

SciMatics SciQSAR version of commercial MultiCASE model A2E for Structural Alerts for DNA Reactivity
(NTP data)

1. QSAR identifier

1.1 QSAR identifier (title)

SciMatics SciQSAR version of commercial MultiCASE model A2E for Structural Alerts for DNA Reactivity
(NTP data), Danish QSAR Group at DTU Food.

1.2 Other related models

MultiCASE CASE Ultra version of commercial MultiCASE model A2E for Structural Alerts for DNA Reactivity
(NTP data), Danish QSAR Group at DTU Food.

Leadscope Enterprise version of commercial MultiCASE model A2E for Structural Alerts for DNA Reactivity
(NTP data), Danish QSAR Group at DTU Food.

1.3. Software coding the model

SciQSAR version 3.1.00.

2. General information

2.1 Date of QMRF

January 2015.

2.2 QMRF author(s) and contact details

QSAR Group at DTU Food;

Danish National Food Institute at the Technical University of Denmark;

<http://qsar.food.dtu.dk/>;

qsar@food.dtu.dk

Eva Bay Wedebye;

National Food Institute at the Technical University of Denmark;

Nikolai Georgiev Nikolov;

National Food Institute at the Technical University of Denmark;

Marianne Dybdahl;

National Food Institute at the Technical University of Denmark;

Sine Abildgaard Rosenberg;

National Food Institute at the Technical University of Denmark;

2.3 Date of QMRF update(s)

2.4 QMRF update(s)

2.5 Model developer(s) and contact details

MultiCASE Inc.;

23811 Chagrin Blvd Ste 305, Beachwood, OH, 44122, USA;

www.multicase.com

MultiCASE Inc. has kindly given their permission that remodelling of their training set for the commercial A61 model in Leadscope was performed by:

Eva Bay Wedebye;

National Food Institute at the Technical University of Denmark;

Nikolai Nikolov;

National Food Institute at the Technical University of Denmark;

Danish QSAR Group at DTU Food;

National Food Institute at the Technical University of Denmark;

<http://qsar.food.dtu.dk/>;

qsar@food.dtu.dk

2.6 Date of model development and/or publication

January 2014.

2.7 Reference(s) to main scientific papers and/or software package

Contrera, J.F., Matthews, E.J., Kruhlak, N.L., and Benz, R.D. (2004) Estimating the safe starting dose in phase I clinical trials and no observed effect level based on QSAR modelling of the human maximum recommended daily dose. *Regulatory Toxicology and Pharmacology*, 40, 185 – 206.

SciQSAR (2009) Reference guide: *Statistical Analysis and Molecular Descriptors*. Included within the SciMatics SciQSAR software.

2.8 Availability of information about the model

The training set is proprietary and commercially available from MultiCASE Inc. It was originally compiled by MultiCASE Inc. and used to train the commercial MultiCASE A2E model. The Danish QSAR Group bought this

model from MultiCASE Inc. in 1999. Permission to remodel the training set in SciQSAR was kindly granted by MultiCASE Inc. The model algorithm is proprietary from commercial software.

2.9 Availability of another QMRF for exactly the same model

3. Defining the endpoint

3.1 Species

Rodents (rats and mice, both sexes, multiple organs).

3.2 Endpoint

QMR4. Human Health Effects

QMR4.10. Mutagenicity

3.3 Comment on endpoint

Structural alerts (SAs) were introduced by Ashby and co-workers back in the late 80's and are chemical moieties that are either electrophilic or can be metabolized to electrophiles and thereby have the potential to cause electrophilic attack on DNA (Ashby & Tennant 1988). Chemicals containing SAs may in this way be genotoxic and therefore potentially carcinogenic. The SAs reflect specific rules, all of which presumably have been enumerated *a priori*.

Ashby and co-workers identified SAs in chemicals with related US National Toxicology Program (NTP) cancer bioassay and *Salmonella typhimurium* mutagenicity assay data (Ashby & Tennant 1988, Ashby *et al.* 1989). They then made a comparison of SA, mutagenicity in *Salmonella* and ability to induce cancer in rats and mice at several sites and in both sexes (this is a typical pattern seen for genotoxic carcinogens as opposed to non-genotoxic carcinogens that are generally restricted in their range of site, sex and species specificity). A strong correlation between the SA and mutagenicity was found and this was not surprising as the *Salmonella* mutagenicity assay identifies mutagenic chemicals that exert their effect through electrophilic attack on DNA, just like SA is expected to do. Also a strong correlation was found between chemicals containing a SA and causing cancer in rats and mice at several sites and in both sexes (i.e. genotoxic carcinogens). In fact, SA appeared to perform as well as the *Salmonella* mutagenicity assay in predicting the genotoxic carcinogens (Ashby & Tennant 1988, Ashby *et al.* 1989). SA is therefore useful to identify genotoxic carcinogens and mutagens but it is important to be aware of the fact that the SAs are not necessarily an exhaustive list of possible alerts for genotoxic carcinogens, and moreover that chemicals that do not contain a SA may be non-genotoxic carcinogens or non-carcinogens. It should be noted that genotoxic chemicals are not necessarily also mutagens (i.e. lead to mutations after the DNA damage) so the presence of a SA in a chemical does not mean it thereby is also mutagenic.

In the fact sheet for the A2E model (personal communication with MultiCASE in 2001), MultiCASE Inc. refer to three publications of previous versions of the model made with smaller training set (Rosenkranz & Klopman 1990a,b,c). The description of the endpoint for this model is based on the assumption that it is similar to the endpoint described in the paper by Rosenkranz & Klopman (1990b):

Chemicals with cancer bioassay results from US NTP and other databases, as well as results from the *Salmonella typhimurium* mutagenicity assay constitute the data in the training set. The training set chemicals are categorised as positive if they contain a SA, as defined by Ashby and co-workers, and as negative if no SA is found in the chemical. Rather than programming the software to recognize the specific SAs, the chemical structures and the final decisions by Ashby and co-workers as to whether or not the chemicals were classified as possessing or lacking a SA were submitted to the program for modelling. In the modelling process the program identified the structural moieties which were found to be related to activity (biophores) or lack of activity (biophobes).

3.4 Endpoint units

No units, 1 for positives and 0 for negatives.

3.5 Dependent variable

Structural alerts for DNA reactivity, positive or negative.

3.6 Experimental protocol

As the training set is proprietary from MultiCASE Inc. and the data sources are unknown an experimental protocol cannot be described, but as mentioned under 3.3 the training set probably consist of rodent carcinogenicity data from US NTP and other data bases categorized by experts using the rules for SA as described by Ashby and co-workers.

3.7 Endpoint data quality and variability

As the training set is commercial by MultiCASE Inc. the quality and variability of the data used is unknown.

4. Defining the algorithm

4.1 Type of model

This is a categorical (Q)SAR model based on calculated molecular descriptors, and if available the modeller's own or third-party descriptors or measured endpoints can be imported and used as descriptors.

4.2 Explicit algorithm

This is a categorical (Q)SAR model made by use of parametric discriminant analysis to create a linear discriminant function (see 4.5). The specific implementation is proprietary within the SciQSAR software.

4.3 Descriptors in the model

Molecular connectivity indices

Molecular shape indices

Topological indices

Electrotopological (Atom E and HE-States) indices

Electrotopological bond types indices

SciQSAR software provides over 400 built-in molecular descriptors. Additionally, SciQSAR makes it possible to import the modeller's own or third-party descriptors or use measured endpoints as custom descriptors.

4.4 Descriptor selection

The initial descriptor set is manually chosen by the model developer from the total set of built-in descriptors. Furthermore, the set of descriptors applied in the modelling by the program is on top of this selection determined by thresholds for descriptor variance and number of nonzero values likewise defined by the model developer.

68 descriptors were selected from the initial pool of descriptors by the system and used to build the model.

4.5 Algorithm and descriptor generation

For a binary classification problem SciQSAR uses discriminant analysis (DA) to make a (Q)SAR model. SciQSAR implements a broad range of discriminant analysis (DA) methods including parametric and non-parametric approaches. The classic parametric method of DA is applicable in the case of approximately

normal within-class distributions. The method generates either a linear discriminant function (the within-class covariance matrices are assumed to be equal) or a quadratic discriminant function (the within-class covariance matrices are assumed to be unequal). When the distribution is assumed to not follow a particular law or is assumed to be other than the multivariate normal distribution, non-parametric DA methods can be used to derive classification criteria. The non-parametric DA methods available within SciQSAR include the kernel and *k*-nearest-neighbor (kNN) methods. The main types of kernels implemented in SciQSAR include uniform, normal, Epanechnikov, bi-weight, or tri-weight kernels, which are used to estimate the group specific density at each observation. Either Mahalanobis or Euclidean distances can be used to determine proximity between compound-vectors in multidimensional descriptor space. When the kNN method is used, the Mahalanobis distances are based on the pooled covariance matrix. When the kernel method is used, the Mahalanobis distances are based on either the individual within-group covariance matrices or the pooled covariance matrix. (Contrera *et al.* 2004)

If the data outcome is continuous, regression analysis is used to build the predictive model. Within SciQSAR several regression methods are available: ordinary multiple regression (OMR), stepwise regression (SWR), all possible subsets regression (PSR), regression on principal components (PCR) and partial least squares regression (PLS). The choice of regression method depends on the number of independent variables and whether correlation or multicollinearity among the independent variables exists: OMR is acceptable with a small number of independent variables, which are not strongly correlated. SWR is used under the same circumstances as OMR but with greater number of variables. PSR is used for problems with a great number of independent variables. PCR and PLS are useful when a high correlation or multicollinearity exist among the independent variables. (SciQSAR 2009)

To test how stable the developed models are, SciQSAR have built-in cross-validation procedures (see 6.).

For this model, the linear method was used.

4.6 Software name and version for descriptor generation

SciQSAR version 3.1.00.

4.7 Descriptors/chemicals ratio

In this model 68 descriptors were used. The training set consists of 781 compounds. The descriptor/chemical ratio is 1:11.5 (68:781).

5. Defining Applicability Domain

5.1 Description of the applicability domain of the model

The definition of the applicability domain consists of two components; the definition in SciQSAR and the in-house further refinement algorithm on the output from SciQSAR to reach the final applicability domain call.

1. SciQSAR

The first criterion for a prediction to be within the models applicability domain is that all of the descriptor values for the test compound can be calculated by SciQSAR. If SciQSAR cannot calculate each descriptor value for the test chemical no prediction value is given by SciQSAR and it is considered outside the model's applicability domain.

2. The Danish QSAR group

The Danish QSAR group has applied a stricter definition of applicability domain for its SciQSAR models. In addition to the applicability domain definition made by SciQSAR a second criterion has been applied for predictions generated from (Q)SAR models with a binary endpoint. For each prediction SciQSAR calculates the probability (p) for the test compound's membership in one of the two outcome classes (positive or negative). The probability of membership in a class is a measure of how well training set knowledge is able to discriminate a positive prediction from a negative prediction within the nearest space of the subject compound-vector. The probability of membership value is also a measure of the degree of confidence of a prediction. The Danish QSAR group uses this probability for a prediction to further define the model's applicability domain. Only positive predictions with a probability equal to or greater than 0.7 and negative predictions with a probability equal to or less than 0.3 are accepted. Positive predictions with a probability between 0.5 and 0.7 as well as negative predictions with a probability between 0.3 and 0.5 are considered outside the model's applicability domain. When these predictions are wed out the accuracy of the model in general increases at the expense of reduced model coverage. Furthermore, as SciQSAR does not define a structural domain, only predictions which were within either Leadscope structural domain (defined as at least one training set chemical within a Tanimoto distance of 0.7) or CASE Ultra structural domain (no unknown fragments for negatives and maximum 1 unknown fragment for positives) were defined as being inside the SciQSAR applicability domain.

5.2 Method used to assess the applicability domain

The system does not generate predictions if it cannot calculate each descriptor value for the test compound.

Only positive predictions with probability equal to or greater than 0.7 and negative predictions with probability equal to or less than 0.3 were accepted.

5.3 Software name and version for applicability domain assessment

SciQSAR version 3.1.00.

5.4 Limits of applicability

The Danish QSAR group applies an overall definition of structures acceptable for QSAR processing which is applicable for all the in-house QSAR software, i.e. not only SciQSAR. According to this definition accepted structures are organic substances with an unambiguous structure, i.e. so-called discrete organics defined as: organic compounds with a defined two dimensional (2D) structure containing at least two carbon atoms, only certain atoms (H, Li, B, C, N, O, F, Na, Mg, Si, P, S, Cl, K, Ca, Br, and I), and not mixtures with two or more 'big components' when analyzed for ionic bonds (for a number of small known organic ions assumed not to affect toxicity the 'parent molecule' is accepted). Structures with less than two carbon atoms or containing atoms not in the list above (e.g. heavy metals) are rendered out as not acceptable for further QSAR processing. Calculation 2D structures (SMILES and/or SDF) are generated by stripping off accepted organic and inorganic ions. Thus, all the training set and prediction set chemicals are used in their non-ionized form. See 5.1 for further applicability domain definition.

6. Internal validation

6.1 Availability of the training set

No

6.2 Available information for the training set

SMILES

6.3 Data for each descriptor variable for the training set

No

6.4 Data for the dependent variable for the training set

No

6.5 Other information about the training set

781 compounds are in the training set: 503 positives and 278 negatives.

6.6 Pre-processing of data before modelling

Only structures acceptable for SciQSAR were used in the final training set. That is, only discrete organic chemicals as described in 5.4 were used. In case of replicate structures, one of the replicates was kept if all the compounds had the same activity and all were removed if they had different activity. No further structures accepted by the software were eliminated (i.e. outliers).

6.7 Statistics for goodness-of-fit

SciQSARs own internal performance test of the model gave the following Cooper's statistics for predictions within the applicability domain as defined by SciQSAR (i.e. the first criterion described in 5.1):

- Sensitivity (true positives / (true positives + false negatives)): 85.7%
- Specificity (true negatives / (true negatives + false positives)): 86.0%
- Concordance ((true positives + true negatives) / (true positives + true negatives + false positives + false negatives)): 85.8%

6.8 Robustness – Statistics obtained by leave-one-out cross-validation

Not performed.

6.9 Robustness – Statistics obtained by leave-many-out cross-validation

SciQSAR's own internal 10-fold cross-validation (10*10% out) procedure was used for predictions within the applicability domain as defined by SciQSAR (i.e. the first criterion described in 5.1). As the probability domain was not applied (i.e. the second criterion described in 5.2) the accuracy of the predictions when applying this domain can be expected to be higher than reflected in these cross-validation results. This gave the following Cooper's statistics:

- Sensitivity (true positives / (true positives + false negatives)): 81.7%
- Specificity (true negatives / (true negatives + false positives)): 80.6%
- Concordance ((true positives + true negatives) / (true positives + true negatives + false positives + false negatives)): 81.1%

6.10 Robustness - Statistics obtained by Y-scrambling

Not performed.

6.11 Robustness - Statistics obtained by bootstrap

Not performed.

6.12 Robustness - Statistics obtained by other methods

Not performed.

7. External validation

7.1 Availability of the external validation set

7.2 Available information for the external validation set

7.3 Data for each descriptor variable for the external validation set

7.4 Data for the dependent variable for the external validation set

7.5 Other information about the training set

7.6 Experimental design of test set

7.7 Predictivity – Statistics obtained by external validation

7.8 Predictivity – Assessment of the external validation set

7.9 Comments on the external validation of the model

External validation has not been performed for this model.

8. Mechanistic interpretation

8.1 Mechanistic basis of the model

The SciQSAR software provides over 400 calculated physico–chemical, electrotopological E-state, connectivity and other molecular descriptors. The descriptors selected for the model may indicate modes of action that are obvious for persons with expert knowledge about the endpoint.

8.2 A priori or posteriori mechanistic interpretation

A posteriori mechanistic interpretation. The descriptors selected for the model may provide a basis for mechanistic interpretation.

8.3 Other information about the mechanistic interpretation

9. Miscellaneous information

9.1 Comments

The model can predict if a chemical contain a structural alert, i.e. a moiety that can cause electrophilic attack on DNA, and thereby have the potential to be a genotoxic carcinogens. A negative prediction means that the chemical does not contain a structural alert and the chemical can be either a non-genotoxic carcinogen or a non-carcinogen.

9.2 Bibliography

Ashby, J., and Tennant, R.W. (1988) Chemical structure, Salmonella mutagenicity and extent of carcinogenicity as indicators of genotoxic carcinogenesis among 222 chemicals tested in rodents by the U.S. NCI/NTP. *Mutat Res.*, 204, 17-115.

Ashby, J., Tennant, R.W., Zeiger, E., and Stasiewicz, S. (1989) Classification according to chemical structure, mutagenicity to Salmonella and level of carcinogenicity of a further 42 chemicals tested for carcinogenicity by the U.S. National Toxicology Program. *Mutat Res.*, 223, 73-103.

Rosenkranz, H.S., and Klopman, G. (1990a) Structural basis of carcinogenicity in rodents of genotoxicants and non-genotoxicants. *Mutat Res.*, 228:2, 105-124.

Rosenkranz, H.S., and Klopman, G. (1990b) Structural alerts to genotoxicity: the interaction of human and artificial intelligence. *Mutagenesis*, 5:4, 333-361.

Rosenkranz, H.S., and Klopman, G. (1990c) Evaluating the ability of CASE, an artificial intelligence structure-activity relational system, to predict structural alerts for genotoxicity. *Mutagenesis*, 5:6, 525-527.

9.3 Supporting information