

This is pre-print author version of an accepted article (Journal of Molecular Graphics and Modelling), available at

<http://www.sciencedirect.com/science/article/pii/S109332631400165X>

DOI: 10.1016/j.jmgn.2014.10.006

## 3D-MoRSE Descriptors Explained

*Oleg Devinyak<sup>†\*</sup>, Dmytro Havrylyuk<sup>‡</sup> and Roman Lesyk<sup>‡</sup>*

<sup>†</sup> Department of Pharmaceutical Disciplines, Uzhgorod National University, 88000, Uzhgorod,  
Ukraine

<sup>‡</sup> Department of Pharmaceutical, Organic and Bioorganic Chemistry, Danylo Halytsky Lviv National  
Medical University, 79010, Lviv, Ukraine

Corresponding author phone: +380 312 612434; e-mail: o.devinyak@gmail.com

### Received date

ABSTRACT: 3D-MoRSE is a very flexible 3D structure encoding framework for chemoinformatics and QSAR purposes due to the range of scattering parameter values and variety of weighting schemes used. While arising in many QSAR studies, up to this time they were considered as hardly interpreted and were treated like a "black box". This study is intended to lift the veil of mystery, providing a comprehensible way to the interpretation of 3D-MoRSE descriptors in QSAR/QSPR studies.

The values of these descriptors are calculated with rather simple equation, but may vary when using differing starting geometries as optimization input. This variation increases with scattering parameter and also is higher for electronegativity weighted and unweighted descriptors. Though each 3D-MoRSE descriptor incorporates the information about the whole molecule structure, its final value is derived

mostly from short-distance (up to 3 Å) atomic pairs. And, if a QSAR study covers structurally similar set of compounds, then the role of 3D-MoRSE descriptor in a model can be interpreted using just several pairs of neighbor atoms. The guide to interpretation process is discussed and illustrated with a case study. Realizing the mathematical concept behind 3D-descriptors and knowing their properties it is easy not only to interpret, but also to predict the importance of 3D-MoRSE descriptors in a QSAR study. The process of prediction is described on the practical example and its accuracy is confirmed with further QSAR modeling.

Keywords: 3D-MoRSE descriptors, QSAR, radial basis function, descriptor interpretation, structure encoding

## INTRODUCTION

QSAR/QSPR modeling is a widely recognized and successful approach to the prediction of compounds' physical and biological properties. Its major goal is to represent the studied attribute as some function of other numerical molecular properties and features called molecular descriptors. And one of the crucial steps in QSAR modeling is the interpretation of obtained models in terms of understandable and clear relationships between structure and activity. The importance of QSAR models interpretation is widely recognized.<sup>1,2</sup> J.C. Dearden et al.<sup>3</sup> defined the lack of mechanistic interpretation as an error in QSAR development. OECD guideline makes this point less strict, stating that "a (Q)SAR should be associated with a mechanistic interpretation, if possible".<sup>4</sup> According to this document, mechanistic interpretation "is the basis for discovery of underlying causal relationships" and its consistency "with other knowledge of fundamental processes in chemistry and toxicology adds to credibility and acceptance of the predictions from the model". The knowledge extracted during interpretation of QSAR model also can be used to design new potent compounds or to describe molecular regions that are involved into interaction with (often unknown) biotarget.<sup>5-9</sup> Citing the OECD guideline further, "mechanistic interpretations of (Q)SARs begin with the number and the

nature of the molecular descriptors used in the model". That is why understanding of the meaning of descriptors is so important during QSAR interpretation step. A lot of descriptors such as partition coefficient logP, HOMO energy, molecular weight, functional group counts and so on are easy interpretable, since their meaning is clear and straightforward. But a major part is much more sophisticated and need sometimes significant efforts during interpretation. One class of such descriptors is represented with 3D-MoRSE. 3D-MoRSE denotes 3D molecular representations of structure based on electron diffraction descriptors and has been introduced in 1996 by J. Schuur, J. Gasteiger and coauthors in two seminal papers.<sup>10,11</sup> These descriptors have found a broad application and have been shown as predominant in a number of QSAR/QSPR studies.<sup>12-22</sup> The majority of these papers describe the effect of 3D-MoRSE values on activity but lacks of interpretation how the values of used 3D-MoRSE descriptors relate to the molecular structure. Actually, the inventors of 3D-MoRSE comprehensively describe computational method used to obtain these descriptors but did not give a key to their interpretation. Two studies of L. Saiz-Urra et al.<sup>14,15</sup> were focused primarily to 3D-MoRSE descriptors giving the conclusion that "the distances among the different atoms will be the principal means to separate the molecules according to their structural features when the other parameters remain constant" and "deeper analysis is necessary to interpret the application of these kinds of descriptors". This paper is intended to evaluate the properties and to disclose the chemical meaning of 3D-MoRSE descriptors.

## THEORY AND INTUITION BEHIND 3D-MORSE DESCRIPTORS

### **Description and calculation of 3D-MoRSE descriptors**

3D-MoRSE (Molecular Representation of Structures based on Electronic diffraction) descriptors were introduced in 1996 by J.H. Schuur, P. Selzer and J. Gasteiger with the motivation for encoding 3D structure of a molecule by a fixed number of variables.<sup>10,11</sup> Indeed, the most obvious way to present 3D structure is its representation within cartesian or internal coordinates. But the statistical methods used in computational chemistry cannot handle such objects; these methods are requiring data with fixed

number of features instead. Simplifying equations used in electron diffraction studies, the authors have got the the function 1:

$$I(s) = \sum_{i=2}^N \sum_{j=1}^{i-1} A_i A_j \frac{\sin sr_{ij}}{sr_{ij}}, \quad (1)$$

where  $s$  is the scattering parameter,  $r_{ij}$  is the euclidean distance between  $i^{\text{th}}$  and  $j^{\text{th}}$  atoms,  $N$  is the total number of atoms and  $A_i$  and  $A_j$  are different atomic properties used as weights. Each term of this function depends on distance and thus may be viewed as a radial basis function itself. Further we will use the words "term", "summand" or "radial basis function" interchangeably. Assigning to  $s$  integer values in the range of 0-31  $\text{\AA}^{-1}$ , 32 values of function 1 can be calculated. These values constitute the 3D-MoRSE code of a molecule. There are few programs that provide possibility to calculate 3D-MoRSE descriptors,<sup>23-25</sup> and we extend this set with our free "3dmorse" program which is written in C++ and currently supports only MOPAC2012 output files as input.<sup>26</sup> 3D-MoRSE descriptors usually are calculated with different weights. For example, widely used DRAGON program (version 5.5) for descriptor calculations proposes five kinds of the descriptors: unweighted 3D-MoRSE and weighted with atomic mass, atomic van der Waals volume, atomic Sanderson electronegativity and atomic polarizability.<sup>24</sup> It is worth to note, that in this program the numeration of 3D-MoRSE descriptors starts from 1, so, for example, Mor01u denotes unweighted descriptor with scattering parameter  $s = 0 \text{\AA}^{-1}$  (since scattering parameter starts from zero), Mor02u denotes descriptor with  $s = 1 \text{\AA}^{-1}$  and so on. This point often is confusing and, due to wide usage of DRAGON in QSAR studies, became the reason for large number of misinterpretations.

For better understanding of 3D-MoRSE descriptors nature, let us calculate several simple examples by hands, starting with methane molecule. So, the simplest descriptor among the studied class is Mor01u. Putting zero instead of  $s$  is giving

$$I(0) = \sum_{i=2}^N \sum_{j=1}^{i-1} A_i A_j \frac{\sin(0 \times r_{ij})}{0 \times r_{ij}} = \sum_{i=2}^N \sum_{j=1}^{i-1} A_i A_j \frac{\sin(0)}{0} \quad (2)$$

Actually, the value of sinus zero divided by zero is undefined, however, the limit (3) is one of the most important limits in trigonometry and is equal to 1.

$$\lim_{\theta \rightarrow 0} \frac{\sin(\theta)}{\theta} = 1 \quad (3)$$

Since the descriptor is unweighted, all  $A_i = A_j = 1$ , and we get simply a number of possible atom pairwise combinations, which is equal to corresponding binomial coefficient and can be calculated with factorials:

$$I(0) = \sum_{i=2}^N \sum_{j=1}^{i-1} 1 = \binom{N}{2} = \frac{N!}{2!(N-2)!} \quad (4)$$

Clearly, this descriptor is a function of number of atoms only, moreover, 3D-MoRSE descriptors with scattering parameter  $s = 0 \text{ \AA}^{-1}$  are always positive for all positive weightings schemes (like atomic mass, van der Waals atom volume, electronegativity, polarizability, but not partial charge).

For methane, which has 5 atoms, Mor01u is equal to

$$Mor01u = \frac{5!}{2!(5-2)!} = 10 \quad (5)$$

For the second descriptor, Mor02u with  $s = 1$ , we should consider another case (keeping in mind, that  $A_i = A_j = 1$  for unweighted 3D-MoRSE):

$$I(1) = \sum_{i=2}^N \sum_{j=1}^{i-1} \frac{\sin(1 \times r_{ij})}{1 \times r_{ij}} \quad (6)$$

To calculate this, the pairwise distances between all atoms in the molecule should be given. Thus the geometry of methane has been modeled with PM7 semiempirical method in MOPAC2012.<sup>27</sup> The results show the same interatomic distance for 4 C-H pairs (1.085 Å) and slightly varying distances for 6 H-H pairs (range 1.771-1.773 Å, mean 1.772 Å). Now, since the methane molecule is highly symmetric and has repeating interatomic distance values, we can write the expression for Mor02u as:

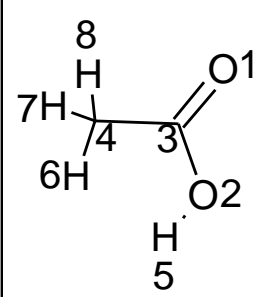
$$Mor02u = 4 \frac{\sin(1.085)}{1.085} + 6 \frac{\sin(1.772)}{1.772} = 3.260 + 3.318 = 6.578. \quad (7)$$

Going to more complicated molecule like acetic acid, we still can measure all interatomic distances with some chemical computation or visualization tool and then provide necessary calculations. Using information from table 1 (geometry has been optimized with PM7 method), we can write the weighted by atomic electronegativities Mor02e as a series of summands. Since the DRAGON program uses carbon-scaled weights, we will also use such convention. Anyway, the difference between the results with original and carbon-scaled weights lays in the constant multiplier, so this point has no effect in QSAR studies. The first term in a sum corresponds to O-O atom pair:

$$1.327 \times 1.327 \times \frac{\sin(1 \times 2.191)}{1 \times 2.191} = 0.654. \quad (8)$$

Similarly, all other terms can be calculated and placed in an upper diagonal matrix, where the diagonal elements will display the contribution of corresponding atom (Table 2). The contribution is represented as a half-sum of all terms where studied atom participates (in such a way diagonal sum is equal to Mor02e value).

Table 1. Parameters of acetic acid 3D structure required to calculate Mor02e

N	Atom pair	Interatomic distance, Å	N	Atom pair	Interatomic distance, Å	N	Atom	Carbon-scaled Sanderson electronegativity
1	O(1)-O(2)	2.191	15	C(3)-H(5)	1.993	1	O	1.327
2	O(1)-C(3)	1.198	16	C(3)-H(6)	2.179	2	C	1
3	O(1)-C(4)	2.411	17	C(3)-H(7)	2.146	3	H	0.942
4	O(1)-H(5)	3.041	18	C(3)-H(8)	2.147			
5	O(1)-H(6)	3.256	19	C(4)-H(5)	2.484			
6	O(1)-H(7)	2.953	20	C(4)-H(6)	1.096			
7	O(1)-H(8)	2.645	21	C(4)-H(7)	1.101			
8	O(2)-C(3)	1.366	22	C(4)-H(8)	1.104			
9	O(2)-C(4)	2.426	23	H(5)-H(6)	2.301			
10	O(2)-H(5)	0.978	24	H(5)-H(7)	2.917			

11	O(2)-H(6)	2.615	25	H(5)-H(8)	3.493	
12	O(2)-H(7)	2.953	26	H(6)-H(7)	1.780	
13	O(2)-H(8)	3.313	27	H(6)-H(8)	1.769	
14	C(3)-C(4)	1.494	28	H(7)-H(8)	1.765	

Table 2. Summands and atom contributions of Mor02e descriptor for acetic acid

	O(1)	O(2)	C(3)	C(4)	H(5)	H(6)	H(7)	H(8)
O(1)	1.178	0.654	1.032	0.367	0.041	-0.044	0.079	0.225
O(2)		1.640	0.951	0.359	1.06	0.24	0.079	-0.064
C(3)			2.086	0.667	0.431	0.355	0.368	0.368
C(4)				1.957	0.232	0.764	0.763	0.762
H(5)					1.016	0.287	0.068	-0.087
H(6)						1.291	0.487	0.492
H(7)							1.169	0.493
H(8)								1.094

This way, a value of Mor02e is set at 11.431, and exactly the same value is given by DRAGON program. Atoms contributions can be further used to identify the most and the less influential atoms or groups in the molecule.

### Theoretical studies of 3D-MoRSE descriptors

Let us consider the roles of all three parameters (pairwise interatomic distance, weighting and scattering) in 3D-MoRSE framework.

#### *Pairwise interatomic distance*

Interatomic distance participates both in the numerator and denominator of radial basis function in eq 1. Thus the result is a periodical function with decreasing amplitude (except Mor01 descriptors where the distance has no any effect). To describe the dependence of 3D-MoRSE summand values on interatomic distance we may visualize the radial basis function of Mor02u. (fig. 1). Since distance is

always positive, the x-axis should be presented with its positive half only (unlike the plot of radial basis function in L. Saiz-Urra et al.<sup>15</sup>).

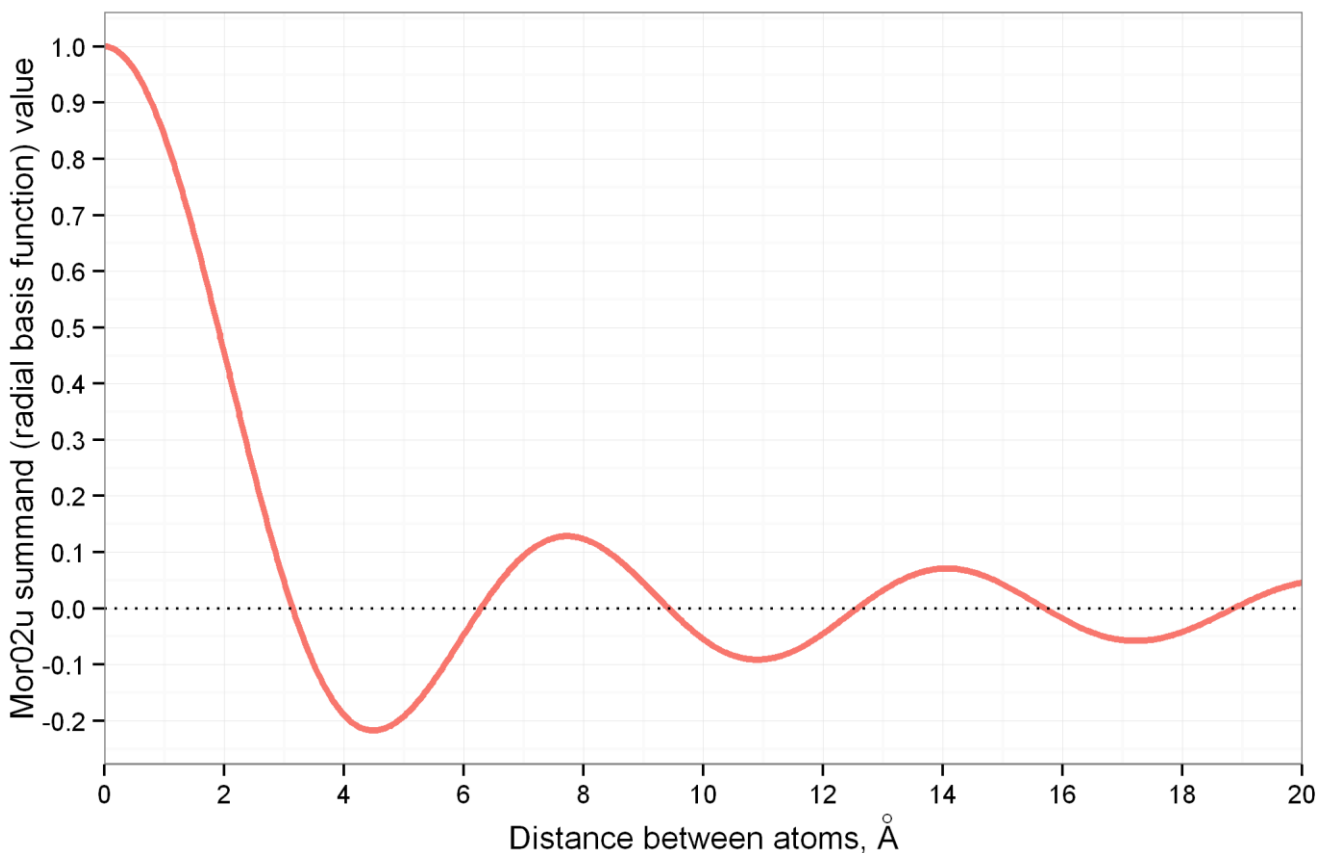


Figure 1. Mor02u radial basis function

When  $s$  is fixed at  $1 \text{ \AA}^{-1}$ , its period is  $2\pi$ . So radial basis function reaches zero at  $\pi$ , has the largest negative value at  $1.5\pi$  (this corresponds to interatomic distance  $4.712 \text{ \AA}$ ), and then begin to increase, crossing zero at  $2\pi$  and so on. Since amplitude drops down quickly, the influential summands of Mor02u are represented by atoms with pairwise distance less than  $2 \text{ \AA}$ . And the effect of atoms located at more than  $5 \text{ \AA}$  apart is tiny.

#### *Weighting scheme*

Being a multiplier at periodic function, weight changes the amplitude of radial basis function. Unweighted descriptors treat each atom equally. Weighting is intended to bring discrimination among atoms, and thus weighted descriptors are sensitive to the presence of specific molecular fragments. The



properties of different weighting schemes can be easily deduced from weighting tables (e.g. table of atomic van der Waals volumes or radii, table of atomic polarizabilities, electronegativities etc), thus only brief summary is presented here:

- **atomic mass** – practically eliminates the role of Hydrogen atoms, while significantly increases the effect of Phosphorus, Sulfur and Chlorine and greatly increases the effect of heavy atoms like Bromine and Iodine on the values of 3D-MoRSE descriptors;
- **van der Waals volume** – significantly decreases the effect of Hydrogen, diminishes the roles of Nitrogen, Oxygen and Fluorine, while giving more influence to Silicon, Phosphorus, Bromine and Iodine;
- **polarizability** – acts similarly to van der Waals volume. The major difference is greatly enhanced effect of metals, but metalorganic compounds are rare in medicinal chemistry;
- **electronegativity** – increases mainly the contributions of Fluorine, Oxygen and Chlorine.
- **atomic partial charge** – unlike the other weighting schemes, this can take both positive and negative values. The weights are not constant across atoms and depend on the surroundings. This weighting may be the most relevant, since it reflects distances between atoms with excessive or deficient electronic density. Unfortunately, it is not implemented in DRAGON, but it is available in ChemoPy<sup>25</sup> and 3dmorse.<sup>26</sup>

Also, since each summand represents some atomic pair, the weights of two atoms are multiplying, so the distance between two influential atoms may define the major part of the descriptor (like two Bromines in a single molecule for weighted by atomic mass 3D-MoRSE descriptors). To understand the effect of weighting, we may look at the plot providing Mor02e summands values on the ordinate axis and corresponding interatomic distance on the abscissa using acetic acid example (fig. 2.)

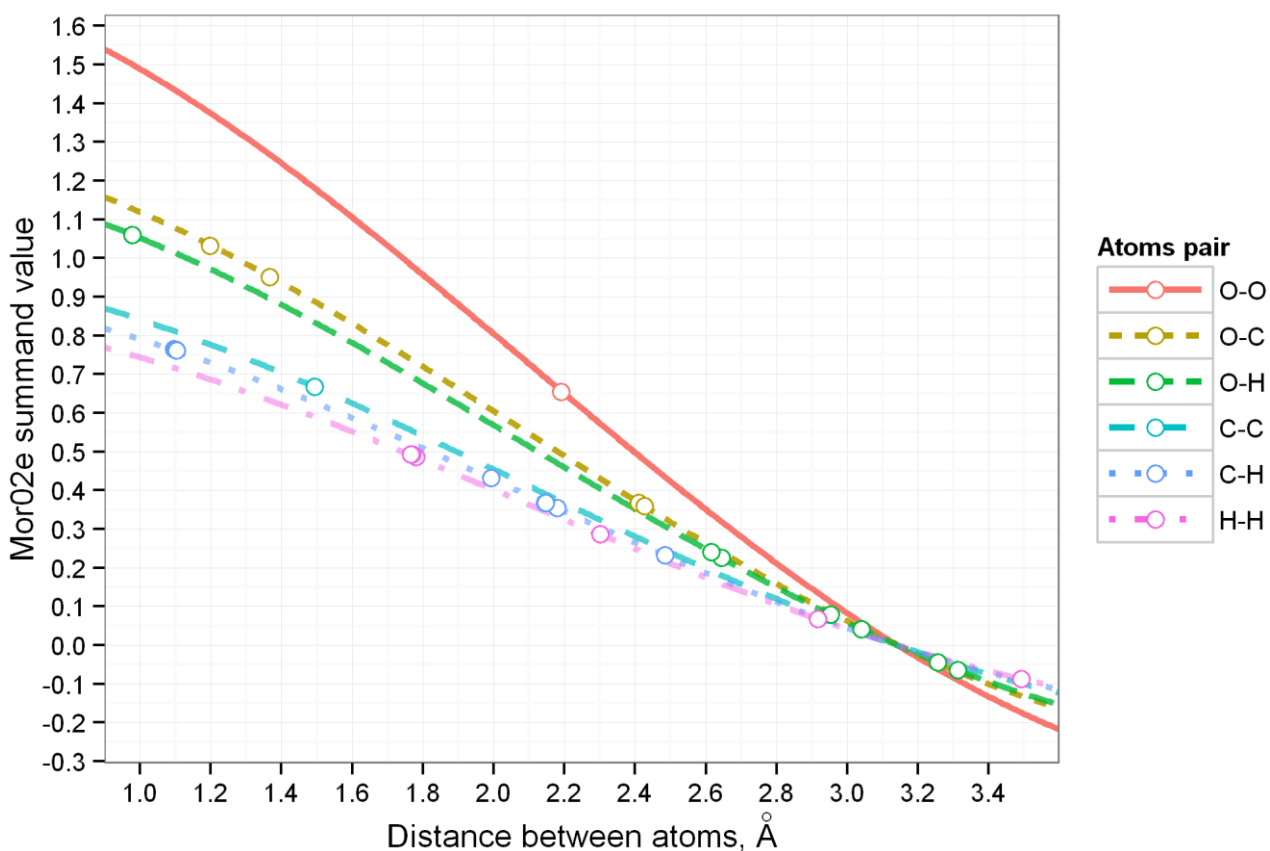


Figure 2. The dependence of Mor02e summand value on interatomic distance and corresponding values for acetic acid (hollow circles)

While electronegativity of Oxygen is only 1.33 times higher the Carbon electronegativity, the effect of O-O pair is 1.77 times higher than that of C-C. However only amplitude is changing, while period remains the same (all curves are crossing at distance that corresponds to  $\pi \text{ \AA}$ ).

### *Scattering parameter*

Like pairwise interatomic distance, the scattering parameter occurs both in numerator and denominator of radial basis function in eq 1. While treating interatomic distance as independent variable, the scattering parameter in numerator can be viewed as angle frequency. So the period of radial basis function decreases with scattering parameter growth. In the same time, the amplitude is also decreasing due to denominator augmentation (fig. 3).

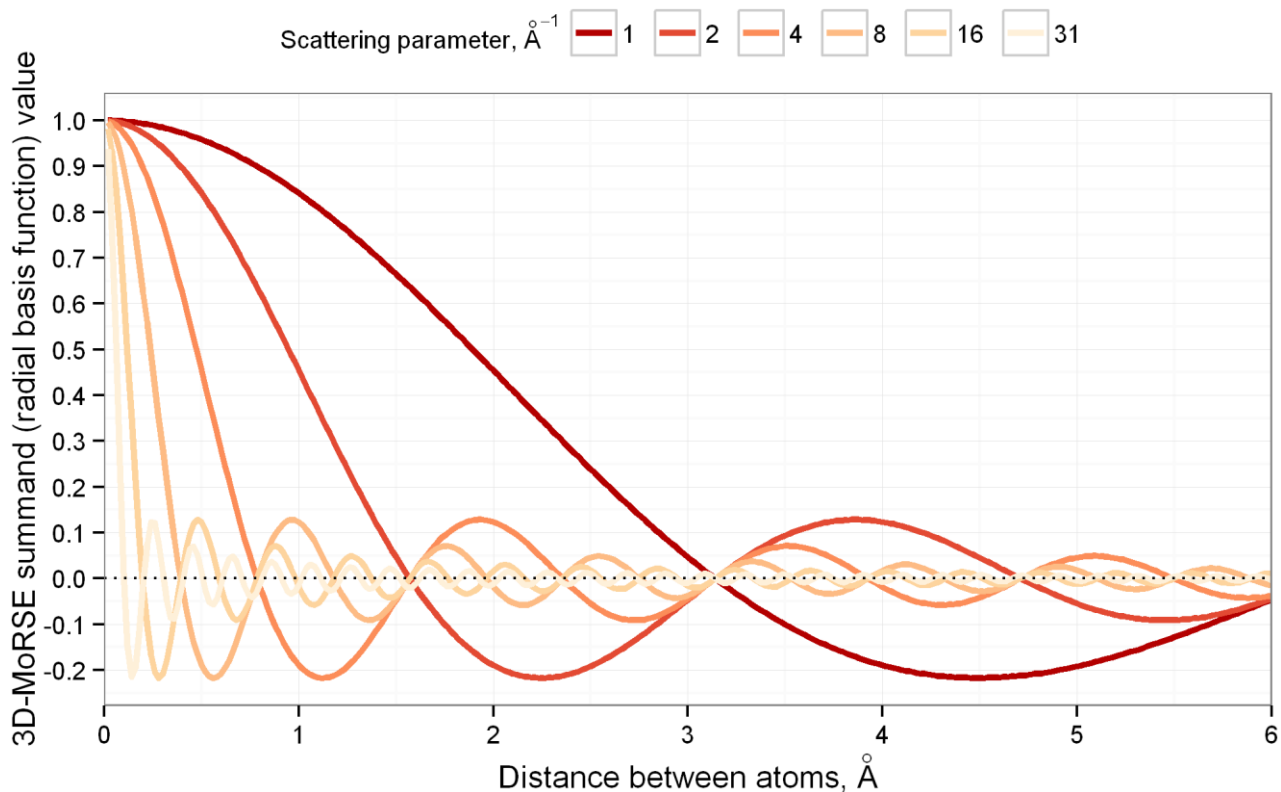


Figure 3. 3D-MoRSE radial basis function at different scattering parameters

### EMPIRICAL PROPERTIES OF 3D-MORSE DESCRIPTORS

In order to test the properties of studied descriptor class we need slightly more complicated model molecule than acetic acid. That molecule should be closer to the structures of modern drugs but also should be as simple as possible to preserve the understanding of the relationship between molecular geometry and descriptors values. So we have chosen  $\gamma$ -aminobutyric acid (GABA) as such model molecule. The next text is organized as follows: for each property the theoretical suggestion are formulated firstly, and then the results of empirical testing are given and discussed.

#### **3D-MoRSE descriptors are random variables.**

The variable in probability and statistics is called random if its value is subject to variations due to chance. 3D-MoRSE descriptors depend on interatomic distances (except the cases when scattering parameter equals zero). Interatomic distances are subject to random variations both in nature (molecular vibration) and in computational studies (due to stochastic methods of molecular geometry

optimization). Consequently, though 3D-MoRSE descriptors are calculated with exact equation 1, indeed their values are drawn from some probability distribution. According to Central Limit Theorem, we can expect this distribution to be normal, since 3D-MoRSE descriptor is represented with sum of terms, each depending on distance between two atoms, which in turn is affected by multiple factors.

To study this point empirically, we have calculated the values of Mor02u for 100 cases of GABA geometry optimization (with PM7 method) using starting preoptimized atom coordinates altered with normally distributed  $N(\mu=0, \sigma^2=0.01)$  random shifts (fig. 4).

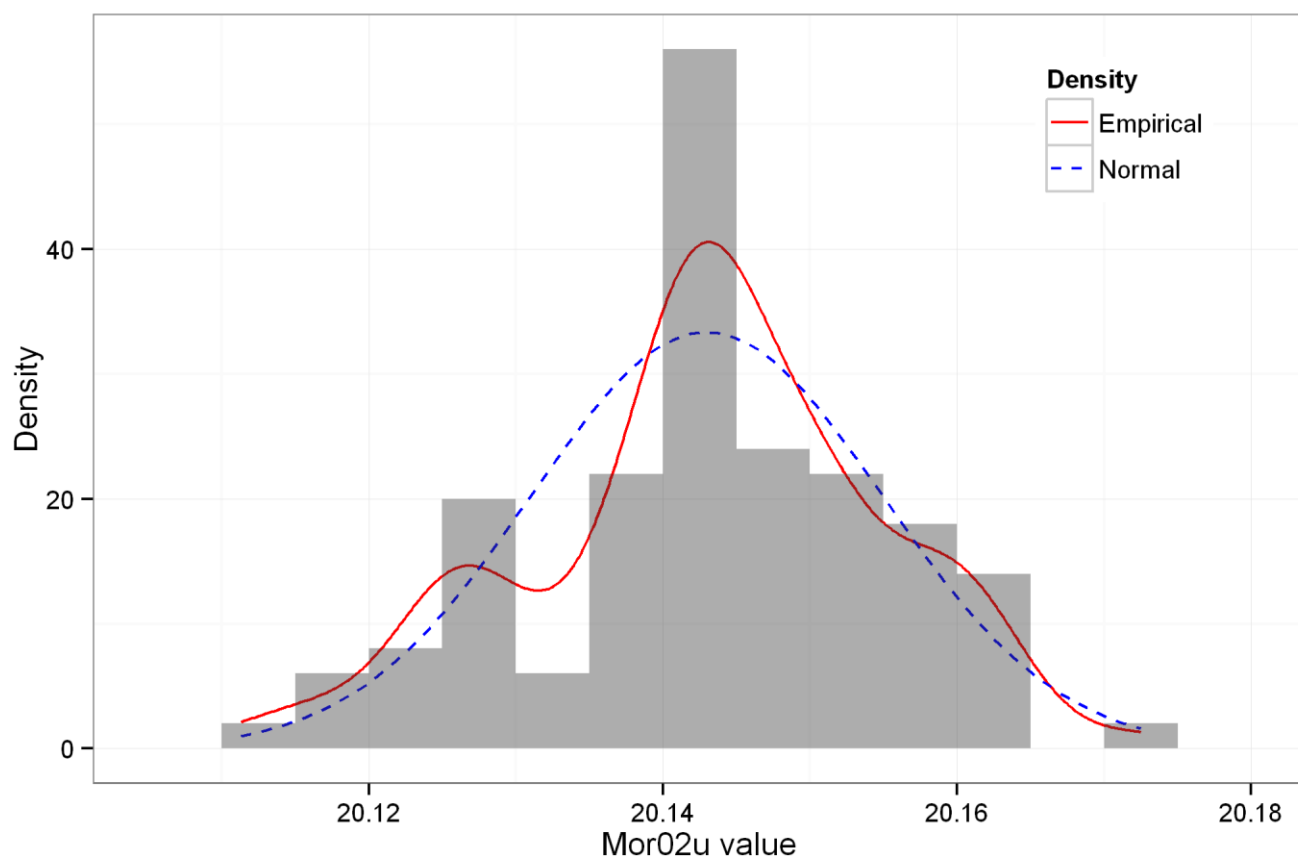


Figure 4. Distribution of Mor02u values for GABA calculated with slightly varying starting atomic coordinates

Standard deviation of the sample of Mor02u values is 0.012. But from practical point of view we are interested not in absolute values of 3D-MoRSE descriptor variation, but in its ability to preserve the discrimination among different structures. That is why the comparison of standard deviation of 3D-

MoRSE values for single compound with the standard deviation for some decoy set of compounds is highly relevant. Clearly, when within-compound standard deviation (noise) comes close to between-compounds standard deviation (signal, i.e. standard deviation of a sample of descriptor values, each corresponding to single compound from some decoy set), it becomes impossible to identify any compound from a set by its 3D-MoRSE descriptor value. For naturally occurring 20  $\alpha$ -amino acids as a decoy set, standard deviation of the sample of 20 Mor02u values is 3.834. In order to score relative variation of Mor02u, we can use the ratio between two standard deviations  $SD_{\text{GABA}}/SD_{\alpha\text{-amino acids}}$  (equivalent to noise-to-signal ratio). When the value of ratio equals to 1 that means complete absence of discriminative power of descriptor. In the studied case the ratio  $SD_{\text{GABA}}/SD_{\alpha\text{-amino acids}}$  is 0.00312 which is small. The visual comparison of between-compounds Mor02u values variation of  $\alpha$ -amino acids and within-compound Mor02u values variation of GABA is provided as supplementary information. Additionally, the shape of distribution (fig. 4.) is not smooth, showing two spikes which may correspond to two close but distinct energetically favorable geometries.

**The relative variation of 3D-MoRSE descriptors increases together with scattering parameter.**

As it has been already shown, the increase in scattering parameter leads to shorter period of radial basis function. When period is rather short, small perturbations in geometry lead to significant changes in the value of radial basis function, which in turn become the source of variation. Since the positions of H-atoms are constrained the least, they contribute much to the variation and thus 3D-MoRSE descriptor weighted with schemes where the role of Hydrogen is diminished should exhibit lower variation. The validity of these theoretical suggestions is clearly demonstrated on the GABA-*vs*-amino acids example by the polynomially smoothed curves of relative variation as a function of scattering parameter (fig. 5.). All weightings show positive correlation between relative variation and scattering parameter with Spearman's  $\rho$  values in 0.77-0.91 range. The lowest relative variation is observed for atomic mass, van der Waals volume and polarizability weightings as it has been predicted. The variation of 3D-MoRSE descriptors with high scattering parameter causes significant decrease in discriminatory ability. For example, the range of GABA Mor32u descriptor values obtained with

slightly modified starting geometry covers the values of three different amino acids, making the differentiation between GABA and those amino acids impossible. The comparison of within-compound and between-compound variations of 3D-MoRSE values at scattering parameters 1 and  $31 \text{ \AA}^{-1}$  is provided as supplementary information.

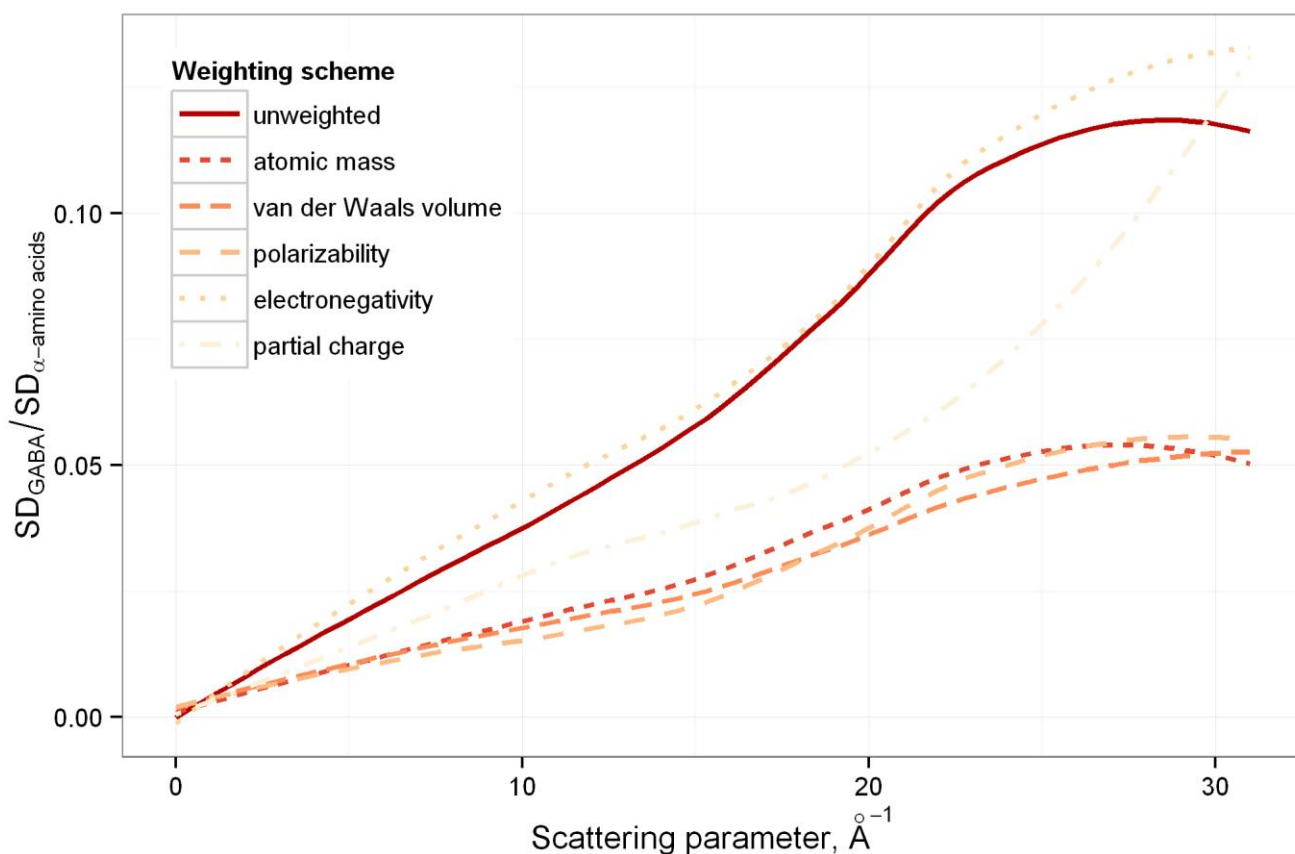


Figure 5. Relationships between relative variation of 3D-MoRSE descriptors and scattering parameter for different weighting schemes

**The effect of distant atom pairs on 3D-MoRSE values is small and can be neglected.**

Since interatomic distance participates in the denominator of radial basis function, the greater distances correspond to smaller summand values. This inference is valuable for the interpretation of 3D-MoRSE descriptors: when looking for atoms and groups responsible for certain 3D-MoRSE value, we can narrow the searching scope to those atomic pairs in which the distance is small. To inspect the impact of different interatomic distances in GABA molecule, we can look at the dynamics of

cumulative sum of 3D-MoRSE terms ordered by interatomic distance (fig. 6.). On the given example, the terms that correspond to atoms located more than  $5\text{\AA}$  away practically has no influence. Moreover, it can be seen that all presented descriptors reach their approximate values before  $3\text{\AA}$ , and the changes of 3D-MoRSE values after this threshold is small. Now how this fact can be utilized? For example, if increasing values of some 3D-MoRSE descriptor lead to increase in biological activity, then atomic groups that contribute to the 3D-MoRSE descriptor mostly are preferential for activity. So searching for these groups we may check only those pairs that are closer than  $3\text{\AA}$  (in most cases such interval corresponds to  $\alpha$  and  $\beta$  atoms relative to the each studied). That is a key point for translating differences in 3D-MoRSE values into structural differences. We have no need to wade through all possible atomic binary combinations: checking the closest neighbors is sufficient, since 3D-MoRSE descriptors are defined predominantly by short-distance atom pairs.

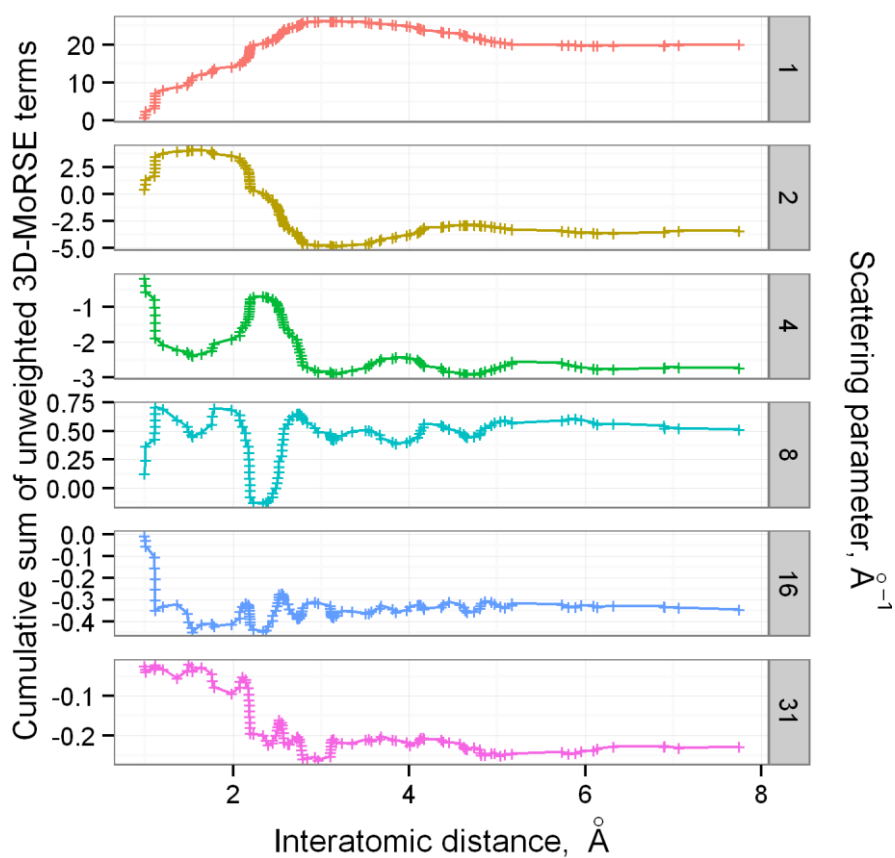


Figure 6. The role of interatomic distance in unweighted 3D-MoRSE descriptors at different scattering parameters

## 3D-MORSE DESCRIPTORS IN ACTION

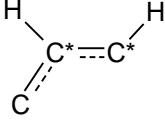
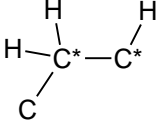
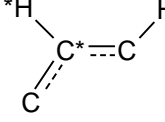
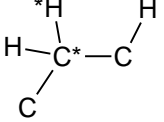
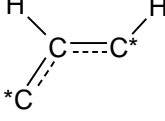
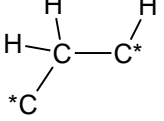
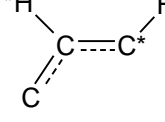
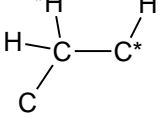
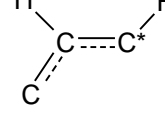
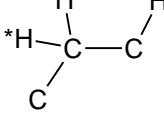
### Prediction of 3D-MoRSE descriptors effects in a QSAR model

One of the key points of scientific method tells that a good theory should make accurate predictions. Let us show how the knowledge about 3D-MoRSE descriptors can help to make reasonable prior suggestions about their effect in a QSAR model. The classification of 20 naturally occurring  $\alpha$ -amino acids into ring-containing (cyclic) and acyclic structures has been utilized as a case study. This classification task is not trivial, since there are four aromatic  $\alpha$ -amino acids (histidine, phenylalanine, tryptophan and tyrosine) and one alicyclic heterocycle (proline). Additionally, only three aromatic  $\alpha$ -amino acids contain phenyl ring, while the fourth (histidine) has imidazole cycle in its structure. Our goal is to predict which 3D-MoRSE descriptors (i.e. which weightings and scattering parameter values) will discriminate ring-containing and acyclic  $\alpha$ -amino acids better. To make a prediction, the first step will be to answer: what are the atomic pairs that can differentiate ring-containing and acyclic structures in the given context? First of all, these are two carbons linked with an aromatic bond, since they are present in aromatic  $\alpha$ -amino acids while in acyclic acids carbons are linked with single bonds. However, not the bond nature is primary for 3D-MoRSE, but bond length. The lengths of aromatic bond (about 1.40 Å) and C(sp<sup>3</sup>)-C(sp<sup>3</sup>) bond (about 1.54 Å) differ, and that contributes to the possibility of distinguishing between the two amino acids classes. The second pair could be C(sp<sup>2</sup>)-H with bond length 1.09 Å, but it is very close to C(sp<sup>3</sup>)-H, which is about 1.11 Å. The difference in 0.02 Å is too small to have some effect. Even the most sensitive to small changes 3D-MoRSE descriptors with the highest scattering parameter 31 Å<sup>-1</sup>, have period  $2\pi/31 = 0.20$  Å. So to obtain effect of amplitude size we need at least 0.05 Å difference. Going to further neighbors, we may select few more atomic pairs that are relevant for classification purposes (table 3.).

Table 3. Closest atomic pairs that could differentiate cyclic  $\alpha$ -amino acids from acyclic ones

Atomic pair in cyclic (aromatic) $\alpha$ -amino acid	Distance, Å	The typical counterpart in acyclic (aliphatic-chain) $\alpha$ -amino acid	Distance, Å	The difference has a tangible effect on 3D-MoRSE
---	-------------	---	-------------	--



				descriptors values	
	1.40			1.54	Yes
	1.09			1.11	No
	2.40			2.50	Yes
	2.15			2.18	No
	2.15			1.76	Yes

Now, to predict the values of scattering parameter that provide the best discrimination, we can sum up the values of radial basis function for the three selected atomic pairs in cyclic amino acids and subtract the values of radial basis function for the corresponding three atomic pairs in acyclic amino acids:

$$\frac{\sin 1.40s}{1.40s} + \frac{\sin 2.40s}{2.40s} + \frac{\sin 2.15s}{2.15s} - \frac{\sin 1.54s}{1.54s} - \frac{\sin 2.50s}{2.50s} - \frac{\sin 1.76s}{1.76s} \quad (9)$$

Carrying out this procedure for each integer value of scattering parameter in the range 1-31 (fig. 7), we may assess the relative discriminative power of 3D-MoRSE descriptors (clearly, the scattering parameter values at which radial basis functions for aromatic acid atomic pairs fall into one phase and radial basis functions for aliphatic acid atomic pairs fall into corresponding counterphase should provide excellent discriminatory ability).

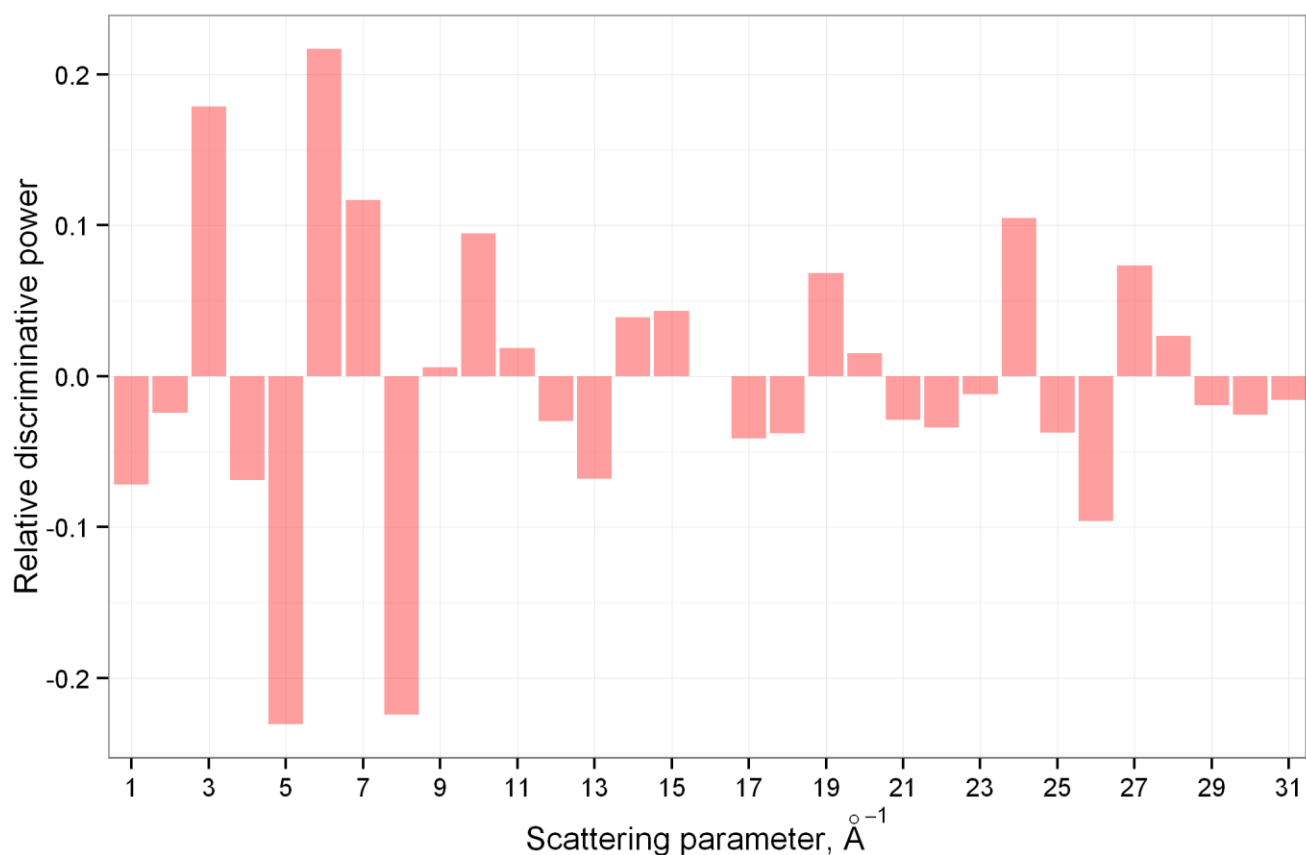


Figure 7. The relative discriminative power of unweighted 3D-MoRSE descriptors predicted for the case classification study

Looking into the obtained results, we may predict that 3D-MoRSE descriptors with scattering parameters 5, 6 and 8 Å<sup>-1</sup> are the most relevant to discriminate cyclic  $\alpha$ -amino acids from acyclic ones. The lowest scattering parameter 0 Å<sup>-1</sup> was not taken into consideration, since corresponding 3D-MoRSE descriptor Mor01u depends only on total number of atoms (see discussion above, in the theoretical part of article). In turn, there is no significant difference in total number of atoms between cyclic and acyclic  $\alpha$ -amino acids. However, ring-containing acids have more carbons, while acyclic have more hydrogens. Thus, when applying weighting schemes which provide carbon atoms with much more influence than hydrogen atoms (atomic mass, polarizability, van der Waals volume), we may expect the rise of discriminative ability. Speaking further about impact of different weightings on relative discriminative ability of 3D-MoRSE descriptors, we may note, that those schemes that give

larger role for carbon atoms should lead to better performance. The weighting by atomic masses greatly diminishes influence of hydrogens, but increase the influence of heteroatoms, the last unlikely contribute to the classification task. In the other hand, weightings by van der Waals volumes and by polarizabilities decrease roles of nitrogen and oxygen atoms as well as hydrogen. Since the major difference between cyclic and acyclic  $\alpha$ -amino acids is presented with carbon atoms, these two weightings are expected to be optimal for the given case study.

In order to test the validity of predictions made, the classification models each using single 3D-MoRSE descriptor have been fitted with logistic regression. Goodness-of-fit has been evaluated as model deviance (lower value of deviance corresponds to better accuracy of model). To see how likely each deviance value can be reached by chance, response permutation test has been carried out for each model with the number of iterations equal to 100.

The results showed good agreement with the prediction (fig. 8). The best classification performance for unweighted 3D-MoRSE descriptors is observed for Mor07u, that corresponds to scattering parameter  $s = 6 \text{ \AA}^{-1}$ . The values of  $s = 5$  and  $8 \text{ \AA}^{-1}$  also lead to reasonable and statistically significant models. As it has been predicted, when  $s = 0 \text{ \AA}^{-1}$  unweighted descriptor Mor01u has no discriminative ability, but Mor01m, Mor01v and Mor01p do have. The suggestion that van der Waals volume and polarizability weighting schemes are favorable is confirmed as well. For sure, there are some other scattering parameter values with discriminative ability, not covered by prediction (especially for van der Waals volume and polarizability weightings). That is because our analysis was rather shallow. It should be emphasized, that all predictions were based entirely on the theory behind 3D-MoRSE and analysis of distances in just few atomic pairs. This case study also illustrates the claim that analysis of the closest atoms (and here we used only the atoms connected directly with covalent bond or through the third atom) is sufficient to obtain the overall picture in 3D-MoRSE framework.

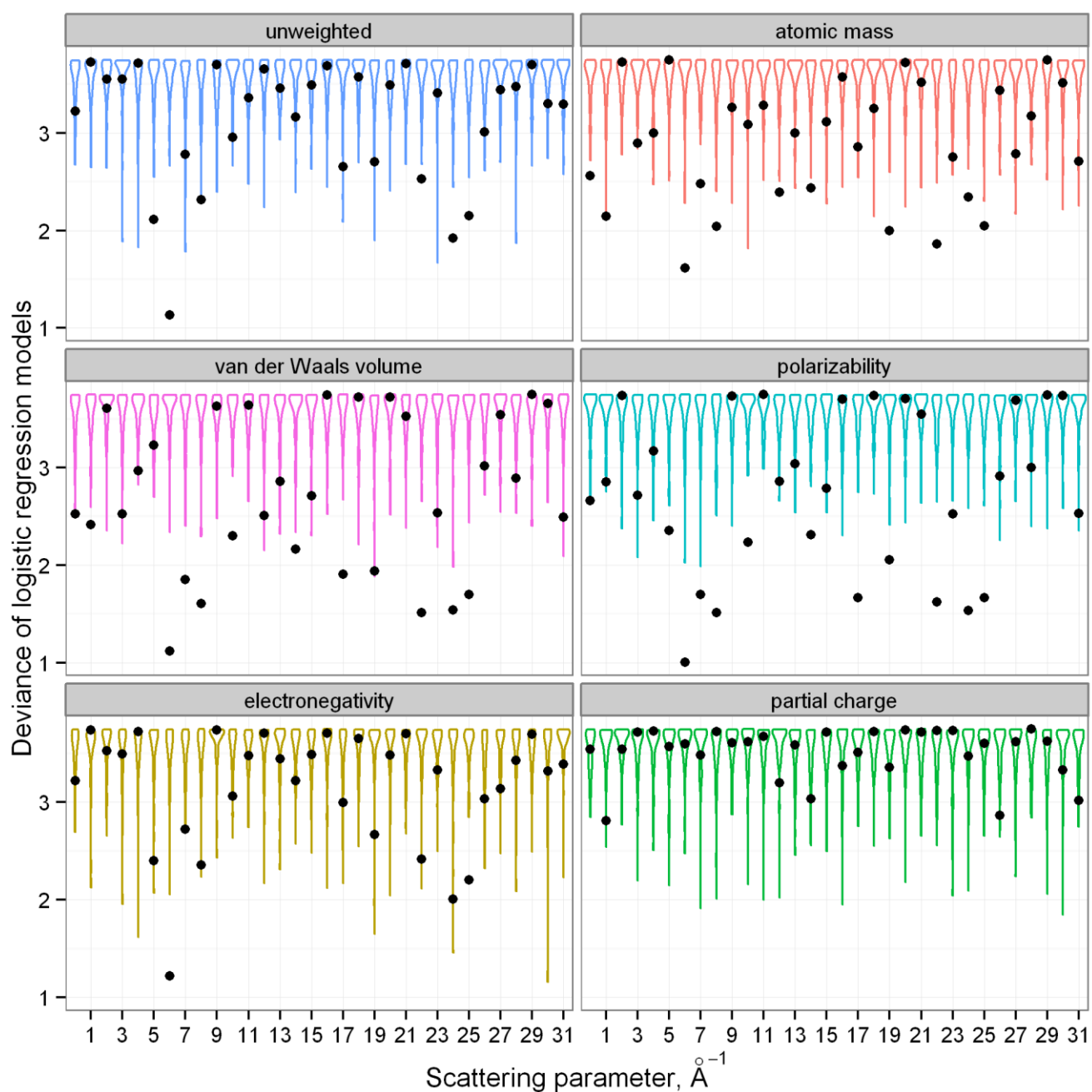


Figure 8. Deviance of logistic regression models of the  $\alpha$ -amino acids classification study developed with different 3D-MoRSE descriptors. The results of response permutation test are illustrated with violin plots

### Interpretation of 3D-MoRSE descriptors in a QSAR model

In the presented above case study the structural features that distinguish two classes were known and the task was to find the most favorable 3D-MoRSE descriptors without using QSAR models. In the

typical QSAR analysis the problem is reverse: the developed QSAR model contains the most favorable 3D-MoRSE descriptors and we should interpretate them. That means we should identify those structural features that are responsible for the desired class (or level) of biological activity.

The road of 3D-MoRSE descriptors interpretation can be passed with the next steps:

1. Estimate the favorable values of 3D-MoRSE descriptors (whether these are high or low values, or the best activity is observed in some optimal range of 3D-MoRSE values).
2. Plot radial basis function with corresponding to studied 3D-MoRSE descriptor scattering parameter and consider weights for different atomic pairs if the descriptor is weighted.
3. Identify on the plot ranges of interatomic distance that lead to desired 3D-MoRSE values.
4. Search in your active compounds the pairs of atoms located at previously found distance. Check whether these pairs are absent in non-active compounds. Those atomic pairs that are present in active compounds and absent in non-active and having interatomic distance in the favorable ranges promote biological activity.

Cramer's steroid data has been used to produce the case QSAR study with 3D-MoRSE interpretation.<sup>28</sup> The data has been downloaded from web resource<sup>29</sup> and the molecular geometries were left as is. The dependent variable was the negative logarithm of corticosteroid-binding globulin (CBG) binding constant. Iterating through the set of 3D-MoRSE descriptors treated as independent variables, simple monoparametric linear regression models were fitted. The whole dataset was employed in model deriving and none validation was carried out since the primary aim is to show 3D-MoRSE descriptor interpretation process, while just simulating the QSAR study to be close to a real one. In such a way, the best monoparametric model (10) has been obtained.

$$\text{pK} = -5.98 - 4.75 \times \text{Mor23v} \quad (10)$$

The determination coefficient of this model is  $R^2 = 0.7385$ , and this value is rather high considering the level of parsimony. The Mor23v descriptor is the 3D-MoRSE descriptor calculated with scattering parameter  $s = 22 \text{ \AA}^{-1}$  and weighted by atomic van der Waals volume. Now to find how this descriptor is related to molecular structure, let us pass the required steps.

*Step 1.* Since high activity in steroid data is expressed with negative numbers (i.e. the less pK values mean the better binding) and  $\beta$  coefficient at the Mor23v descriptors is negative too, the higher Mor23v values correspond to higher activity levels.

*Step 2.* In this stage the plots of corresponding radial basis functions (11) should be produced.

$$f(r) = A_1 A_2 \frac{\sin 22r}{22r}, \quad (11)$$

where  $A_1$  and  $A_2$  are corresponding carbon-scaled atomic van der Waals volumes,  $r$  is interatomic distance.

The steroid dataset include molecules with common atoms C, O, and H (additionally, the compound **31** contains flourine). The carbon-scaled van der Waals volumes for these atoms are 1, 0.715 and 0.263 respectively. The low weight for hydrogen atom makes the contribution of hydrogen-containing atomic pairs negligible (and we would plot radial basis function for the most influential hydrogen-containing atomic pair C-H just for visual support of this claim). Since it is reasonable to search the favorable atomic pairs restricting maximum interatomic distance to 3 Å, and there are no C-C or C-O bonds shorter than 1 Å, we specify the proper limits for the x-axis. The result is presented in figure 9.

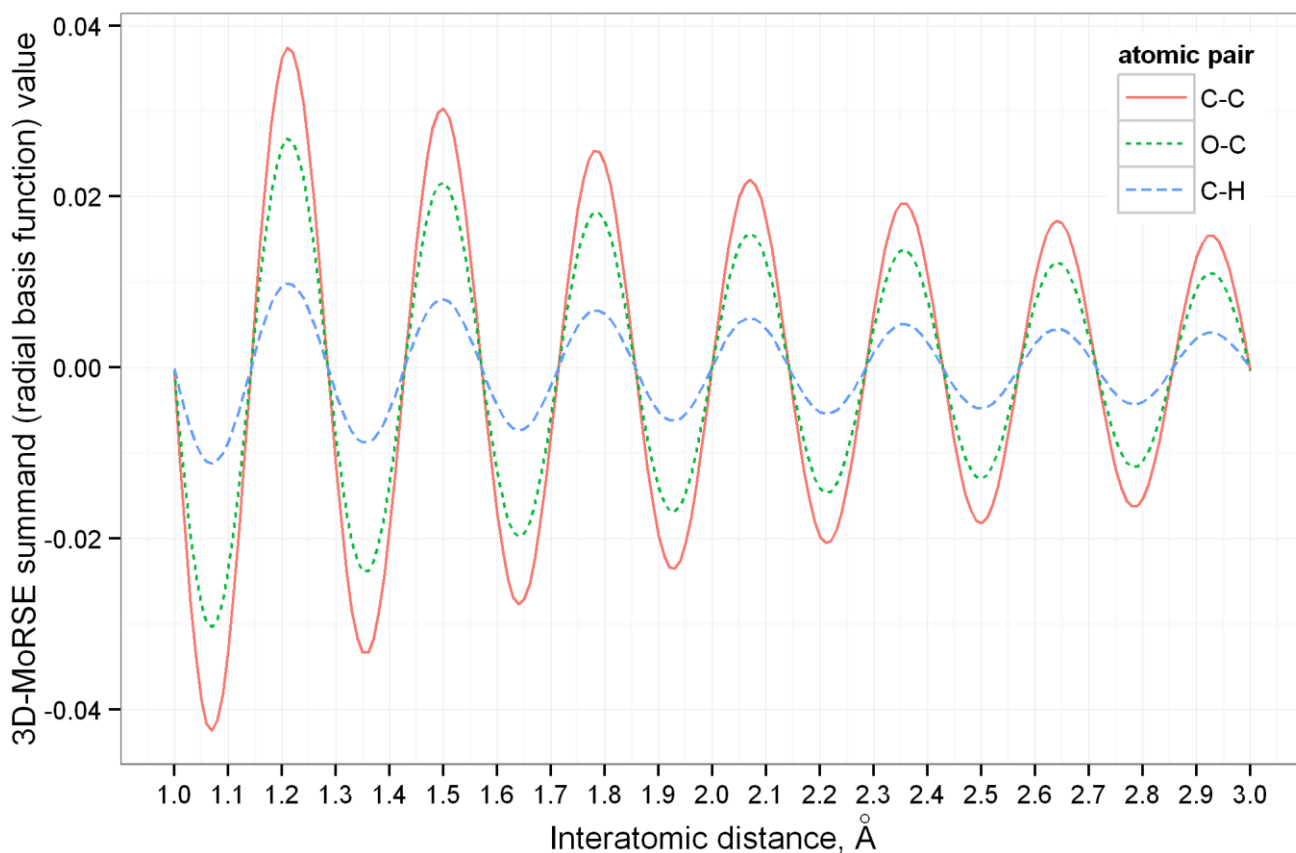


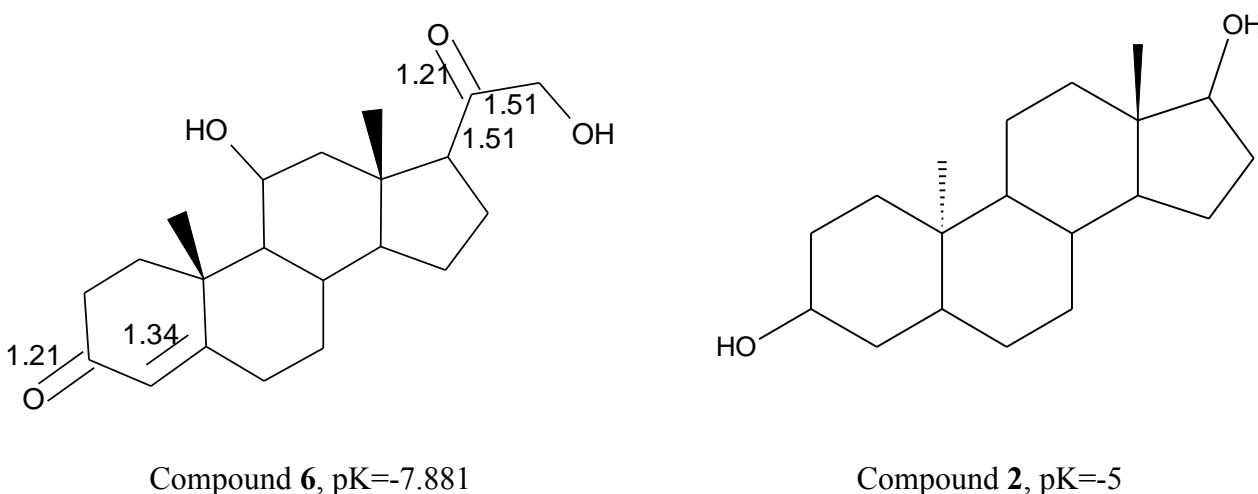
Figure 9. Radial basis functions of Mor23v descriptor corresponding to different atomic pairs

*Step 3.* Looking at the figure 9, the most favorable atomic pairs are located at the distances about 1.21, 1.50, 1.78, 2.07, 2.36, 2.64 and 2.93 Å, while the most detrimental pairs are located at 1.07, 1.36, 1.64, 1.93, 2.21, 2.50 and 2.78 Å. Not only these exact interatomic distances should be searched in studied molecules, but the deviations up to 0.02 Å are acceptable also (for smaller scattering parameter values the acceptable region is larger). The favorable and detrimental distances also can be obtained more accurately with direct calculations. Indeed, the coordinates of extremums for  $\sin sr$  function (other components do not have impact on the period) are

$$\frac{(\frac{\pi}{2} + k\pi)}{s}, \quad (12)$$

where  $s$  is scattering parameter and  $k$  takes all non-negative integer values.

*Step 4.* To search the prespecified distances in studied molecules, we can safely take into consideration only few the most and the least active compounds. Here only two compounds are utilized: compound **6** and compound **2** that are one of the most and the least active compounds respectively (fig. 10). First of all, the single C-C bonds lengths in cyclopentanepiperhydrophenantrene fragment can be omitted since that part is common for all steroid molecules (ignoring the influence of substituents on the steroid core geometry). In the compound **6** C(3)=O, C(20)=O, C(17)-C(20) and C(20)-C(22) bonds have favorable lengths, while C(4)=C(5) bond has unfavorable length (fig. 10). Considering the different amplitude for C-C and O-C atomic pairs, the positive contribution of C-C is higher, despite the double bond in C=O is shorter. C(10)-C(19) and C(13)-C(18) are located at 1.53 Å in both molecules. Such distance is not so far from positive peak at 1.50 Å, so the methyl groups in 10<sup>th</sup> and 13<sup>th</sup> positions of steroid core have small positive effect on the activity predicted. Single C-O bond has length 1.43 Å, which, according to fig. 9, has zero impact. Compound **2** do not show specific favorable or unfavorable close pairs of atoms. Summarizing, the interpretation of Mor23v descriptor in the current case study can be formulated as follows: the presence of carbonyl groups and carbon chain in the 17<sup>th</sup> position have positive impact on the CBG binding affinity, and the presence of double C=C bonds is detrimental.



\* The compounds are depicted in 2D representation just for easy understanding, while the distances are measured in 3D for sure.

Figure 10. Representatives of the most and the least active steroid compounds with the most relevant for Mor23v descriptor interatomic distances.



What about improving the monoparametric model with one more 3D-MoRSE descriptor? Making an exhaustive search, the best two-parametric QSAR model ( $R^2= 0.793$ , the improvement is not so high) is:

$$\text{pK} = -7.61 -5.06 \times \text{Mor23v} -1.36 \times \text{Mor13m} \quad (13)$$

Let us find what is the sense of Mor13m in this QSAR model. Using the known properties of 3D-MoRSE descriptors and making some model diagnostics it is possible to interpretate the descriptor in a more elegant way than four-steps distance analysis. To know which atomic pairs have the greatest influence on Mor13m, the identification of those compounds that have decreased their prediction error significantly when moving from monoparametric to two-parametric model may be used due to the causal chain (fig. 11).

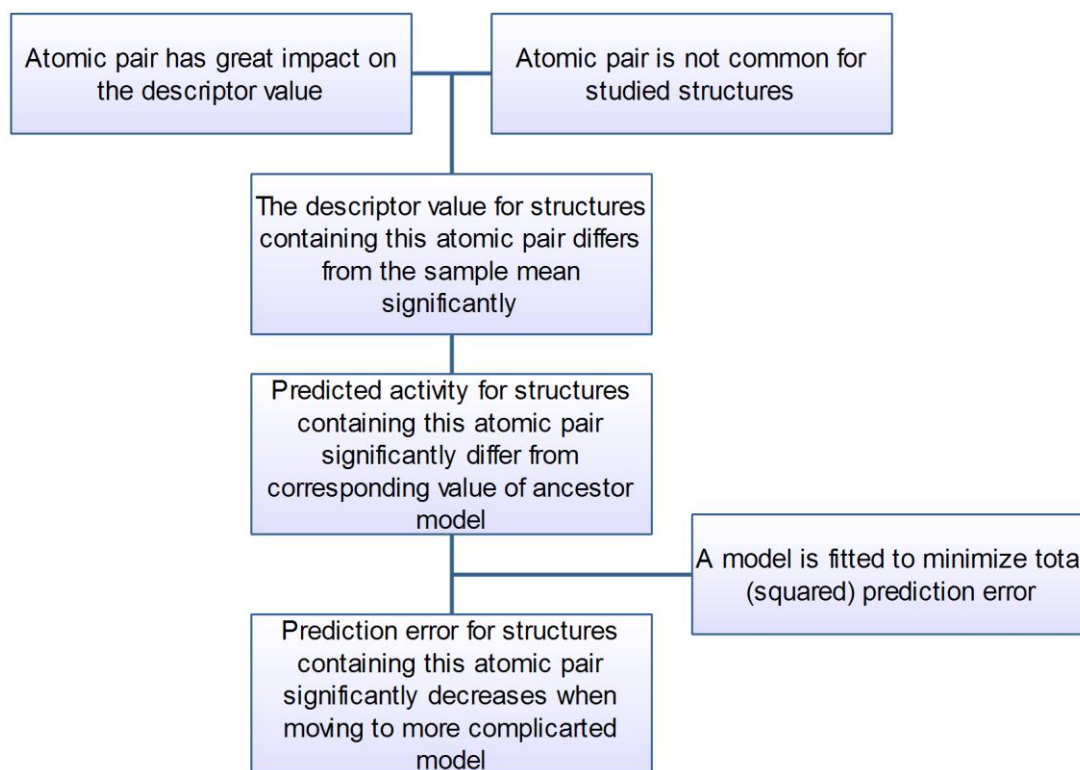


Figure 11. The relationship between the presence of influential atomic pairs in molecular structure and the decrease of prediction error for this structure

The plot of squared residual dynamics (fig. 12) shows that compound **31** has the greatest decrease of prediction error. The term that corresponds to the compound **31** holds 68% of total decrease in the sum

of squared errors (*SSE*). Taking into account, that compound **31** contains fluorine, which is the heaviest atom in the studied dataset, and the chosen second descriptor is weighted accordingly, i.e. with atomic masses, it may be argued that one of Mor13m roles is to penalize the compound **31** for the fluorine atom. Indeed, the C-F distance (1.40 Å) is close to one of the local minimums on the corresponding radial basis function plot (fig. 13.). Additionally, single C-O bonds (each having the length of 1.43 Å) created by three attached hydroxyl groups also fall into the local minimum, thus penalizing the predicted activity of **31** further (fig. 13.).

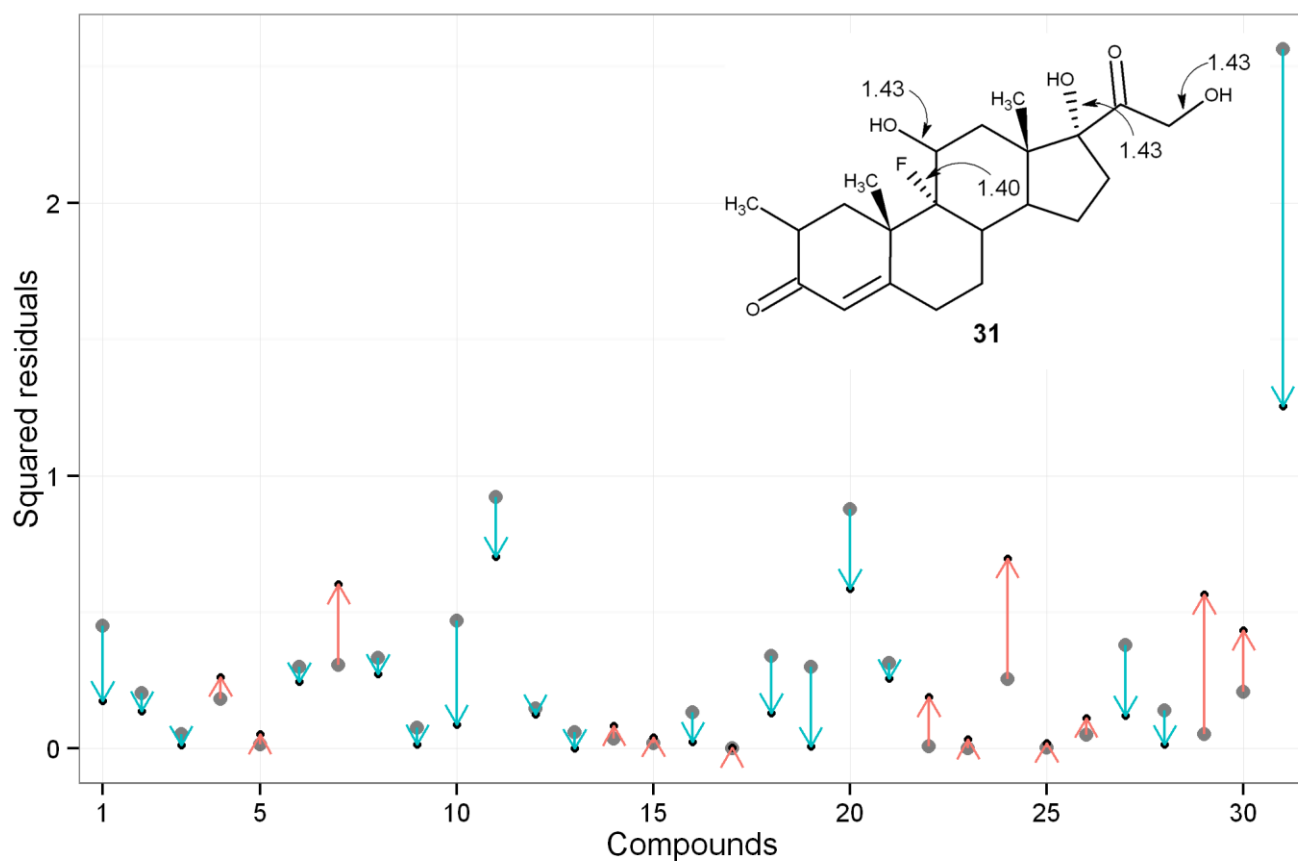


Figure 12. The dynamics of squared residuals when moving from monoparametric to two-parametric QSAR model

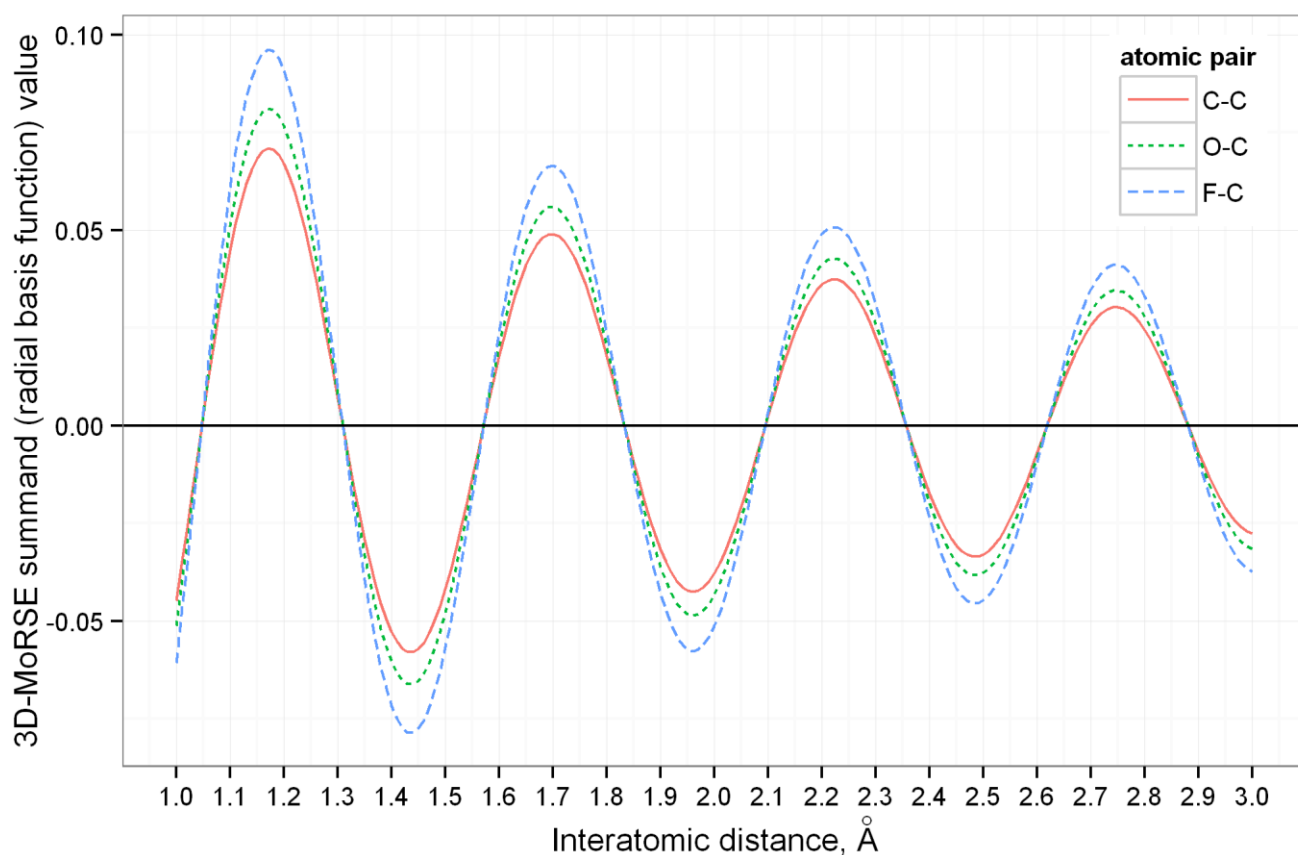


Figure 13. Radial basis functions of Mor13m descriptor corresponding to different atomic pairs

The interpretation of Mor13m descriptor in this case study can be formulated as follows: the presence of fluorine atoms and hydroxyl groups in steroid structure is unfavorable and decreases CBG binding affinity. However, we would like to note, that this inference is caused mostly by the single compound **31**, and thus much credits should not be given to it. Apparently, when a model begin to optimize the prediction error just for single compounds, that is stopping criterion and further complication of such model is inappropriate.

## METHODS

All structures used in the theoretical analysis and case studies were downloaded in three-dimensional representation and further optimized in MOPAC2012 with PM7 method<sup>27</sup> (except the steroids, which were used as-is for the consistency purposes). To calculate 3D-MoRSE descriptors from MOPAC2012 output files the short program "3dmorse" has been written in C++. The possibility to return 3D-MoRSE

summands in a separate file is implemented as a specific feature of the program. The source code and the executable file compiled for Windows environment can be freely downloaded from the web<sup>26</sup>. All other calculations has been carried out in R 3.0.1 environment<sup>30</sup> with additional package for graphics preparation "ggplot2".<sup>31</sup>

## CONCLUSIONS

Lifting the veil of 3D-MoRSE descriptors, we have showed both their strength and weakness. The range of scattering parameter values and variety of weighting schemes provide 3D-MoRSE with significant flexibility, ensuring good discriminatory power even for similar structures. The values of the descriptors are calculated with rather simple equation, but when using differing starting geometries as optimization input the values were unstable and exhibited variation. This variation increased with scattering parameter and also was higher for electronegativity weighted and unweighted descriptors. Though each 3D-MoRSE descriptor incorporates the information about the whole molecule structure, it has been shown that its final value is derived mostly from short-distance (up to 3 Å) atomic pairs. Thus, if a QSAR study covers structurally similar set of compounds, then the role of 3D-MoRSE descriptor in a model can be interpreted using just several pairs of neighbor atoms. As a brief synopsis, to explain the effect of 3D-MoRSE descriptor its radial basis function should be plotted and then atoms that are located at distances that coincide with maxima or minima of the function should be identified. In turn, the presence or absence of the found atomic pairs in a molecule determine its biological activity value (at least this is suggested by the interpreted descriptor). Realizing the mathematical concept behind 3D-descriptors and knowing their properties it is easy not only to interpret, but also to predict the importance of 3D-MoRSE descriptors in a QSAR study.

Finishing the paper, we would like to compare any descriptor class with a human language. We, the medicinal chemists, speak about compounds in terms of their structure: atoms, bonds, groups, templates, scaffolds and substituents. But machine learning methods do not understand such language. Thus we have to project chemical structures into a space of numbers following the predefined rules and equations. Each molecular descriptor is a word, and QSAR models choose the best words to explain

activity phenomenon. Thus, the ability to translate these words into the familiar language of structural entities is crucial for successful application of QSAR method. And, like there are no better or worse languages, there no better or worse descriptor classes. The languages of Alaska Natives can perfectly describe snow, but are poor to describe, for example, technology. In such a way some descriptors are suitable in one case and not relevant in another. 3D-MoRSE descriptors cannot describe complex atomic groups or regions with high or low electron density or some quantum-chemical properties etc, but result in a good model performance when activity variation coincides with variation in interatomic distances due to changes of bonds order and the introduction of new atoms.

#### SUPPLEMENTARY INFORMATION

A figure with the visual comparison of within-compound and between-compound variations in 3D-MoRSE values at scattering parameters 1 and  $31 \text{ \AA}^{-1}$  is provided as supplementary information.

#### AUTHOR INFORMATION

Corresponding Author

\*Phone: +380 312 612434. E-mail: o.devinyak@gmail.com.

Notes

The authors declare no competing financial interest.

#### ACKNOWLEDGMENT

The authors support all people of good will currently struggling for sovereign and unified Ukraine.

#### REFERENCES

- (1) Stanton, D. T. On the physical interpretation of QSAR models. *J. Chem. Inf. Comput. Sci.* **2003**, 43, 1423-1433.
- (2) Todeschini, R.; Consonni, V. *Handbook of molecular descriptors*. Wiley-VCH: 2008. p 688.
- (3) Dearden, J. C.; Cronin, M. T.; Kaiser, K. L. How not to develop a quantitative structure-activity or structure-property relationship (QSAR/QSPR). *SAR QSAR Environ. Res.* **2009**, 20, 241-266.

(4) Organization for Economic Co-operation and Development, Guidance document on the validation of (quantitative) structure-activity relationship ((Q)SAR) models. In *OECD series on testing and assessment 69*, OECD Document ENV/JM/MONO: 2007; p 154.

(5) Todeschini, R.; Gramatica, P. SD-modelling and prediction by WHIM descriptors. Part 5. Theory development and chemical meaning of WHIM descriptors. *Quant. Struct.-Act. Relat.* **1997**, *16*, 113-119.

(6) Randić, M.; Zupan, J. On interpretation of well-known topological indices. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 550-560.

(7) Hu, Q.-N.; Liang, Y.-Z.; Yin, H.; Peng, X.-L.; Fang, K.-T. Structural interpretation of the topological index. 2. The molecular connectivity index, the kappa index, and the atom-type E-state index. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1193-1201.

(8) Devinyak, O.; Slivka, M.; Slivka, M.; Vais, V.; Lendel, V. Quantitative structure-activity relationship study and directed synthesis of thieno[2,3-d]pyrimidine-2,4-diones as monocarboxylate transporter 1 inhibitors. *Med. Chem. Res.* **2012**, *21*, 2263-2272.

(9) Devinyak, O.; Zimenkovsky, B.; Lesyk, R. Biologically active 4-thiazolidinones: A review of QSAR studies and QSAR modeling of antitumor activity. *Curr. Top. Med. Chem.* **2012**, *12*, 2763-2784.

(10) Schuur, J. H.; Selzer, P.; Gasteiger, J. The coding of the three-dimensional structure of molecules by molecular transforms and its application to structure-spectra correlations and studies of biological activity. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 334-344.

(11) Gasteiger, J.; Sadowski, J.; Schuur, J.; Selzer, P.; Steinhauer, L.; Steinhauer, V. Chemical information in 3d space. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 1030-1037.

(12) Kušić, H.; Rasulev, B.; Leszczynska, D.; Leszczynski, J.; Koprivanac, N. Prediction of rate constants for radical degradation of aromatic pollutants in water matrix: A QSAR study. *Chemosphere* **2009**, *75*, 1128-1134.

(13) Schuur, J.; Gasteiger, J. Infrared spectra simulation of substituted benzene derivatives on the basis of a 3d structure representation. *Anal. Chem.* **1997**, *69*, 2398-2405.

- (14) Saíz-Urra, L.; Pérez-Castillo, Y.; Pérez González, M.; Molina Ruiz, R.; Cordeiro, M.; Rodríguez-Borges, J. E.; García-Mera, X. Theoretical prediction of antiproliferative activity against murine leukemia tumor cell line (L1210). 3D MoRSE descriptor and its application in computational chemistry. *QSAR & Comb. Sci.* **2009**, 28, 98-110.
- (15) Saíz-Urra, L.; González, M. P.; Teijeira, M. QSAR studies about cytotoxicity of benzophenazines with dual inhibition toward both topoisomerases i and ii: 3d-morse descriptors and statistical considerations about variable selection. *Biorg. Med. Chem.* **2006**, 14, 7347-7358.
- (16) Duchowicz, P. R.; Vitale, M. G.; Castro, E. A.; Fernández, M.; Caballero, J. QSAR analysis for heterocyclic antifungals. *Biorg. Med. Chem.* **2007**, 15, 2680-2689.
- (17) Cheng, Z.; Zhang, Y.; Fu, W. QSAR study of carboxylic acid derivatives as HIV-1 integrase inhibitors. *Eur. J. Med. Chem.* **2010**, 45, 3970-3980.
- (18) Arab Chamjangali, M.; Beglari, M.; Bagherian, G. Prediction of cytotoxicity data (CC50) of anti-HIV 5-phenyl-1-phenylamino-1H-imidazole derivatives by artificial neural network trained with Levenberg–Marquardt algorithm. *J. Mol. Graph. Model.* **2007**, 26, 360-367.
- (19) Caballero, J.; Fernández, M. Linear and nonlinear modeling of antifungal activity of some heterocyclic ring derivatives using multiple linear regression and bayesian-regularized neural networks. *J. Mol. Model.* **2006**, 12, 168-181.
- (20) Caballero, J.; Vergara-Jaque, A.; Fernández, M.; Coll, D. Docking and quantitative structure–activity relationship studies for sulfonyl hydrazides as inhibitors of cytosolic human branched-chain amino acid aminotransferase. *Mol. Divers.* **2009**, 13, 493-500.
- (21) Sun, M.; Chen, J.; Cai, J.; Cao, M.; Yin, S.; Ji, M. Simultaneously optimized support vector regression combined with genetic algorithm for QSAR analysis of KDR/VEGFR-2 inhibitors. *Chem. Biol. Drug. Des.* **2010**, 75, 494-505.
- (22) Di Tullio, M.; Maccallini, C.; Ammazalorso, A.; Giampietro, L.; Amoroso, R.; De Filippis, B.; Fantacuzzi, M.; Wiczling, P.; Kaliszan, R. QSAR, QSPR and QSRR in terms of 3D-MoRSE descriptors for in silico screening of clofibric acid analogues. *Mol. Inf.* **2012**, 31, 453-458.

- (23) Tetko, I.; Gasteiger, J.; Todeschini, R.; Mauri, A.; Livingstone, D.; Ertl, P.; Palyulin, V.; Radchenko, E.; Zefirov, N.; Makarenko, A.; Tanchuk, V.; Prokopenko, V. Virtual computational chemistry laboratory – design and description. *J. Comput. Aided Mol. Des.* **2005**, *19*, 453-463.
- (24) Dragon program. <http://www.taletе.mi.it> (accessed July 2014).
- (25) Cao, D.-S.; Xu, Q.-S.; Hu, Q.-N.; Liang, Y.-Z. ChemoPy: Freely available python package for computational biology and chemoinformatics. *Bioinformatics* **2013**, *29*(8), 1092-1094.
- (26) 3dmorse program. <http://github.com/devinyak/3dmorse> (accessed July 2014).
- (27) Stewart, J. J. Optimization of parameters for semiempirical methods VI: More modifications to the NDDO approximations and re-optimization of parameters. *J. Mol. Model.* **2013**, *19*, 1-32.
- (28) Coats, E. A.; The COMFA steroids as a benchmark dataset for development of 3D QSAR methods. In *3D QSAR in drug design*, Springer: 1998; pp 199-213.
- (29) 31 steroids binding to the corticosteroid binding globulin (CBG) receptor. <http://www2.chemie.uni-erlangen.de/services/steroids/> (accessed July 2014).
- (30) R Core Team *R: A language and environment for statistical computing*. R Foundation for Statistical Computing: Vienna, Austria, 2012.
- (31) Wickham, H. *ggplot2: Elegant graphics for data analysis*. Springer New York: 2009. p 213.