



## **Guide to MLogP Model version 1.0.0**

### **Table of Contents**

1. Model explanation.....	2
1.1 Introduction.....	2
1.2 Model details.....	2
1.3 Applicability Domain.....	2
1.4 Model statistics.....	4
2. Model usage.....	5
2.1 Input.....	5
2.2 Output.....	5

# 1. Model explanation

## 1.1 Introduction

The model provides a quantitative prediction of water/octanol partition coefficient (LogP). It is implemented inside the VEGA online platform, accessible at: <http://www.vega-qsar.eu/>

## 1.2 Model details

The model is based on the the Moriguchi LogP (MLogP) and consists of a regression equation based on 13 structural parameters as described in: I.Moriguchi, S.Hirono, Q.Liu, I.Nakagome, and Y.Matsushita, Chem.Pharm.Bull. 1992, 40, 127-130; I.Moriguchi, S.Hirono, I.Nakagome, H.Hirano, Chem.Pharm.Bull. 1994, 42, 976-978.

For the purpose of applicability domain assessment, the training set of the Meylan LogP model (9,961 compounds) has been considered, setting all molecules as belonging to the test set.

## 1.3 Applicability Domain

The applicability domain of predictions is assessed using an Applicability Domain Index (ADI) that has values from 0 (worst case) to 1 (best case). The ADI is calculated by grouping several other indices, each one taking into account a particular issue of the applicability domain. Most of the indices are based on the calculation of the most similar compounds found in the training and test set of the model, calculated by a similarity index that consider molecule's fingerprint and structural aspects (count of atoms, rings and relevant fragments).

For each index, including the final ADI, three intervals for its values are defined, such that the first interval corresponds to a positive evaluation, the second one corresponds to a suspicious evaluation and the last one corresponds to a negative evaluation.

Following, all applicability domain components are reported along with their explanation and the intervals used.

- **Similar molecules with known experimental value.** This index takes into account how similar are the first two most similar compounds found. Values near 1 mean that the predicted compound is well represented in the dataset used to build the model, otherwise the prediction could be an extrapolation. Defined intervals are:

$1 \geq \text{index} > 0.9$	strongly similar compounds with known experimental value in the training set have been found
$0.9 \geq \text{index} > 0.75$	only moderately similar compounds with known experimental value in the training set have been found
$\text{index} \leq 0.75$	no similar compounds with known experimental value in the training set have been found

- **Accuracy (average error) of prediction for similar molecules.** This index takes into account the error in prediction for the two most similar compounds found. Values near 0 mean that the predicted compounds falls in an area of the model's space where the model gives reliable predictions, otherwise the greater is the value, the worse the model behaves. Defined intervals are:

index < 0.5	accuracy of prediction for similar molecules found in the training set is good
0.5 <= index < 1.0	accuracy of prediction for similar molecules found in the training set is not optimal
index > 1.0	accuracy of prediction for similar molecules found in the training set is not adequate

- **Concordance with similar molecules (average difference between target compound prediction and experimental values of similar molecules)** . This index takes into account the difference between the predicted value and the experimental values of the two most similar compounds. Values near 0 mean that the prediction made disagrees with the values found in the model's space, thus the prediction could be unreliable. Defined intervals are:

index < 0.5	similar molecules found in the training set have experimental values that agree with the target compound predicted value
0.5 <= index < 1.0	similar molecules found in the training set have experimental values that slightly disagree with the target compound predicted value
index > 1.0	similar molecules found in the training set have experimental values that completely disagree with the target compound predicted value

- **Maximum error of prediction among similar molecules.** This index takes into account the maximum error in prediction among the two most similar compounds. Values near 0 means that the predicted compounds falls in an area of the model's space where the model gives reliable predictions without any outlier value. Defined intervals are:

index < 0.5	the maximum error in prediction of similar molecules found in the training set has a low value, considering the experimental variability
0.5 <= index < 1.0	the maximum error in prediction of similar molecules found in the training set has a moderate value, considering the experimental variability
index >= 1.0	the maximum error in prediction of similar molecules found in the training set has a high value, considering the experimental variability

- **Global AD Index.** The final global index takes into account all the previous indices, in order to give a general global assessment on the applicability domain for the predicted compound. Defined intervals are:

1 >= index > 0.85	predicted substance is into the Applicability Domain of the model
0.85 >= index > 0.75	predicted substance could be out of the Applicability Domain of the model
index <= 0.75	predicted substance is out of the the Applicability Domain of the model

## 1.4 Model statistics

On the pruned training set from EPI Suite KowWin module (9,961 compounds), the logP model has the following statistics:

- Test set:  $n = 9961$ ;  $R^2 = 0.73$ ; RMSE = 0.96

## 2. Model usage

### 2.1 Input

The model accepts as input two molecule formats: SDF (multiple MOL file) and SMILES. All molecules found as input are preprocessed before the calculation of molecular descriptors, in order to obtain a standardized representation of compound. For this reason, some cautions should be taken.

- **Hydrogen atoms.** In SDF files, hydrogen atoms should be explicit. As some times SDF file store only skeleton atoms, and hydrogen atoms are implicit, during the processing of the molecule the system tries to add implicit hydrogens on the basis of the known standard valence of each atom (for example, if a carbon atoms has three single bonds, an hydrogen atom will be added such to reach a valence of four). In SMILES molecules, the default notation uses implicit hydrogen. Anyway please note that in some cases it is necessary to explicitly report an hydrogen; this happens when the conformation is not unambiguous. For example, when a nitrogen atom is into an aromatic ring with a notation like "cnc" it is not clear whether it corresponds to C-N=C or to C-[NH]-C, thus if the situation is the latter, it should be explicitly reported as "c[nH]c".

- **Aromaticity.** The system calculates aromaticity using the basic Hueckel rule. Note that each software for drawing and storing of molecules can use different approaches to aromaticity (for instance, commonly the user can choose between the basic Hueckel rule and a loose approach that lead to considering aromatic a greater number of rings). As in the input files aromaticity can be set explicitly (for instance, in SMILES format by using lowercase letters), during the processing of the molecule the system removes aromaticity from rings that don't satisfy the Hueckel rule. Please note that when aromaticity is removed from a ring, it is not always possible to rebuild the original structure in Kekule form (i.e. with an alternation of single and double bonds, like in the SMILES for benzene, C=1C=CC=CC1), in this case all bonds are set to single. Furthermore, please note that aromaticity detection is a really relevant issue, some molecular descriptors can have significantly different values whether a ring is perceived as aromatic or not. For this reason it is strongly recommended:

- Always use explicit hydrogens in SDF file.
- Avoid explicit aromaticity notation in original files; in this way, the perception of aromaticity is left to the preprocessing step and there is no chance of mistakes due to the transformation of rings that were set to aromatic in the original format but not recognized as aromatic in VEGA.

Note that when some modification of the molecule are performed during the preprocessing (e.g. adding of lacking hydrogens, correction of aromaticity), a warning is given in the remark field of the results.

### 2.2 Output

Results given as text file consist of a plain-text tabbed file (easily importable and processable by any spreadsheet software) containing in each row all the information about the prediction of a molecule. Note that if some problems were encountered while processing the molecule structure, some warning are reported in the last field (Remarks).

Results given as PDF file consists of a document containing all the information about the prediction. For each molecule, results are organized in sections with the following order:

#### *1 – Prediction summary*

Here is reported a depiction of the compound and the final assessment of the prediction (i.e. the prediction made together with the analysis of the applicability domain). Following, all information related to the prediction are reported (the values of the two logP descriptors). Note that if some problems were encountered while processing the molecule structure, some warning are reported in the last field (Remarks).

A graphical representation of the evaluation of the prediction and of its reliability is also provided, using the following elements:



Compound has a low logP value (less than 3.0)



Compound has a high logP value (more than 3.0 and less than 8.0)



Compound has a very high logP value (more than 8.0)



Prediction has low reliability (compound out of the AD)



Prediction has moderate reliability (compound could be out of the AD)



Prediction has high reliability (compound into the AD)

#### *3.1 – Applicability Domain: Similar compounds, with predicted and experimental values*

Here it is reported the list of the six most similar compounds found in the training and test set of the model, along with their depiction and relevant information (mainly experimental value and predicted value).

#### *3.2 – Applicability Domain: Measured Applicability Domain scores*

Here it is reported the list of all Applicability Domain scores, starting with the global Applicability Domain Index (ADI). Note that the final assessment on prediction reliability is given on the basis of the value of the ADI. For each index, it is reported its value and a brief explanation of the meaning of that value.

#### *4.1 – Reasoning: Relevant chemical fragments and moieties*

If some rare and/or missing Atom Centered Fragments are found, they are reported here with a depiction of each fragment.