

1 A comparison of nine machine learning models based on an
2 expanded mutagenicity dataset and their application for
3 predicting pyrrolizidine alkaloid mutagenicity

4 Christoph Helma^{*1}, Verena Schöning⁴, Philipp Boss³, and Jürgen Drewe²

5 ¹in silico toxicology gmbh, Rastatterstrasse 41, 4057 Basel, Switzerland

6 ²Zeller AG, Seeblickstrasse 4, 8590 Romanshorn, Switzerland

7 ³Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular
8 Medicine in the Helmholtz Association, Robert-Rössle-Strasse 10, Berlin, 13125, Germany

9 ⁴Clinical Pharmacology and Toxicology, Department of General Internal Medicine,
10 Bern University Hospital, University of Bern, Inselspital, 3010 Bern, Switzerland

11 ^{*} Correspondence: Christoph Helma <helma@in-silico.ch>

12 Random forest, support vector machine, logistic regression, neural
13 networks and k-nearest neighbor (**lazar**) algorithms, were applied to new
14 *Salmonella* mutagenicity dataset with 8290 unique chemical structures.

15 **TODO:** PA results

16 Introduction

17 **TODO:** rationale for investigation

18 The main objectives of this study were

- 19 • to generate a new mutagenicity training dataset, by combining the most compre-
20 hensive public datasets

- to compare the performance of MolPrint2D (*MP2D*) fingerprints with Chemistry Development Kit (*CDK*) descriptors
- to compare the performance of global QSAR models (random forests (*RF*), support vector machines (*SVM*), logistic regression (*LR*), neural nets (*NN*)) with local models (*lazar*)
- to apply these models for the prediction of pyrrolizidine alkaloid mutagenicity

Materials and Methods

Data

Mutagenicity training data

An identical training dataset was used for all models. The training dataset was compiled from the following sources:

- Kazius/Bursi Dataset (4337 compounds, Kazius, McGuire, and Bursi (2005)): http://cheminformatics.org/datasets/bursi/cas_4337.zip
- Hansen Dataset (6513 compounds, Hansen et al. (2009)): http://doc.ml.tu-berlin.de/toxbenchmark/Mutagenicity_N6512.csv
- EFSA Dataset (695 compounds EFSA (2016)): <https://data.europa.eu/euodp/data/storage/f/2017-0719T142131/GENOTOX%20data%20and%20dictionary.xls>

Mutagenicity classifications from Kazius and Hansen datasets were used without further processing. To achieve consistency with these datasets, EFSA compounds were classified as mutagenic, if at least one positive result was found for TA98 or T100 Salmonella strains.

Dataset merges were based on unique SMILES (*Simplified Molecular Input Line Entry Specification*) strings of the compound structures. Duplicated experimental data with

the same outcome was merged into a single value, because it is likely that it originated from the same experiment. Contradictory results were kept as multiple measurements in the database. The combined training dataset contains 8290 unique structures and 8309 individual measurements.

Source code for all data download, extraction and merge operations is publicly available from the git repository <https://git.in-silico.ch/mutagenicity-paper> under a GPL3 License. The new combined dataset can be found at <https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mutagenicity.csv>.

Pyrrolizidine alkaloid (PA) dataset

The pyrrolizidine alkaloid dataset was created from five independent, necine base substructure searches in PubChem (<https://pubchem.ncbi.nlm.nih.gov/>) and compared to the PAs listed in the EFSA publication EFSA (2011) and the book by Mattocks Mattocks (1986), to ensure, that all major PAs were included. PAs mentioned in these publications which were not found in the downloaded substances were searched individually in PubChem and, if available, downloaded separately. Non-PA substances, duplicates, and isomers were removed from the files, but artificial PAs, even if unlikely to occur in nature, were kept. The resulting PA dataset comprised a total of 602 different PAs.

The PAs in the dataset were classified according to structural features. A total of 9 different structural features were assigned to the necine base, modifications of the necine base and to the necic acid:

For the necine base, the following structural features were chosen:

- Retronecine-type (1,2-unsaturated necine base, 392 compounds)
- Otonecine-type (1,2-unsaturated necine base, 46 compounds)
- Platynecine-type (1,2-saturated necine base, 140 compounds)

68 For the modifications of the necine base, the following structural features were chosen:

- 69 • N-oxide-type (84 compounds)
- 70 • Tertiary-type (PAs which were neither from the N-oxide- nor DHP-type, 495 com-
71 pounds)
- 72 • Dehydropyrrolizidine-type (pyrrolic ester, 23 compounds)

73 For the necic acid, the following structural features were chosen:

- 74 • Monoester-type (154 compounds)
- 75 • Open-ring diester-type (163 compounds)
- 76 • Macrocyclic diester-type (255 compounds)

77 The compilation of the PA dataset is described in detail in Schöning et al. (2017).

78 Descriptors

79 MolPrint2D (*MP2D*) fingerprints

80 MolPrint2D fingerprints (O’Boyle et al. (2011)) use atom environments as molecular
81 representation. They determine for each atom in a molecule, the atom types of its
82 connected atoms to represent their chemical environment. This resembles basically the
83 chemical concept of functional groups.

84 In contrast to predefined lists of fragments (e.g. FP3, FP4 or MACCs fingerprints) or
85 descriptors (e.g CDK) they are generated dynamically from chemical structures. This
86 has the advantage that they can capture unknown substructures of toxicological relevance
87 that are not included in other descriptors. In addition they allow the efficient calculation
88 of chemical similarities (e.g. Tanimoto indices) with simple set operations.

89 MolPrint2D fingerprints were calculated with the OpenBabel cheminformatics library
90 (O’Boyle et al. (2011)). They can be obtained from the following locations:

91 *Training data:*

- 92 • sparse representation ([https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/fingerprints.mp2d)
93 [mp2d/fingerprints.mp2d](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/fingerprints.mp2d))
- 94 • descriptor matrix ([https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/mutagenicity-fingerprints.csv.gz)
95 [mp2d/mutagenicity-fingerprints.csv.gz](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/mutagenicity-fingerprints.csv.gz))

96 *Pyrrolizidine alkaloids:*

- 97 • sparse representation ([https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/fingerprints.mp2d)
98 [mp2d/fingerprints.mp2d](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/fingerprints.mp2d))
- 99 • descriptor matrix ([https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/pa-fingerprints.csv.gz)
100 [mp2d/pa-fingerprints.csv.gz](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/pa-fingerprints.csv.gz))

101 **Chemistry Development Kit (CDK) descriptors**

102 Molecular 1D and 2D descriptors were calculated with the PaDEL-Descriptors program
103 (<http://www.yapcsoftware.com> version 2.21, Yap (2011)). PaDEL uses the Chemistry De-
104 velopment Kit (CDK, <https://cdk.github.io/index.html>) library for descriptor calcula-
105 tions.

106 As the training dataset contained 8290 instances, it was decided to delete instances
107 with missing values during data pre-processing. Furthermore, substances with equivocal
108 outcome were removed. The final training dataset contained 1442 descriptors for 8083
109 compounds.

110 CDK training data can be obtained from [https://git.in-silico.ch/mutagenicity-paper/](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv)
111 [tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv).

112 The same procedure was applied for the pyrrolizidine dataset yielding descriptors for
113 compounds. CDK features for pyrrolizidine alkaloids are available at [https://git.in-silico.](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/PA-Padel-2D_m2.csv)
114 [ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/PA-Padel-2D_m2.csv](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/PA-Padel-2D_m2.csv).

115 Algorithms

116 **lazar**

117 **lazar** (*lazy structure activity relationships*) is a modular framework for read-across model
118 development and validation. It follows the following basic workflow: For a given chemical
119 structure **lazar**:

- 120 • searches in a database for similar structures (neighbours) with experimental data,
- 121 • builds a local QSAR model with these neighbours and
- 122 • uses this model to predict the unknown activity of the query compound.

123 This procedure resembles an automated version of read across predictions in toxicology,
124 in machine learning terms it would be classified as a k-nearest-neighbour algorithm.

125 Apart from this basic workflow, **lazar** is completely modular and allows the researcher to
126 use arbitrary algorithms for similarity searches and local QSAR (*Quantitative structure–*
127 *activity relationship*) modelling. Algorithms used within this study are described in the
128 following sections.

129 Neighbour identification

130 Utilizing this modularity, similarity calculations were based both on MolPrint2D finger-
131 prints and on CDK descriptors.

132 For MolPrint2D fingerprints chemical similarity between two compounds a and b is
133 expressed as the proportion between atom environments common in both structures
134 $A \cap B$ and the total number of atom environments $A \cup B$ (Jaccard/Tanimoto index).

$$sim = \frac{|A \cap B|}{|A \cup B|}$$

135 For CDK descriptors chemical similarity between two compounds a and b is expressed
136 as the cosine similarity between the descriptor vectors A for a and B for b .

$$sim = \frac{A \cdot B}{|A||B|}$$

137 Threshold selection is a trade-off between prediction accuracy (high threshold) and the
138 number of predictable compounds (low threshold). As it is in many practical cases
139 desirable to make predictions even in the absence of closely related neighbours, we follow
140 a tiered approach:

- 141 • First a similarity threshold of 0.5 is used to collect neighbours, to create a local
142 QSAR model and to make a prediction for the query compound. This are predic-
143 tions with *high confidence*.
- 144 • If any of these steps fails, the procedure is repeated with a similarity threshold
145 of 0.2 and the prediction is flagged with a warning that it might be out of the
146 applicability domain of the training data (*low confidence*).
- 147 • Similarity thresholds of 0.5 and 0.2 are the default values chosen by the software
148 developers and remained unchanged during the course of these experiments.

149 Compounds with the same structure as the query structure are automatically eliminated
150 from neighbours to obtain unbiased predictions in the presence of duplicates.

151 Local QSAR models and predictions

152 Only similar compounds (neighbours) above the threshold are used for local QSAR
153 models. In this investigation, we are using a weighted majority vote from the neigh-
154 bour's experimental data for mutagenicity classifications. Probabilities for both classes
155 (mutagenic/non-mutagenic) are calculated according to the following formula and the

156 class with the higher probability is used as prediction outcome.

$$p_c = \frac{\sum \text{sim}_{n,c}}{\sum \text{sim}_n}$$

157 p_c Probability of class c (e.g. mutagenic or non-mutagenic)

158 $\sum \text{sim}_{n,c}$ Sum of similarities of neighbours with class c

159 $\sum \text{sim}_n$ Sum of all neighbours

160 **Applicability domain**

161 The applicability domain (AD) of **lazar** models is determined by the structural diver-
162 sity of the training data. If no similar compounds are found in the training data no
163 predictions will be generated. Warnings are issued if the similarity threshold had to be
164 lowered from 0.5 to 0.2 in order to enable predictions. Predictions without warnings
165 can be considered as close to the applicability domain (*high confidence*) and predictions
166 with warnings as more distant from the applicability domain (*low confidence*). Quantita-
167 tive applicability domain information can be obtained from the similarities of individual
168 neighbours.

169 **Availability**

- 170 • Source code for this manuscript (GPL3): [https://git.in-silico.ch/lazar/tree/?h=](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
171 [mutagenicity-paper](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
- 172 • Crossvalidation experiments (GPL3): [https://git.in-silico.ch/lazar/tree/models/](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
173 [?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
- 174 • Pyrrolizidine alkaloid predictions (GPL3): [https://git.in-silico.ch/lazar/tree/](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
175 [predictions/?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
- 176 • Public web interface: <https://lazar.in-silico.ch>

177 **Tensorflow models**

178 **TODO: Philipp** Kannst Du bitte die folgenden Absätze ergänzen und die Vor-
179 gangsweise für MP2D/CDK bzw CV/PA Vorhersagen beschreiben.

180 **Random forests (*RF*)**

181 **Logistic regression (SGD) (*LR-sgd*)**

182 **Logistic regression (scikit) (*LR-scikit*)**

183 **Neural Nets (*NN*)**

184 **Support vector machines (*SVM*)**

185 **Validation**

186 10-fold cross-validation was used for all Tensorflow models.

187 **Availability**

188 Jupyter notebooks for these experiments can be found at the following locations

189 *Crossvalidation:*

- 190 • MolPrint2D fingerprints: [https://git.in-silico.ch/mutagenicity-paper/tree/](https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/mp2d/tensorflow)
191 crossvalidations/mp2d/tensorflow
192 • CDK descriptors: [https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/](https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/cdk/tensorflow)
193 cdk/tensorflow

194 *Pyrrolizidine alkaloids:*

- MolPrint2D fingerprints: <https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/tensorflow>
- CDK descriptors: <https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/tensorflow>
- CDK desc

Results

10-fold crossvalidations

Crossvalidation results are summarized in the following tables: Table 1 shows results with MolPrint2D descriptors and Table 2 with CDK descriptors.

Table 1: Summary of crossvalidation results with MolPrint2D descriptors (lazar-HC: lazar with high confidence, lazar-all: all lazar predictions, RF: random forests, LR-sgd: logistic regression (stochastic gradient descent), LR-scikit: logistic regression (scikit), NN: neural networks, SVM: support vector machines)

	lazar-HC	lazar-all	RF	LR-sgd	LR-scikit	NN	SVM
Accuracy	84	82	80	84	84	84	84
True positive rate	88	85	78	83	83	82	83
True negative rate	78	79	82	84	85	85	86
Positive predictive value	82	80	81	84	84	84	85
Negative predictive value	85	84	80	84	84	83	84
Nr. predictions	6300	7777	8303	8303	8303	8303	8303

Table 2: Summary of crossvalidation results with CDK descriptors (lazar-HC: lazar with high confidence, lazar-all: all lazar predictions, RF: random forests, LR-sgd: logistic regression (stochastic gradient descent), LR-scikit: logistic regression (scikit), NN: neural networks, SVM: support vector machines)

	lazar-HC	lazar-all	RF	LR-sgd	LR-scikit	NN	SVM
Accuracy	52	52	84	79	80	85	82
True positive rate	90	90	81	81	80	85	82
True negative rate	14	14	86	78	80	85	82
Positive predictive value	52	52	85	79	80	85	82
Negative predictive value	56	56	82	80	80	85	82
Nr. predictions	811	811	8077	8077	8077	8077	8077

Figure 1 depicts the position of all crossvalidation results in receiver operating characteristic (ROC) space.

Confusion matrices for all models are available from the git repository <https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/confusion-matrices/>, individual predictions can be found in <https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/predictions/>.

With exception of lazar/CDK all investigated algorithm/descriptor combinations give accuracies between (80 and 85%) which is equivalent to the experimental variability of the *Salmonella typhimurium* mutagenicity bioassay (80-85%, Benigni and Giuliani (1988)). Sensitivities and specificities are balanced in all of these models.

Pyrrolizidine alkaloid mutagenicity predictions

Mutagenicity predictions from all investigated models for 602 pyrrolizidine alkaloids (PAs) can be downloaded from <https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/pa-predictions.csv>. A visual representation of all PA predictions

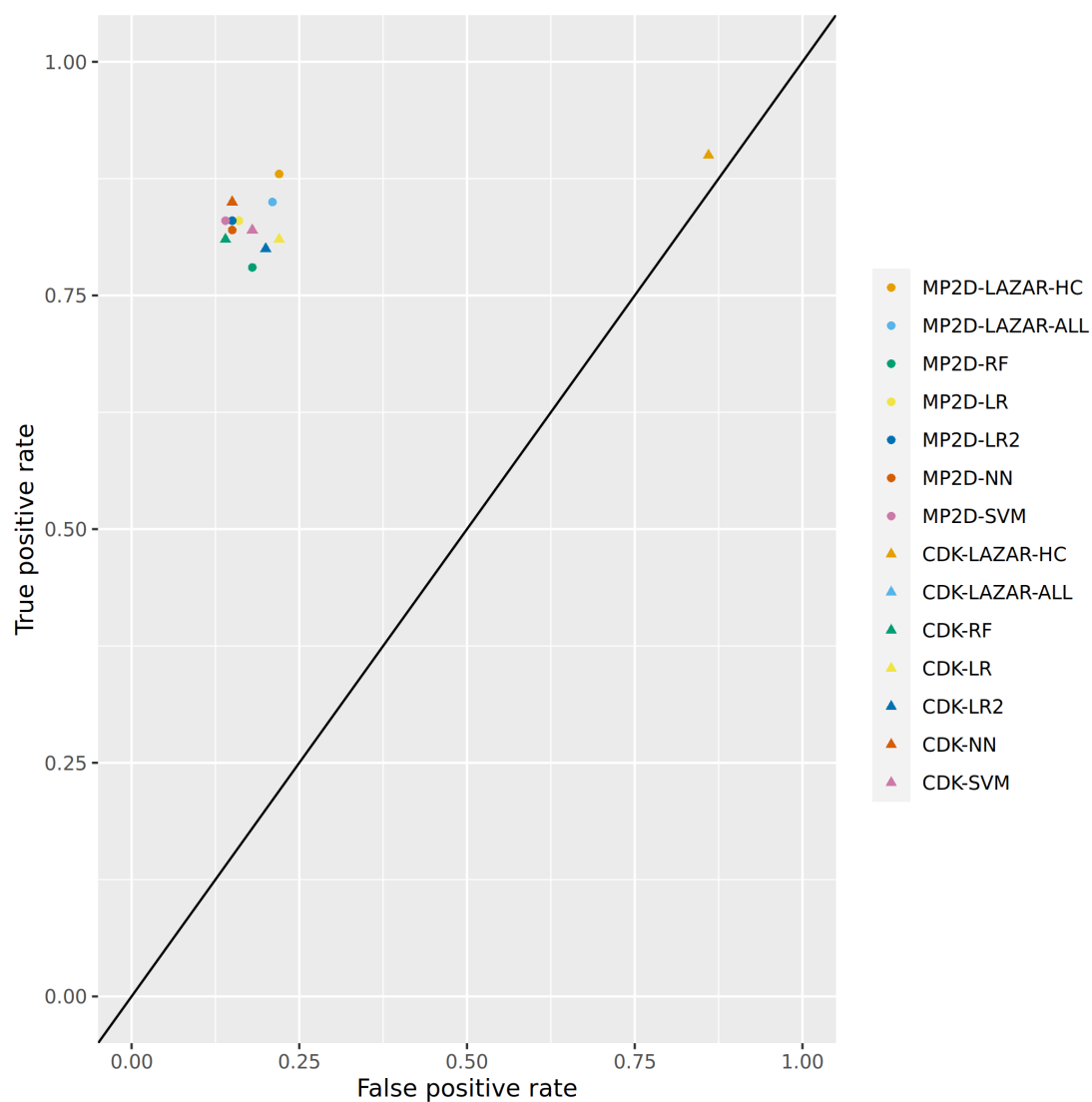


Figure 1: ROC plot of crossvalidation results (lazar-HC: lazar with high confidence, lazar-all: all lazar predictions, RF: random forests, LR-sgd: logistic regression (stochastic gradient descent), LR-scikit: logistic regression (scikit), NN: neural networks, SVM: support vector machines).

can be found at <https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/pa-predictions.pdf>.

Table 3 and Table 4 summarise the outcome of pyrrolizidine alkaloid predictions from all models with MolPrint2D and CDK descriptors.

Table 3: Summary of MolPrint2D pyrrolizidine alkaloid predictions

Model	mutagenic	non-mutagenic	Nr. predictions
lazar-all	20% (111)	80% (449)	93% (560)
lazar-HC	25% (76)	75% (225)	50% (301)
RF	5% (28)	95% (574)	100% (602)
LR-sgd	21% (127)	79% (475)	100% (602)
LR-scikit	20% (118)	80% (484)	100% (602)
NN	21% (124)	79% (478)	100% (602)
SVM	14% (82)	86% (520)	100% (602)

Table 4: Summary of CDK pyrrolizidine alkaloid predictions

Model	mutagenic	non-mutagenic	Nr. predictions
lazar-all	20% (111)	80% (449)	93% (560)
lazar-HC	25% (76)	75% (225)	50% (301)
RF	2% (10)	98% (592)	100% (602)
LR-sgd	16% (97)	84% (505)	100% (602)
LR-scikit	15% (88)	85% (514)	100% (602)
NN	25% (150)	75% (452)	100% (602)
SVM	3% (19)	97% (583)	100% (602)

Figure 2 - Figure 10 display the proportion of positive mutagenicity predictions from all

222 models for the different pyrrolizidine alkaloid groups.

223 For the visualisation of the position of pyrrolizidine alkaloids in respect to the train-
224 ing data set we have applied t-distributed stochastic neighbor embedding (t-SNE,
225 Maaten and Hinton (2008)) for MolPrint2D and CDK descriptors. t-SNE maps
226 each high-dimensional object (chemical) to a two-dimensional point, maintaining the
227 high-dimensional distances of the objects. Similar objects are represented by nearby
228 points and dissimilar objects are represented by distant points.

229 Figure 11 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training
230 data in MP2D space (Tanimoto/Jaccard similarity).

231 Figure 12 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training
232 data in CDK space (Euclidean similarity).

233 Discussion

234 Data

235 A new training dataset for *Salmonella* mutagenicity was created from three different
236 sources (Kazius, McGuire, and Bursi (2005), Hansen et al. (2009), EFSA (2016)). It
237 contains 8290 unique chemical structures, which is according to our knowledge the
238 largest public mutagenicity dataset presently available. The new training data can
239 be downloaded from [https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mutagenicity.csv)
240 [mutagenicity.csv](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mutagenicity.csv).

241 Algorithms

242 **lazar** is formally a *k-nearest-neighbor* algorithm that searches for similar structures
243 for a given compound and calculates the prediction based on the experimental data for
244 these structures. The QSAR literature calls such models frequently *local models*, because

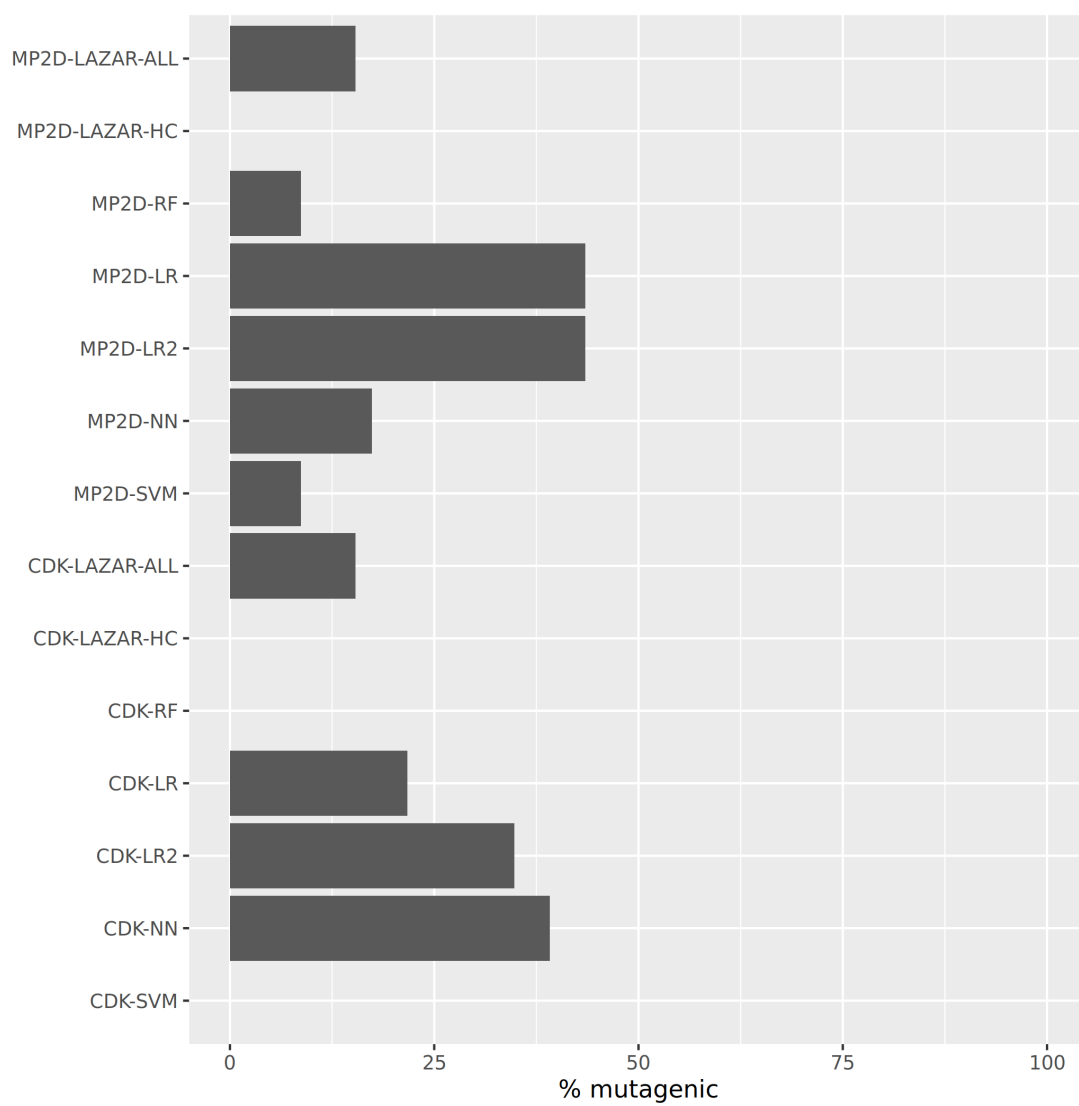


Figure 2: Summary of Dehydropyrrolizidine predictions

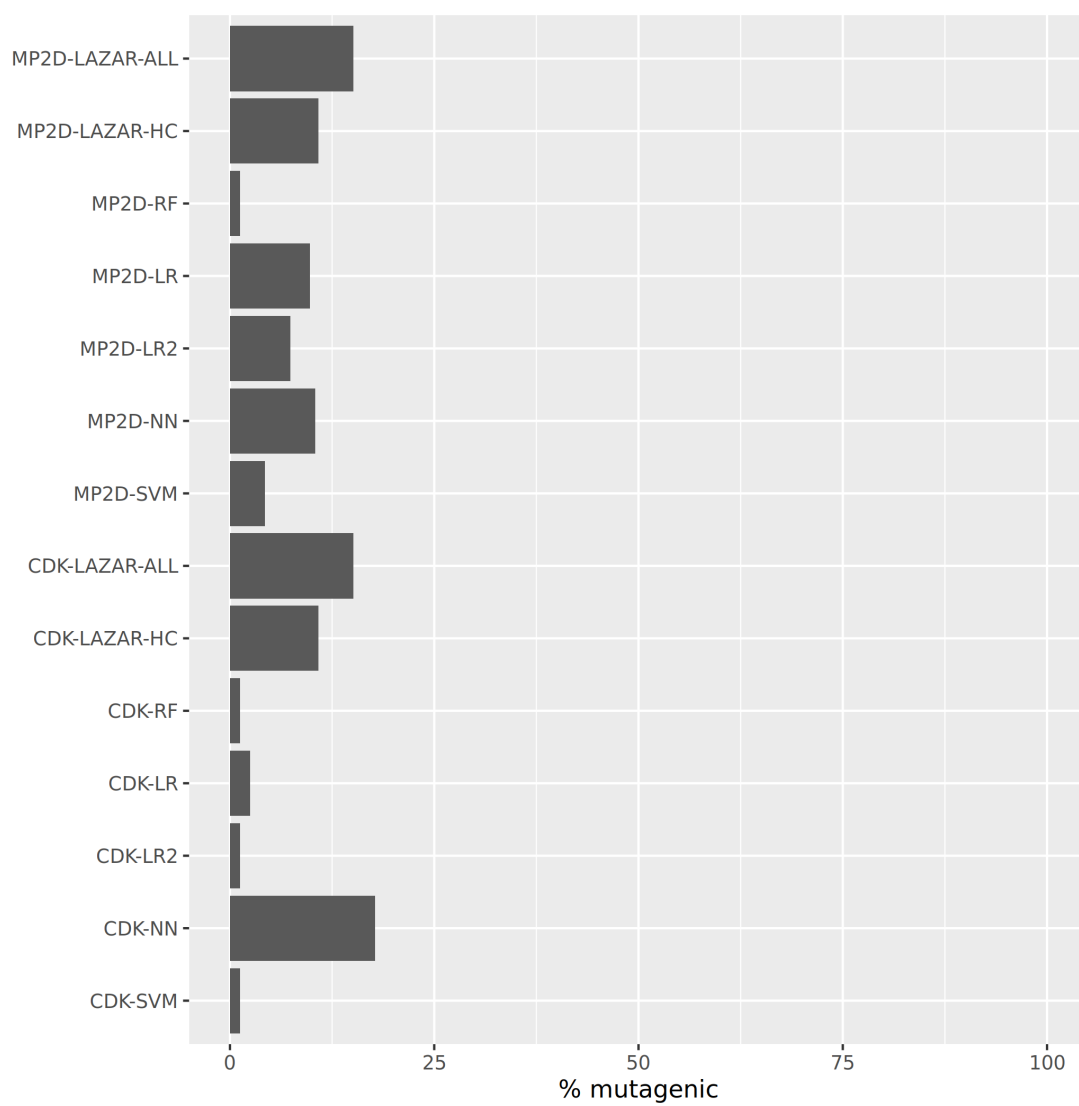


Figure 3: Summary of Diester predictions

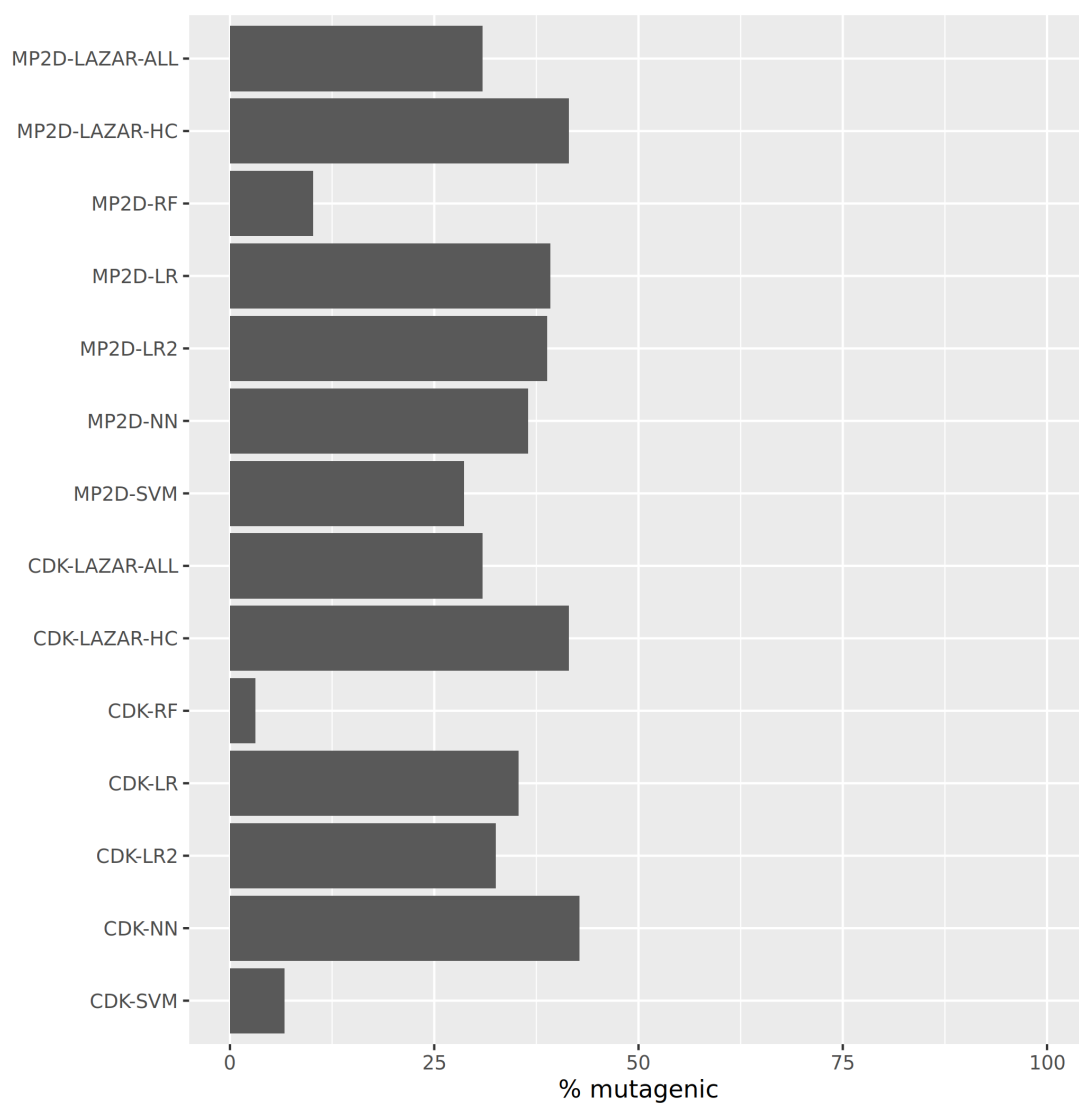


Figure 4: Summary of Macroyclic-diester predictions

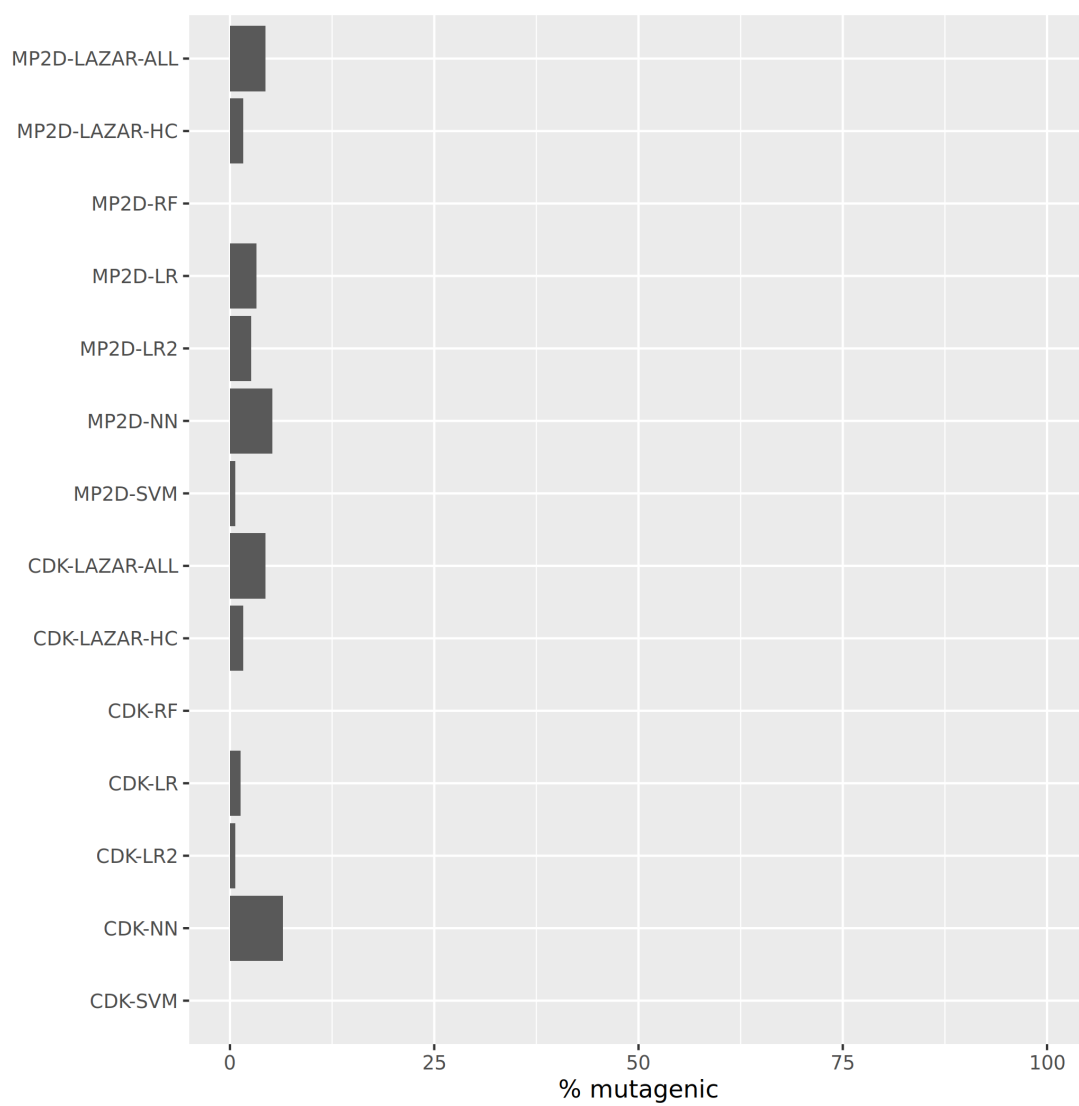


Figure 5: Summary of Monoester predictions

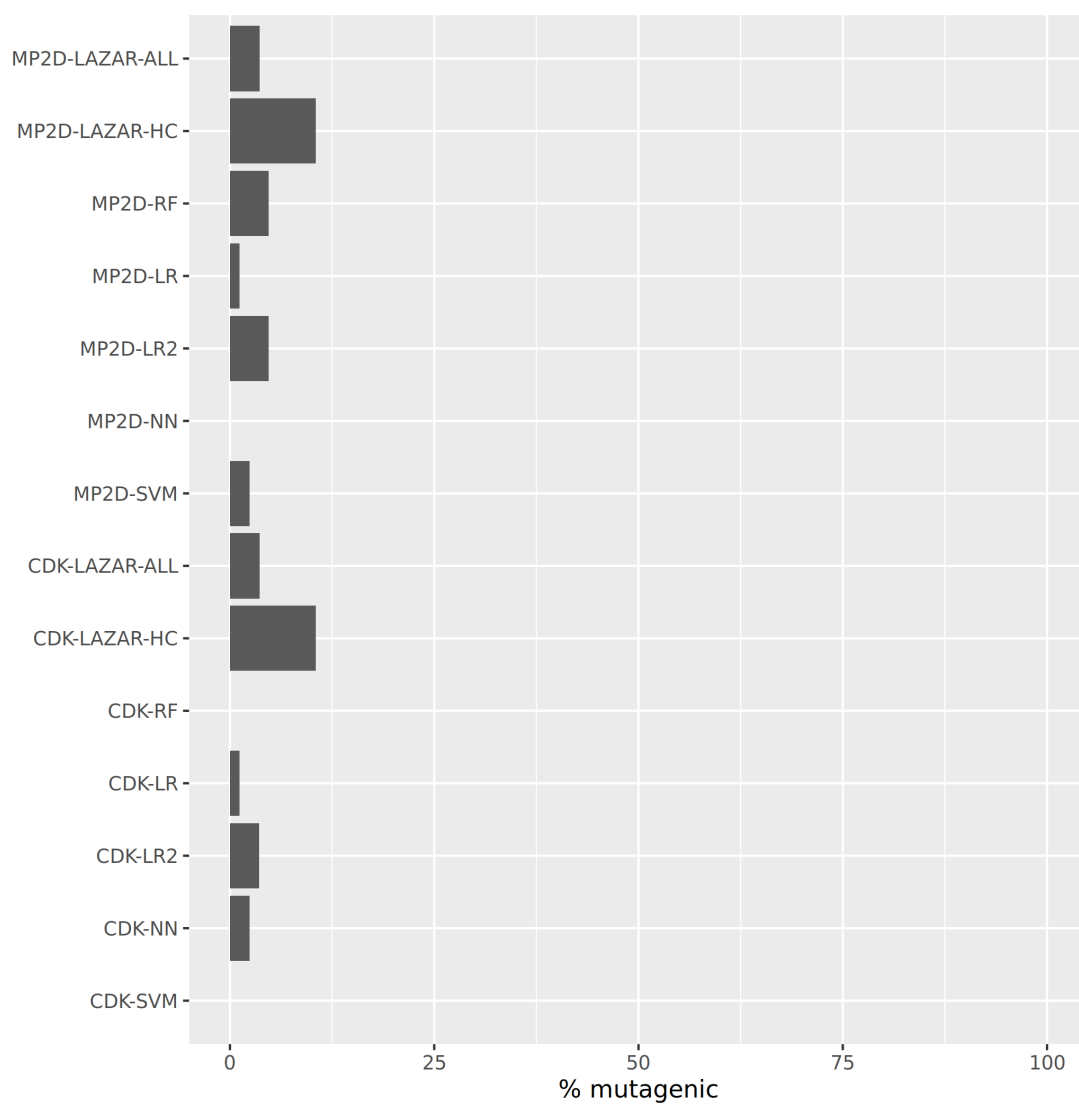


Figure 6: Summary of N-oxide predictions

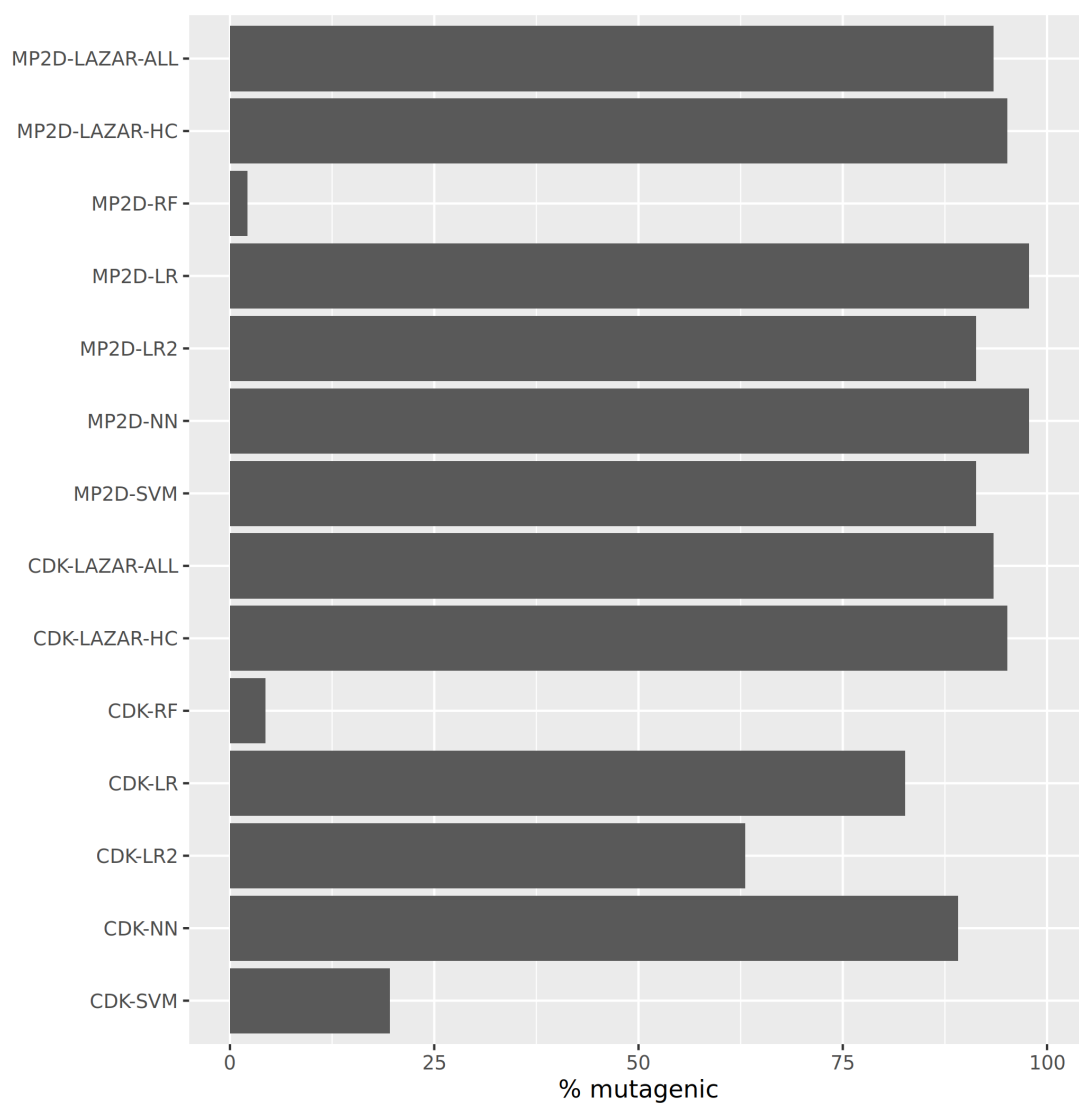


Figure 7: Summary of Otonecine predictions

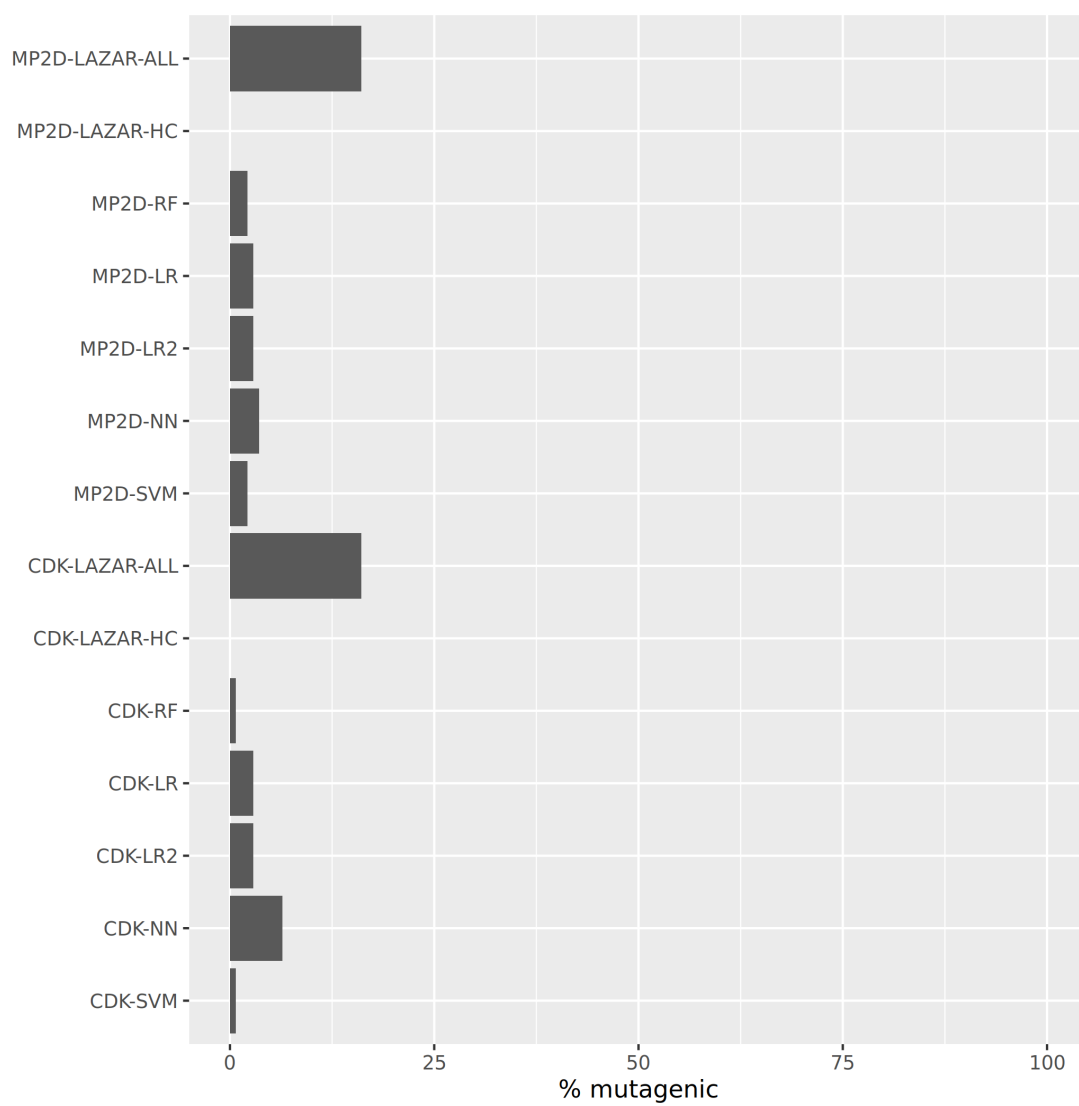


Figure 8: Summary of Platynecine predictions

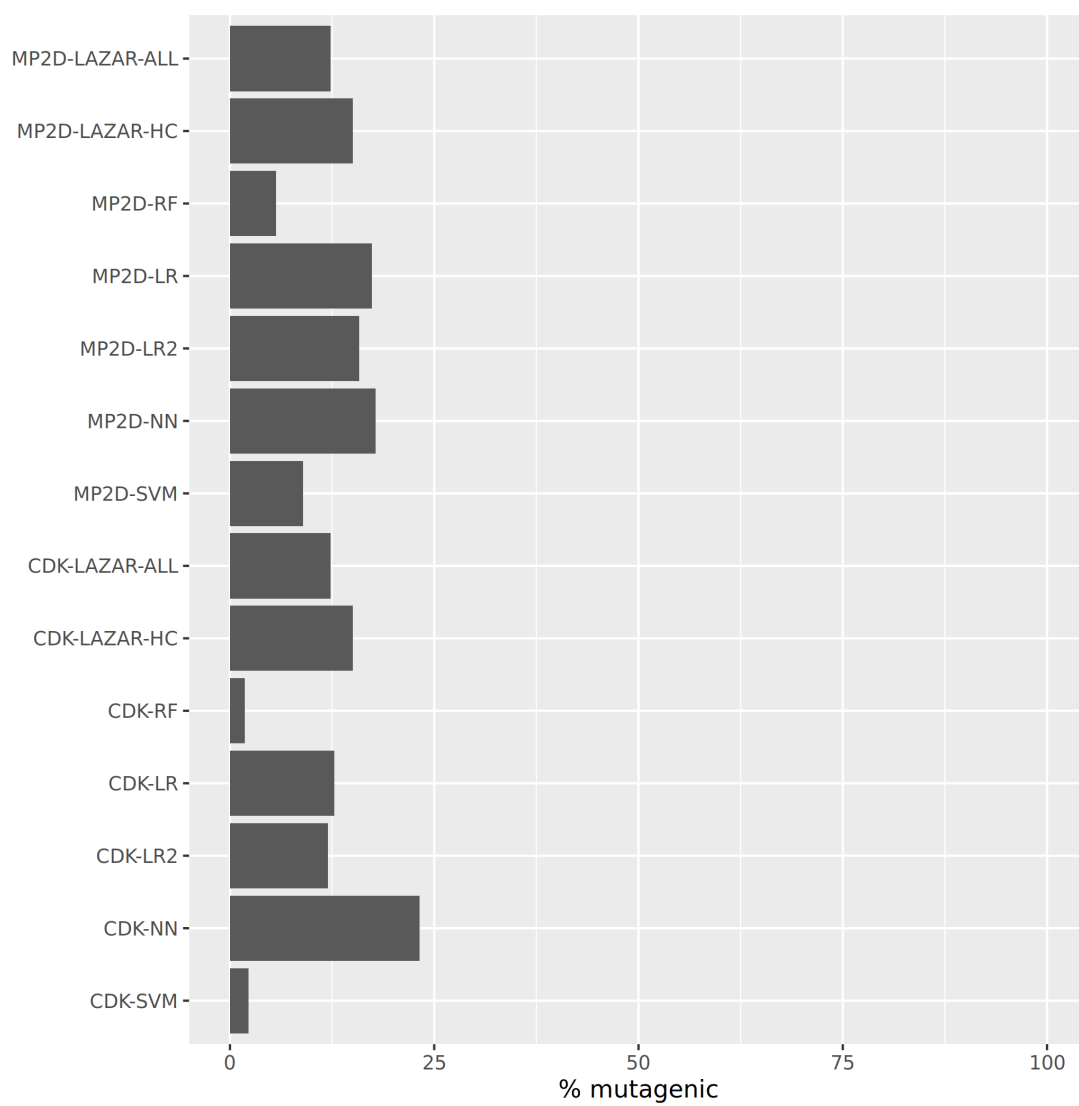


Figure 9: Summary of Retronecine predictions

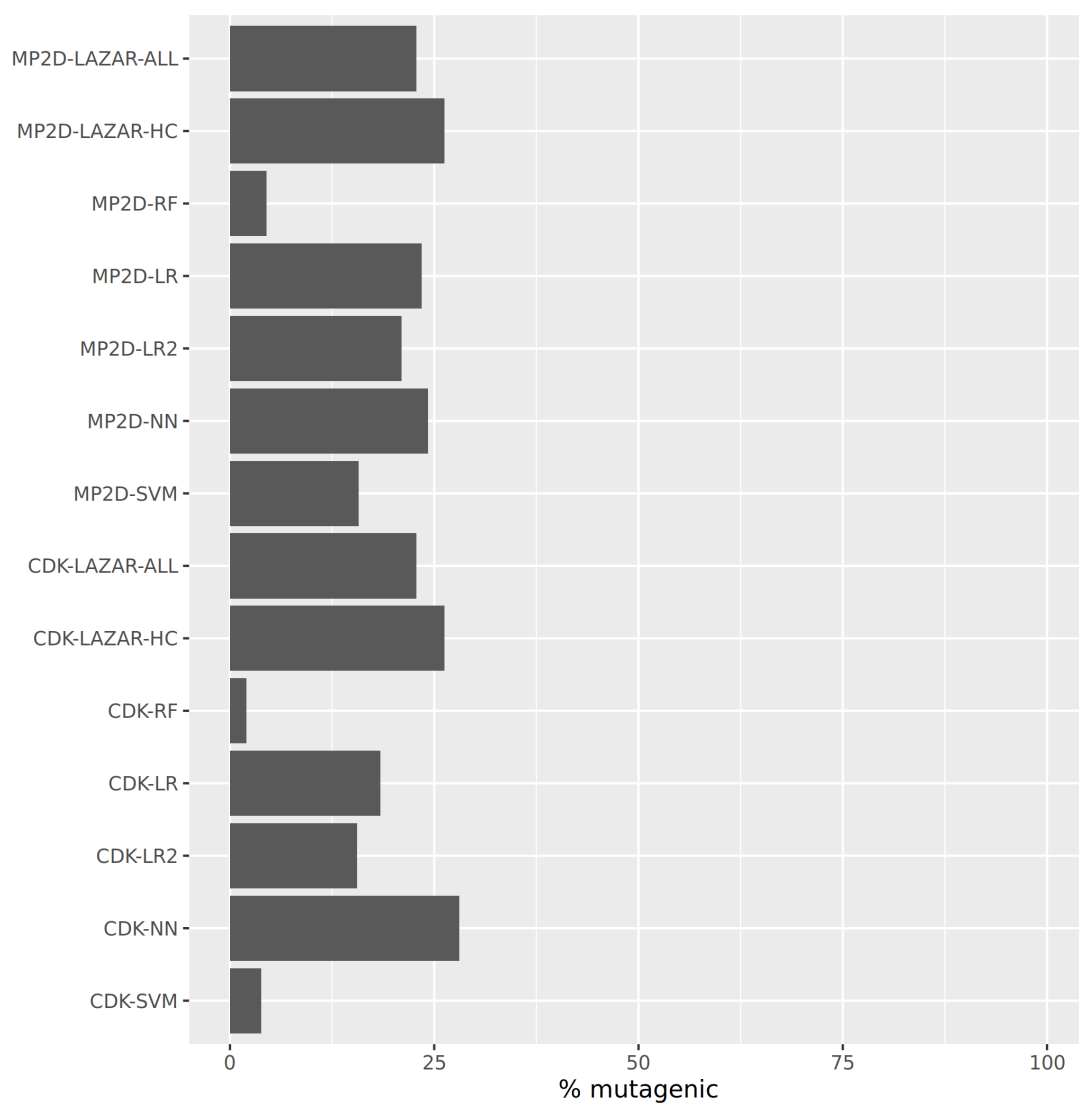


Figure 10: Summary of Tertiary PA predictions



Figure 11: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA)

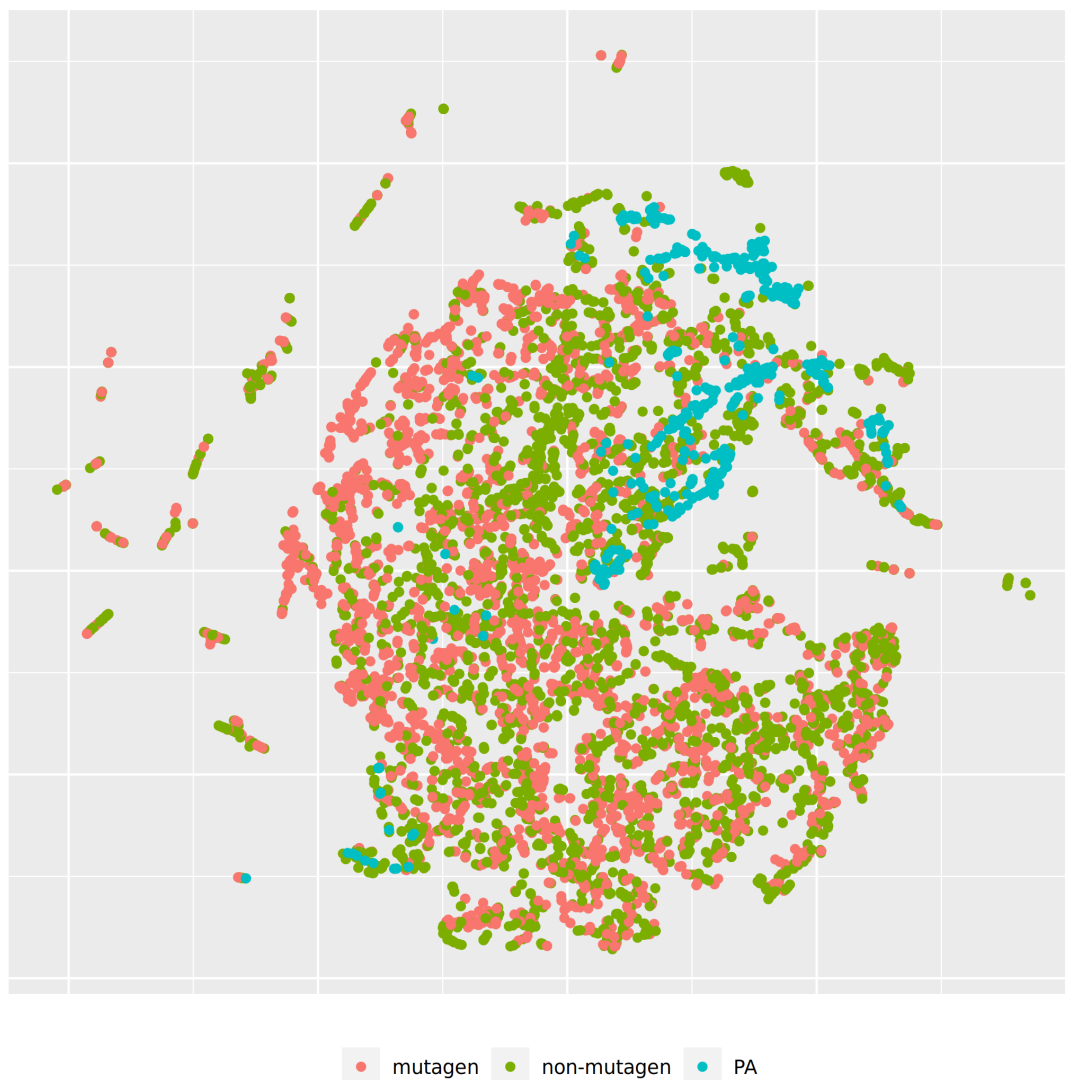


Figure 12: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA)

models are generated specifically for each query compound. The investigated tensorflow models are in contrast *global models*, i.e. a single model is used to make predictions for all compounds. It has been postulated in the past, that local models are more accurate, because they can account better for mechanisms, that affect only a subset of the training data.

Table 1, Table 2 and Figure 1 show that all models with the exception of *lazar*-CDK have similar crossvalidation accuracies that are comparable to the experimental variability of the *Salmonella typhimurium* mutagenicity bioassay (80-85% according to Benigni and Giuliani (1988)). All of these models have balanced sensitivity (true position rate) and specificity (true negative rate) and provide highly significant concordance with experimental data (as determined by McNemar’s Test). This is a clear indication that *in-silico* predictions can be as reliable as the bioassays. Given that the variability of experimental data is similar to model variability it is impossible to decide which model gives the most accurate predictions, as models with higher accuracies (e.g. NN-CDK) might just approximate experimental errors better than more robust models.

lazar predictions with CDK descriptors are a notable exception, as it has a much lower overall accuracy () than all other models. *lazar* uses basically a k-nearest-neighbor (with variable k) and it seems that CDK descriptors are not very well suited for chemical similarity calculations. We have confirmed this independently by validating k-nn models from the R *caret* package, which give also sub-par accuracies (data not shown).

Figure 12 is another indication that similarity calculations with CDK descriptors are not as useful as fingerprint based similarities, because it shows a less clearer separation between chemical classes and mutagens/non-mutagens than Figure 11. It seems that more complex models than simple k-nn are required to utilize CDK descriptors efficiently.

Our results do not support the assumption that local models are superior to global models for classification purposes. For regression models (lowest observed effect level)

we have found however that local models may outperform global models (Helma et al. (2018)) with accuracies similar to experimental variability.

Descriptors

This study uses two types of descriptors for the characterisation of chemical structures: *MolPrint2D* fingerprints (MP2D, Bender et al. (2004)) use atom environments (i.e. connected atom types for all atoms in a molecule) as molecular representation, which resembles basically the chemical concept of functional groups. MP2D descriptors are used to determine chemical similarities in the default **lazar** settings, and previous experiments have shown, that they give more accurate results than predefined fingerprints (e.g. MACCS, FP2-4).

Chemistry Development Kit (CDK, Willighagen E. L. (2017)) descriptors were calculated with the PaDEL graphical interface (Yap (2011)). They include 1D and 2D topological descriptors as well as physical-chemical properties.

With exception of **lazar** all investigated algorithms obtained models within the experimental variability for both types of descriptors. As discussed before CDK descriptors seem to be less suitable for chemical similarity calculations than MolPrint2D descriptors.

Given that similar predictive accuracies are obtainable from both types of descriptors the choice depends more on practical considerations:

MolPrint2D fragments can be calculated very efficiently for every well defined chemical structure with OpenBabel (O’Boyle et al. (2011)). CDK descriptor calculations are in contrast much more resource intensive and may fail for a significant number of compounds (from 8290).

MolPrint2D fragments are generated dynamically from chemical structures and can be used to determine if a compound contains structural features that are absent in training

data. This feature can be used to determine applicability domains. CDK descriptors contain in contrast a predefined set of descriptors with unknown toxicological relevance. MolPrint2D fingerprints can be represented very efficiently as sets of features that are present in a given compound which makes similarity calculations very efficient. Due to the large number of substructures present in training compounds, they lead however to large and sparsely populated datasets, if they have to be expanded to a binary matrix (e.g. as input for tensorflow models). CDK descriptors contain in contrast in every case matrices with 1442 columns.

Pyrrolizidine alkaloid mutagenicity predictions

Figure 2 - Figure 10 show a clear differentiation between the different pyrrolizidine alkaloid groups. The largest proportion of mutagenic predictions was observed for Otonecines 72% (458/634), the lowest for Monoesters 2% (45/1940) and N-Oxides 2% (27/1044).

Although most of the models show similar accuracies, sensitivities and specificities in crossvalidation experiments some of the models (MPD-RF, CDK-RF and CDK-SVM) predict a lower number of mutagens (2-5%) than the majority of the models (14-25% Table 3, Table 4, Figure 2 - Figure 10).

From a practical point we still have to face the question, how to choose model predictions, if no experimental data is available (we found two PAs in the training data, but this number is too low, to draw any general conclusions).

TODO: Verena Hier ist ein alter Text von Dir zum Recylen:

Necic acid

The rank order of the necic acid is comparable in the four models considered (LAZAR, RF and DL (R-project and Tensorflow). PAs from the monoester type had the lowest genotoxic potential, followed by PAs from the open-ring diester type. PAs with macro-

cyclic diesters had the highest genotoxic potential. The result fit well with current state of knowledge: in general, PAs, which have a macrocyclic diesters as necic acid, are considered more toxic than those with an open-ring diester or monoester EFSA 2011Fu et al. 2004Ruan et al. 2014b(; ;).

Necine base

The rank order of necine base is comparable in LAZAR, RF, and DL (R-project) models: with platynecine being less or as genotoxic as retronecine, and otonecine being the most genotoxic. In the Tensorflow-generate DL model, platynecine also has the lowest genotoxic probability, but are then followed by the otonecines and last by retronecine. These results partly correspond to earlier published studies. Saturated PAs of the platynecine-type are generally accepted to be less or non-toxic and have been shown in *in vitro* experiments to form no DNA-adducts Xia et al. 2013(). Therefore, it is striking, that 1,2-unsaturated PAs of the retronecine-type should have an almost comparable genotoxic potential in the LAZAR and DL (R-project) model. In literature, otonecine-type PAs were shown to be more toxic than those of the retronecine-type Li et al. 2013().

Modifications of necine base

The group-specific results of the Tensorflow-generated DL model appear to reflect the expected relationship between the groups: the low genotoxic potential of *N*-oxides and the highest potential of dehydropyrrolizidines Chen et al. 2010().

In the LAZAR model, the genotoxic potential of dehydropyrrolizidines (DHP) (using the extended AD) is comparable to that of tertiary PAs. Since, DHP is regarded as the toxic principle in the metabolism of PAs, and known to produce protein- and DNA-adducts Chen et al. 2010(), the LAZAR model did not meet this expectation it predicted the majority of DHP as being not genotoxic. However, the following issues need to be considered. On the one hand, all DHP were outside of the stricter AD of 0.5. This indicates that in general, there might be a problem with the AD. In addition, DHP has

345 two unsaturated double bonds in its necine base, making it highly reactive. DHP and
346 other comparable molecules have a very short lifespan, and usually cannot be used in *in*
347 *vitro* experiments. This might explain the absence of suitable neighbours in LAZAR.

348 Furthermore, the probabilities for this substance groups needs to be considered, and
349 not only the consolidated prediction. In the LAZAR model, all DHPs had probabilities
350 for both outcomes (genotoxic and not genotoxic) mainly below 30%. Additionally, the
351 probabilities for both outcomes were close together, often within 10% of each other. The
352 fact that for both outcomes, the probabilities were low and close together, indicates a
353 lower confidence in the prediction of the model for DHPs.

354 In the DL (R-project) and RF model, *N*-oxides have a by far more genotoxic potential
355 than tertiary PAs or dehydropyrrolizidines. As PA *N*-oxides are easily conjugated for
356 extraction, they are generally considered as detoxification products, which are *in vivo*
357 quickly renally eliminated Chen et al. 2010(). On the other hand, *N*-oxides can be also
358 back-transformed to the corresponding tertiary PA Wang et al. 2005(). Therefore, it
359 may be questioned, whether *N*-oxides themselves are generally less genotoxic than the
360 corresponding tertiary PAs. However, in the groups of modification of the necine base,
361 dehydropyrrolizidine, the toxic principle of PAs, should have had the highest genotoxic
362 potential. Taken together, the predictions of the modifications of the necine base from
363 the LAZAR, RF and R-generated DL model cannot - in contrast to the Tensorflow DL
364 model - be considered as reliable.

365 Overall, when comparing the prediction results of the PAs to current published knowl-
366 edge, it can be concluded that the performance of most models was low to moderate.
367 This might be contributed to the following issues:

- 368 1. In the LAZAR model, only 26.6% PAs were within the stricter AD. With the
369 extended AD, 92.3% of the PAs could be included in the prediction. Even though
370 the Jaccard distance between the training dataset and the PA dataset for the RF,

371 SVM, and DL (R-project and Tensorflow) models was small, suggesting a high
372 similarity, the LAZAR indicated that PAs have only few local neighbours, which
373 might adversely affect the prediction of the mutagenic potential of PAs.

374 2. All above-mentioned models were used to predict the mutagenicity of PAs. PAs
375 are generally considered to be genotoxic, and the mode of action is also known.
376 Therefore, the fact that some models predict the majority of PAs as not genotoxic
377 seems contradictory. To understand this result, the basis, the training dataset, has
378 to be considered. The mutagenicity of in the training dataset are based on data of
379 mutagenicity in bacteria. There are some studies, which show mutagenicity of PAs
380 in the AMES test Chen et al. 2010(). Also, Rubiolo et al. (1992) examined several
381 different PAs and several different extracts of PA-containing plants in the AMES
382 test. They found that the AMES test was indeed able to detect mutagenicity of
383 PAs, but in general, appeared to have a low sensitivity. The pre-incubation phase
384 for metabolic activation of PAs by microsomal enzymes was the sensitivity-limiting
385 step. This could very well mean that this is also reflected in the QSAR models.

386 Conclusions

387 A new public *Salmonella* mutagenicity training dataset with 8309 compounds was cre-
388 ated and used it to train `lazar` and Tensorflow models with MolPrint2D and CDK
389 descriptors.

390 References

391 Bender, Andreas, Hamse Y. Mussa, Robert C. Glen, and Stephan Reiling. 2004. "Molec-
392 ular Similarity Searching Using Atom Environments, Information-Based Feature Selec-
393 tion, and a Naïve Bayesian Classifier." *Journal of Chemical Information and Computer*

394 *Sciences* 44 (1): 170–78. <https://doi.org/10.1021/ci034207y>.

395 Benigni, R., and A. Giuliani. 1988. “Computer-assisted Analysis of Interlaboratory
396 Ames Test Variability.” *Journal of Toxicology and Environmental Health* 25 (1): 135–48.
397 <https://doi.org/10.1080/15287398809531194>.

398 EFSA. 2011. “Scientific Opinion on Pyrrolizidine Alkaloids in Food and Feed.” *EFSA*
399 *Journal*, no. 9: 1–134.

400 ———. 2016. “Guidance on the Establishment of the Residue Definition for Dietary
401 Assessment: EFSA Panel on Plant Protect Products and Their Residues (PPR).” *EFSA*
402 *Journal*, no. 14: 1–12.

403 Hansen, Katja, Sebastian Mika, Timon Schroeter, Andreas Sutter, Antonius ter Laak,
404 Thomas Steger-Hartmann, Nikolaus Heinrich, and Klaus-Robert Müller. 2009. “Bench-
405 mark Data Set for in Silico Prediction of Ames Mutagenicity.” *Journal of Chemical*
406 *Information and Modeling* 49 (9): 2077–81. <https://doi.org/10.1021/ci900161g>.

407 Helma, Christoph, David Vorgrimmler, Denis Gebele, Martin Gütlein, Barbara Engeli,
408 Jürg Zarn, Benoit Schilter, and Elena Lo Piparo. 2018. “Modeling Chronic Toxicity: A
409 Comparison of Experimental Variability with (Q)SAR/Read-Across Predictions.” *Fron-*
410 *tiers in Pharmacology*, no. 9: 413.

411 Kazius, J., R. McGuire, and R. Bursi. 2005. “Derivation and Validation of Toxicophores
412 for Mutagenicity Prediction.” *J Med Chem*, no. 48: 312–20.

413 Maaten, L. J. P. van der, and G. E. Hinton. 2008. “Visualizing Data Using T-Sne.”
414 *Journal of Machine Learning Research*, no. 9: 2579–2605.

415 Mattocks, AR. 1986. *Chemistry and Toxicology of Pyrrolizidine Alkaloids*. Academic
416 Press.

417 O’Boyle, Noel, Michael Banck, Craig James, Chris Morley, Tim Vandermeersch, and

- 418 Geoffrey Hutchison. 2011. "Open Babel: An open chemical toolbox." *J. Cheminf.* 3 (1):
419 33. <https://doi.org/doi:10.1186/1758-2946-3-33>.
- 420 Schöning, Verena, Felix Hammann, Mark Peinl, and Jürgen Drewe. 2017. "Editor's
421 Highlight: Identification of Any Structure-Specific Hepatotoxic Potential of Different
422 Pyrrolizidine Alkaloids Using Random Forests and Artificial Neural Networks." *Toxicol.*
423 *Sci.*, no. 160: 361–70.
- 424 Willighagen E. L., Alvarsson J. et al., Mayfield J. W. 2017. "The Chemistry Devel-
425 opment Kit (Cdk) V2.0: Atom Typing, Depiction, Molecular Formulas, and Substruc-
426 ture Searching." *J. Cheminform.*, no. 9(33). [https://doi.org/https://doi.org/10.1186/](https://doi.org/https://doi.org/10.1186/s13321-017-0220-4)
427 [s13321-017-0220-4](https://doi.org/https://doi.org/10.1186/s13321-017-0220-4).
- 428 Yap, CW. 2011. "PaDEL-Descriptor: An Open Source Software to Calculate Molecular
429 Descriptors and Fingerprints." *Journal of Computational Chemistry*, no. 32: 1466–74.