



Guide to Carcinogenicity IRFMN/Antares 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 Structural Alerts for carcinogen compounds.....	4
1.5 Model statistics.....	7
2. Model usage.....	8
2.1 Input.....	8
2.2 Output.....	8

1. Model explanation

1.1 Introduction

The model provides a qualitative prediction of carcinogenicity (presence of carcinogenic effects in male or female rats). It is implemented inside the VEGA online platform, accessible at: <http://www.vega-qsar.eu/>

1.2 Model details

The model has been built as a set of rules, extracted with Sarpy software from a dataset obtained from the carcinogenicity database of EU-funded project ANTARES. This database is a collection of chemical rat carcinogenesis data (presence of carcinogenic effects in male or female rats) obtained from the CAESAR project database and the “FDA 2009 SAR Carcinogenicity - SAR Structures” database. The CAESAR toxicity values originated from the distributed structure-searchable toxicity DSSTox database, which was built from the Lois Gold’s carcinogenic potency database (CPDB).

The Sarpy software has been used with a cross-validated procedure, ending with the extraction of a set of 127 rules (structural alerts) related to carcinogenic activity. These rules are expressed SMARTS representing molecular fragments.

If at least one rule is matching with the given compound, a “Carcinogen” prediction is given. Otherwise, a “Possible NON-Carcinogen” prediction is given.

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. Note that for purpose of evaluating accuracy and concordance indices, prediction of "Possible NON-Carcinogen" are considered as "NON-Carcinogen".

- **Similar molecules with known experimental value.** This index takes into account how similar are the first three 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:

$1 \geq \text{index} > 0.8$	strongly similar compounds with known experimental value in the training set have been found
$0.8 \geq \text{index} > 0.6$	only moderately similar compounds with known experimental value in the training set have been found
$\text{index} \leq 0.6$	no similar compounds with known experimental value in the training set have been found

- **Accuracy of prediction for similar molecules.** This index takes into account the classification accuracy in prediction for the three most similar compounds found. Values near 1 mean that the predicted compounds falls in an area of the model's space where the model gives reliable predictions (no misclassifications), otherwise the lower is the value, the worse the model behaves. Defined intervals are:

$1 \geq \text{index} > 0.8$	accuracy of prediction for similar molecules found in the training set is good
$0.8 \geq \text{index} > 0.6$	accuracy of prediction for similar molecules found in the training set is not optimal
$\text{index} \leq 0.6$	accuracy of prediction for similar molecules found in the training set is not adequate

- **Concordance for similar molecules .** This index takes into account the difference between the predicted value and the experimental values of the three 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:

$1 \geq \text{index} > 0.8$	similar molecules found in the training set have experimental values that agree with the predicted value
$0.8 \geq \text{index} > 0.6$	some similar molecules found in the training set have experimental values that disagree with the predicted value
$\text{index} \leq 0.6$	similar molecules found in the training set have experimental values that disagree with the predicted value

- **Atom Centered Fragments similarity check.** This index takes into account the presence of one or more fragments that aren't found in the training set, or that are rare fragments. First order atom centered fragments from all molecules in the training set are calculated, then compared with the first order atom centered fragments from the predicted compound; then the index is calculated as following: a first index RARE takes into account rare fragments (those who occur less than three times in the training set), having value of 1 if no such fragments are found, 0.85 if up to 2 fragments are found, 0.7 if more than 2 fragments are found; a second index NOTFOUND takes into account not found fragments, having value of 1 if no such fragments are found, 0.6 if a fragments is found, 0.4 if more than 1 fragment is found. Then, the final index is given as the product RARE * NOTFOUND. Defined intervals are:

index = 1	all atom centered fragment of the compound have been found in the compounds of the training set
$1 > \text{index} \geq 0.7$	some atom centered fragment of the compound have not been found in the compounds of the training set or are rare fragments
$\text{index} < 0.7$	a prominent number of atom centered fragments of the compound have not been found in the compounds of the training set or are rare fragments

- **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:

$1 \geq \text{index} \geq 0.8$	predicted substance is into the Applicability Domain of the model
$0.8 > \text{index} \geq 0.6$	predicted substance could be out of the Applicability Domain of the model
$\text{index} < 0.6$	predicted substance is out of the the Applicability Domain of the model

1.4 Structural Alerts for carcinogen compounds

Following, the list of the 127 rules for carcinogenicity, expressed as SMARTS strings:

SA 1: CN[N+]=O
 SA 2: NNC=O
 SA 3: CN(C=O)N=O
 SA 4: CCCCCCN(C)N=O
 SA 5: CCCN(CCC)N=O
 SA 6: CNCCNN=O
 SA 7: CNCCN(C)N=O
 SA 8: CCNN=O
 SA 9: CCCCCN(C)N=O
 SA 10: CCCCN(C)N=O
 SA 11: CC(O)CNN=O
 SA 12: CN(N=O)C(=O)NCCO
 SA 13: NC(=O)N(CCO)N=O
 SA 14: CN(N=O)C(N)=O
 SA 15: CCCN=O
 SA 16: O(c1ccccc1)c2ccccc2
 SA 17: COc1cccc(O)c1
 SA 18: CCc1ccc(OC)cc1O
 SA 19: CCCN=O
 SA 20: CCOC(=O)C(C)(C)O
 SA 21: CC(C)(O)C(O)=O
 SA 22: Cc1ccccc1-c2ccc(N)cc2
 SA 23: Nc1ccc(cc1)-c2ccc(N)cc2
 SA 24: Nc1ccc(cc1)-c2ccccc2
 SA 25: Nc1ccc(C=C)cc1

SA 26: Nc1cccc(c1)-c2ccccc2
 SA 27: Cc1ccc(N)c(C)c1
 SA 28: Cc1ccccc1N
 SA 29: Nc1ccc(Cc2ccc(N)cc2)cc1
 SA 30: CN(C)c1ccc(Cc2ccccc2)cc1
 SA 31: Cc1ccc(N)cc1
 SA 32: Cc1ccc(NO)cc1
 SA 33: Cc1cccc(N)c1
 SA 34: CNc1ccc(C=C)cc1
 SA 35: Nc1ccc2ccccc2c1
 SA 36: CNc1ccc(N)cc1
 SA 37: Nc1ccccc1O
 SA 38: COc1ccccc1N
 SA 39: Oc1cccc2ccccc12
 SA 40: Nc1ccc(O)c(N)c1
 SA 41: CC(C)C(C)(O)CO
 SA 42: Nc1cccc(c1)S(O)(=O)=O
 SA 43: Cc1cccc2ccccc12
 SA 44: NNCO
 SA 45: CN(N)CO
 SA 46: CC(O)CNN
 SA 47: NNCCO
 SA 48: NNc1ccccc1
 SA 49: NNCC=C
 SA 50: CCNN
 SA 51: CCCNN
 SA 52: CC(=O)NN
 SA 53: CCCN(N)CCC
 SA 54: CCCN(C)N
 SA 55: NN
 SA 56: ClCCCCl
 SA 57: CCl
 SA 58: CCBr
 SA 59: CBr
 SA 60: c1ccc2cc3c(ccc4ccccc34)cc2c1
 SA 61: c1ccc-2c(c1)-c3cccc4cccc-2c34
 SA 62: CN=O
 SA 63: N=O
 SA 64: NO
 SA 65: Oc1cccc(c1)-c2ccccc(O)c2
 SA 66: OS(O)(=O)=O
 SA 67: COS(O)=O
 SA 68: COS(=O)=O
 SA 69: ClC1CCCC(Cl)C1Cl
 SA 70: CC(Cl)CCCCCl
 SA 71: c1ccc(cc1)N=Nc2ccccc2
 SA 72: CCNCCCCl

SA 73: ClCCNCCCCl
 SA 74: Cc1ccnncn1
 SA 75: c1cnncnc1
 SA 76: CC(=C)C(O)=O
 SA 77: CC=C(C)CO
 SA 78: CC(C)=NO
 SA 79: CC(C)=N
 SA 80: Cn1cnncn1
 SA 81: c1nncnn1
 SA 82: COc1ccc(CC=C)cc1
 SA 83: Nc1ncc2ncn(CCCCO)c2n1
 SA 84: OCC1OC(CC1O)n2cnc3cnncnc23
 SA 85: CC1CCC=C(C)C1
 SA 86: C1C=CCC=C1
 SA 87: O=C(OCc1cccc1)c2cccc2
 SA 88: CCOCc1cccc1C
 SA 89: C(=Cc1cccc1)c2cccc2
 SA 90: [O-][N+](=O)c1ccco1
 SA 91: CCNCC(C)=O
 SA 92: N=[N+]
 SA 93: Cc1ccc(cc1)S(O)(=O)=O
 SA 94: O=C1c2cccc2C(=O)c3cccc13
 SA 95: Cc1cccn1
 SA 96: CCCCC(O)CCCCC(O)CCC
 SA 97: Clc1cccc(Cl)c1Cl
 SA 98: COP=O
 SA 99: CC(CN)c1cccc1
 SA 100: OCC#C
 SA 101: NNCc1cccc1
 SA 102: C1CCc2cccc2C1
 SA 103: c1ccsc1
 SA 104: Nc1ccccn1
 SA 105: C1CO1
 SA 106: CC(O)CCCC=O
 SA 107: C[S]=O
 SA 108: c1csn1
 SA 109: CC1COCO1
 SA 110: Nc1ccc([S]c2cccc2)cc1
 SA 111: Cc1ccc(cc1)C(N)=O
 SA 112: CN(C)P(N(C)C)N(C)C
 SA 113: [N+]c1cnnc1
 SA 114: O=C1CCO1
 SA 115: OCCNCC=C
 SA 116: CCNCCCC(C)C
 SA 117: c1ccoc1
 SA 118: CCOC(N)=O
 SA 119: C=CCCCC=O

SA 120: C1CN1
SA 121: c1cc2ccccc2s1
SA 122: Cc1ncc[nH]1
SA 123: [O-][N+](=O)c1ccc(o1)-c2csn2
SA 124: C#C
SA 125: CCF
SA 126: CN=[N+]
SA 127: CCCN=CN

1.5 Model statistics

Following, statistics obtained applying the model to its original dataset:

- Training set: n = 1543; Accuracy = 0.66; Specificity = 0.48; Sensitivity = 0.82

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 warnings

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). Note that if some problems were encountered while processing the molecule structure, some warnings 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 is classified as non-carcinogen



Compound is classified as carcinogen



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.

If some relevant fragments are found (see section 1.4 of this guide), they are reported here (one for each page) with a brief explanation of their meaning and the list of the three most similar compounds that contain the same fragment. Note that if no relevant fragments are found, this section is not shown.