobromine | 0.086 | 0.346 | 0.022 | 0.048 | 0.498 | Neoplastic |
| Tretoquinol | 0.127 | 0.366 | 0.041 | 0.102 | 0.364 | Bronchodilators |
| Therapeutical category: diuretics | | | | | | |
| Amiloride | 0.342 | 0.034 | 0.489 | 0.09 | 0.045 | Diuretic |
| Esplacto | 0.153 | 0.215 | 0.077 | 0.488 | 0.067 | Hypolipidemic |
| Acetazolamide | 0.214 | 0.005 | 0.757 | 0.019 | 0.005 | Diuretic |
| Bendroflumethiazide | 0.2 | 0.004 | 0.784 | 0.005 | 0.007 | Diuretic |
| Bumetanide | 0.311 | 0.021 | 0.609 | 0.025 | 0.034 | Diuretic |
| Chlorthalidone | 0.312 | 0.027 | 0.553 | 0.084 | 0.023 | Diuretic |
| Ethacrynate | 0.308 | 0.141 | 0.253 | 0.205 | 0.093 | Analgesic |
| Etozolin | 0.109 | 0.438 | 0.033 | 0.089 | 0.33 | Bronchodilators |
| Furosemide | 0.351 | 0.092 | 0.387 | 0.037 | 0.133 | Diuretic |
| Hydrochlorothiazide | 0.333 | 0.04 | 0.546 | 0.032 | 0.05 | Diuretic |
| Piretanide | 0.347 | 0.05 | 0.497 | 0.033 | 0.074 | Diuretic |
| Quinethazone | 0.353 | 0.084 | 0.411 | 0.057 | 0.095 | Diuretic |
| Torasemide | 0.297 | 0.151 | 0.212 | 0.216 | 0.124 | Analgesic |
| Triamterene | 0.349 | 0.053 | 0.407 | 0.129 | 0.063 | Diuretic |
| Therapeutical category: hypolipidemics | | | | | | |
| Acifran. | 0.264 | 0.165 | 0.163 | 0.276 | 0.132 | Hypolipidemic |
| Benzalbutyramide | 0.269 | 0.026 | 0.447 | 0.246 | 0.012 | Diuretic |
| Bezofibrate | 0.098 | 0.288 | 0.026 | 0.398 | 0.189 | Hypolipidemic |
| Clofibrate | 0.035 | 0.459 | 0.005 | 0.144 | 0.357 | Bronchodilators |
| Clomestrone | 0.12 | 0.137 | 0.062 | 0.643 | 0.038 | Hypolipidemic |
| Etiroxate | 0.067 | 0.277 | 0.019 | 0.561 | 0.076 | Hypolipidemic |
| Fenofibrate | 0.077 | 0.352 | 0.017 | 0.321 | 0.232 | Bronchodilators |
| Halofenate | 0.343 | 0.043 | 0.496 | 0.018 | 0.099 | Diuretic |
| Mevastatin | 0.222 | 0.168 | 0.144 | 0.395 | 0.071 | Hypolipidemic |
| Nicoclonate | 0.106 | 0.242 | 0.036 | 0.511 | 0.105 | Hypolipidemic |
| Nicomol | 0.118 | 0.356 | 0.036 | 0.28 | 0.21 | Bronchodilators |
| Oxiniacic acid | 0.258 | 0.104 | 0.179 | 0.015 | 0.444 | Neoplastic |
| Pirozadil | 0.065 | 0.313 | 0.015 | 0.027 | 0.58 | Neoplastic |
| Pravastatin | 0.277 | 0.155 | 0.212 | 0.277 | 0.079 | Hypolipidemic |
| Probucol | 0.069 | 0.205 | 0.022 | 0.651 | 0.053 | Hypolipidemic |
| Sinvastatin | 0.134 | 0.201 | 0.06 | 0.534 | 0.07 | Hypolipidemic |
| Sitosterol | 0.095 | 0.022 | 0.108 | 0.771 | 0.004 | Hypolipidemic |
| Thyrotropic acid | 0.093 | 0.178 | 0.039 | 0.647 | 0.043 | Hypolipidemic |
| Thyroxine | 0.068 | 0.187 | 0.023 | 0.679 | 0.043 | Hypolipidemic |
| Tiadenol | 0.02 | 0.053 | 0.007 | 0.915 | 0.005 | Hypolipidemic |
| Therapeutical category: antineoplastics | | | | | | |
| Altretamine | 0.021 | 0.469 | 0.002 | 0.062 | 0.446 | Bronchodilators |
| Amsacrine | 0.347 | 0.11 | 0.33 | 0.096 | 0.117 | Analgesic |
| Busulfan | 0.329 | 0.036 | 0.525 | 0.01 | 0.1 | Diuretic |
| Carboquone | 0.098 | 0.376 | 0.028 | 0.05 | 0.447 | Neoplastic |
| Carmustine | 0.257 | 0.007 | 0.69 | 0.002 | 0.043 | Diuretic |
| Citarabine | 0.173 | 0.293 | 0.066 | 0.125 | 0.343 | Neoplastic |
| Chlorambucil | 0.342 | 0.096 | 0.369 | 0.12 | 0.073 | Diuretic |
| Dacarbazine | 0.065 | 0.299 | 0.014 | 0.043 | 0.579 | Neoplastic |
| Ddticarbazine | 0.111 | 0.421 | 0.035 | 0.236 | 0.197 | Bronchodilators |
| Etoposide | 0.106 | 0.4 | 0.03 | 0.103 | 0.361 | Bronchodilators |

Table 5 (*Continued*)

| Drug | Probability of classification | | | | | |
|------|----------|----------------|----------|----------|-----------|----------------|
|      | Analgesic | Bronchodilators | Diuretic | Hypolipid | Neoplastic | Classification |
| Fluorouracil | 0.252 | 0.071 | 0.152 | 0.026 | 0.499 | Neoplastic |
| Hidroure | 0.005 | 0.017 | 0 | 0.023 | 0.955 | Neoplastic |
| Mechlorethamine | 0.02 | 0.576 | 0.002 | 0.038 | 0.363 | Bronchodilators |
| Mercaptoacetamide | 0.1 | 0.365 | 0.026 | 0.094 | 0.415 | Neoplastic |
| Metrotexate | 0.348 | 0.104 | 0.333 | 0.106 | 0.11 | Analgesic |
| Pipobroman | 0.055 | 0.399 | 0.012 | 0.025 | 0.509 | Neoplastic |
| Procarbazole | 0.124 | 0.237 | 0.043 | 0.467 | 0.129 | Hypolipidemic |
| Vinblastine | 0.105 | 0.428 | 0.03 | 0.191 | 0.246 | Bronchodilators |
| Vincristine | 0.112 | 0.413 | 0.033 | 0.209 | 0.234 | Bronchodilators |
| Vindesine | 0.11 | 0.395 | 0.033 | 0.256 | 0.206 | Bronchodilators |

be pointed out that, in other contexts, these approaches have already shown to be structurally subtle enough to discriminate the pharmacological activity, even where the position of a single methyl group makes the difference [19].

## 4. Conclusions

In spite of the low size of the data set, the results outlined here suggest that it is possible to achieve a pattern of general pharmacological activity based on molecular topology. Although the set of compounds were designed with a high structural heterogeneity, the results obtained require confirmation with huge compound databases.

The topological charge indices seem to play an important role in such discrimination as well as the Wiener index and the quotients between valence and non-valence connectivity indices.

The success of these models, together with the achievements outlined in the introduction, point out that mere structural extra-mechanistic similarity patterns can play a sound role in drug discovery.

## References

[1] C. Hansch, A. Leo, D. Hoekman, Exploring QSAR, Hydrophobic, Electronic, and Steric Constants, ACS Professional Reference Book, American Chemical Society, Washington, DC, 1995.

[2] R. García-Domenech, J.V. de Julián-Ortiz, Antimicrobial activity characterization in a heterogeneous group of compounds, J. Chem. Inf. Comput. Sci. 38 (1998) 445–449.

[3] J. Devillers, D. Domine, A non-congeneric model for predicting toxicity of organic molecules to *Vibrio fischeri*, SAR QSAR Environ. Res. 10 (1999) 61–70.

[4] J. Apostolakis, A. Caflisch, Computational ligand design, Comb. Chem. High Throughput Screen. 2 (1999) 91–104.

[5] G.W. Bemis, M.A. Murcko, The properties of known drugs. Part 1. molecular frameworks, J. Med. Chem. 39 (1996) 2887–2893.

[6] G.W. Bemis, M.A. Murcko, The properties of known drugs. Part 2. Side chains, J. Med. Chem. 42 (1999) 5095–5099.

[7] V.J. Gillet, P. Willett, J. Bradshaw, Identification of biological profiles using substructural analysis and genetic algorithms, J. Chem. Inf. Comput. Sci. 38 (1998) 165–179.

[8] L.B. Kier, Indexes of molecular shape from chemical graphs, Med. Res. Rev. 7 (1987) 417–440.

[9] R.D. Cramer, G. Redl, C.E. Berkoff, Substructural analysis. A novel approach to the problem of drug design, J. Med. Chem. 17 (1974) 533–535.

[10] A. Ajay, W.P. Walters, M.A. Murcko, Can we learn to distinguish between "drug-like" and "non-drug-like" molecules? J. Med. Chem. 41 (1998) 3314–3324.

[11] R. Brown, Y. Martin, Use of structure-activity data to compare structure-based clustering methods and descriptors for use in compound selection, J. Chem. Inf. Comput. Sci. 36 (1996) 572–584.

[12] J. Sadowski, H. Kubinyi, A scoring for discriminating between drugs and non-drugs, J. Med. Chem. 41 (1998) 3325–3329.

[13] A. Ghose, G. Crippen, Atomic physicochemical parameters for three-dimensional-structure-directed quantitative structure-activity relationships. Part 2. Modeling dispersive and hydrophobic interactions, J. Chem. Inf. Comput. Sci. 27 (1987) 21–35.

[14] J. Wang, K. Ramnarayan, Toward designing drug-like libraries: a novel computational approach for prediction of drug feasibility of compounds, J. Comb. Chem. 1 (1999) 524–533.

[15] T.M. Frimurer, R. Bywater, L. Nærum, L.N. Lauritsen, S. Brunak, Improving the odds in discriminating "drug-like" from "non-drug-like" compounds, J. Chem. Inf. Comput. Sci. 40 (2000) 1315–1324.

[16] S.J. Teague, A.M. Davis, P.D. Leeson, T.I. Oprea, The design of leadlike combinatorial libraries, Angew. Chem. Int. Ed. Engl. 38 (1999) 3743–3748.

[17] T.I. Oprea, Property distribution of drug-related chemical databases, J. Comput. Aided Mol. Des. 14 (2000) 251–264.

[18] P.A. Hunt, QSAR using 2D descriptors and TRIPOS' SIMCA, J. Comput. Aided Mol. Des. 13 (1999) 453–467.

[19] R. Gozalbes, M. Brun-Pascaud, R. García-Domenech, J. Gálvez, P.M. Girard, J.P. Doucet, F. Derouin, Prediction of quinolone activity against *Mycobacterium avium* by molecular topology and virtual computational screening, Antimicrob. Agents Chemother. 44 (2000) 2764–2770.

[20] G.A. Bakken, P.C. Jurs, Classification of multidrug-resistance reversal agents using structure-based descriptors and linear discriminant analysis, J. Med. Chem. 43 (2000) 4534–4541.

[21] R.P. Sheridan, R.B. Nachbar, B.L. Bush, Extending the trend vector: the trend matrix and sample-based partial least squares, J. Comput. Aided Mol. Des. 8 (1994) 323–340 (Erratum p. 634).

[22] J.M. Luco, Prediction of the brain-blood distribution of a large set of drugs from structurally derived descriptors using partial least-squares (PLS) modeling, J. Chem. Inf. Comput. Sci. 39 (1999) 396–404.

[23] J.V. Julián-Ortiz, J. de Gálvez, C. Muñoz-Collado, R. García-Domenech, C. Gimeno-Cardona, Virtual combinatorial syntheses and computational screening of new potential anti-herpes compounds, J. Med. Chem. 42 (1999) 3308–3314.

[24] J. Jaén-Oltra, M.T. Salabert-Salvador, F.J. García-March, F. Pérez-Giménez, F. Tomás-Vert, Artificial neural network applied to prediction of fluorquinolone antibacterial activity by topological methods, J. Med. Chem. 43 (2000) 1143–1148.

[25] E. Estrada, E. Uriarte, A. Montero, M. Teijeira, L. Santana, E.A. De Clercq, A novel approach for the virtual screening and rational design of anticancer compounds, J. Med. Chem. 43 (2000) 1975–1985.

[26] C. de Gregorio, L.B. Kier, L.H. Hall, QSAR modeling with the electrotopological state indices: corticosteroids, J. Comput. Aided Mol. Des. 12 (1998) 557–561.

[27] L.B. Kier, L.H. Hall, The nature of structure-activity relationships and their relation to molecular connectivity, Eur. J. Med. Chem. 12 (1977) 307–312.

[28] L.B. Kier, L.H. Hall, General definition of valence delta-values for molecular connectivity, J. Pharm. Sci. 72 (1983) 1170–1173.

[29] J. Gálvez, R. García-Domenech, M.T. Salabert-Salvador, R. Soler, Charge indices. New topological descriptors, J. Chem. Inf. Comput. Sci. 34 (1994) 520–525.

[30] J. Gálvez, R. García-Domenech, J.V. de Julián-Ortiz, R. Soler, Topological approach to drug design, J. Chem. Inf. Comput. Sci. 35 (1995) 272–284.

[31] H. Wiener, Structural determination of paraffin boiling points, J. Am. Chem. Soc. 69 (1947) 17–20.

[32] BMDP Statistical Program, University of California, Los Angeles, 1990.

[33] J. Gálvez, R. García-Domenech, J.V. de Julián-Ortiz, L. Popa, Pharmacological distribution diagrams: a tool for de novo drug design, J. Mol. Graphics 14 (1996) 272–276.