

1 A comparison of nine machine learning mutagenicity models
2 and their application for predicting pyrrolizidine alkaloids

3 Christoph Helma^{*1}, Verena Schöning⁵, Jürgen Drewe^{*2,4}, and Philipp Boss³

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

5 ²Max Zeller Söhne AG, Seeblickstrasse 4, 8590 Romanshorn, Switzerland

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

8 ⁴Clinical Pharmacology, Department of Pharmaceutical Sciences, University Hospital
9 Basel, University of Basel, Petersgraben 4, 4031 Basel, Switzerland

10 ⁵Clinical Pharmacology and Toxicology, Department of General Internal Medicine,
11 University Hospital Bern, University of Bern, Inselspital, 3010 Bern, Switzerland

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

13 Jürgen Drewe <juergendrewe@zellerag.ch>

14 Random forest, support vector machine, logistic regression, neural
15 networks and k-nearest neighbor (**lazar**) algorithms, were applied to new
16 *Salmonella* mutagenicity dataset with 8290 unique chemical structures
17 utilizing MolPrint2D and Chemistry Development Kit (CDK) descriptors.
18 Crossvalidation accuracies of all investigated models ranged from 80-85%
19 which is comparable with the interlaboratory variability of the *Salmonella*
20 mutagenicity assay. Pyrrolizidine alkaloid predictions showed a clear
21 distinction between chemical groups, where otonecines had the highest
22 proportion of positive mutagenicity predictions and monoesters the lowest.

23 Introduction

24 The assessment of mutagenicity is an important part in the safety assessment of chemical
25 structures, because genomic changes may lead to cancer and germ cells damage. The
26 *Salmonella typhimurium* bacterial reverse mutation test (Ames test) is capable to identify
27 substances that cause mutations (e.g., base-pair substitutions, frameshifts, insertions,
28 deletions) and is generally used as the first step in genotoxicity and carcinogenicity
29 assessments.

30 Computer based (*in silico*) mutagenicity predictions can be used in the early screening of
31 novel compounds (e.g. drug candidates), but they are also gaining regulatory acceptance
32 e.g. for the registration of industrial chemicals within REACH ((ECHA) (2017)) or
33 the assessment of impurities in pharmaceuticals (ICH M7 guideline, Harmonisation of
34 Technical Requirements for Pharmaceuticals for Human Use International Council for
35 Harmonisation of Technical Requirements for Pharmaceuticals for Human Use (ICH)
36 (2017)).

37 Currently, *Salmonella* mutagenicity is the toxicological endpoint with the largest amount
38 of public data for almost 10000 structures, whereas datasets for other endpoints contain
39 typically only a few hundred compounds. The Ames test itself is relatively reproducible
40 with an interlaboratory variability of 80-85% (Piegorsch and Zeiger (1991)).

41 This makes the development of mutagenicity models also interesting from a computa-
42 tional chemistry and machine learning point of view. The relatively large amount of
43 public data reduces the probability of chance effects due to small sample sizes and the
44 reliability of the underlying assay reduces the risk of overfitting experimental errors.

45 Within this study we attempted

- 46 • to generate a new public mutagenicity training dataset, by combining the most
47 comprehensive public datasets

- to compare the performance of MolPrint2D (*MP2D*) fingerprints with Chemistry Development Kit (*CDK*) descriptors for mutagenicity predictions
- 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*)

In order to highlight potentials and problems with the application of mutagenicity models to compounds with limited experimental data we decided to apply these mutagenicity models to Pyrrolizidine alkaloids (PAs).

Pyrrolizidine alkaloids (PAs) are characteristic metabolites of some plant families, mainly: *Asteraceae*, *Boraginaceae*, *Fabaceae* and *Orchidaceae* (Hartmann and Witte (1995), Langel, Ober, and Pelser (2011)) and form a powerful defence mechanism against herbivores. PAs are heterocyclic ester alkaloids composed of a necine base (two fused five-membered rings joined by a single nitrogen atom) and a necic acid (one or two carboxylic ester arms), occurring principally in two forms, tertiary base PAs and PA N-oxides. Several *in vitro* studies have shown the mutagenic potential of PAs, which seems highly dependent on structure of necine base and necic acid (Hadi et al. (2021); Allemang et al. (2018), Louisse et al. (2019)). However, due to limited availability of pure substances, only a limited number of PAs have been investigated with regards to their structure-specific mutagenicity. To overcome this bottleneck, the prediction of structure-specific mutagenic potential of PAs with different machine learning models could provide further insight into the mechanisms.

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*, Weininger, Weininger, and Weininger (1989)) 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://>

93 //git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mutagenicity.csv.

94 **Pyrrolizidine alkaloid (PA) dataset**

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

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

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

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

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

- 111 • N-oxide-type (84 compounds)
- 112 • Tertiary-type (PAs which were neither from the N-oxide- nor DHP-type, 495 com-
113 pounds)
- 114 • Dehydropyrrolizidine-type (pyrrolic ester, 23 compounds)

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

- 116 • Monoester-type (154 compounds)

- Open-ring diester-type (163 compounds)
- Macrocyclic diester-type (255 compounds)

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

Descriptors

MolPrint2D (*MP2D*) fingerprints

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

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

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

Training data:

- sparse representation (<https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/fingerprints.mp2d>)
- descriptor matrix (<https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mp2d/mutagenicity-fingerprints.csv.gz>)

Pyrrolizidine alkaloids:

- sparse representation (<https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/>)

140 mp2d/fingerprints.mp2d)
141 • 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)
142 mp2d/pa-fingerprints.csv.gz)

143 Chemistry Development Kit (*CDK*) descriptors

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

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

152 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)
153 tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv.

154 The same procedure was applied for the pyrrolizidine dataset yielding descriptors for
155 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)
156 ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/PA-Padel-2D_m2.csv.

157 Algorithms

158 **lazar**

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

- 162 • searches in a database for similar structures (neighbours) with experimental data,

- builds a local QSAR model with these neighbours and
- uses this model to predict the unknown activity of the query compound.

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

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

Feature preprocessing

MolPrint2D features were used without preprocessing. Near zero variance and strongly correlated CDK descriptors were removed and the remaining descriptor values were centered and scaled. Preprocessing was performed with the R **caret** `preProcess` function using the methods “nzv”, “corr”, “center” and “scale” with default settings.

Neighbour identification

Utilizing this modularity, similarity calculations were based both on MolPrint2D fingerprints and on CDK descriptors.

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

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

For CDK descriptors chemical similarity between two compounds a and b is expressed

183 as the cosine similarity between the descriptor vectors A for a and B for b .

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

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

- 188 • First a similarity threshold of 0.5 (MP2D/Tanimoto) or 0.9 (CDK/Cosine) is used
189 to collect neighbours, to create a local QSAR model and to make a prediction for
190 the query compound. This are predictions with *high confidence*.
- 191 • If any of these steps fails, the procedure is repeated with a similarity threshold of
192 0.2 (MP2D/Tanimoto) or 0.7 (CDK/Cosine) and the prediction is flagged with a
193 warning that it might be out of the applicability domain of the training data (*low*
194 *confidence*).
- 195 • These similarity thresholds are the default values chosen by software developers
196 and remained unchanged during the course of these experiments.

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

199 **Local QSAR models and predictions**

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

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

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

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

206 $\sum \text{sim}_{n,c}$ Sum of similarities of neighbours with class c

207 $\sum \text{sim}_n$ Sum of all neighbours

208 **Applicability domain**

209 The applicability domain (AD) of **1azar** models is determined by the structural diver-
210 sity of the training data. If no similar compounds are found in the training data no
211 predictions will be generated. Warnings are issued if the similarity threshold had to be
212 lowered from 0.5 to 0.2 in order to enable predictions. Predictions without warnings
213 can be considered as close to the applicability domain (*high confidence*) and predictions
214 with warnings as more distant from the applicability domain (*low confidence*). Quantita-
215 tive applicability domain information can be obtained from the similarities of individual
216 neighbours.

217 **Validation**

218 10-fold cross validation was performed for model evaluation.

219 **Pyrrolizidine alkaloid predictions**

220 For the prediction of pyrrolizidine alkaloids models were generated with the MP2D and
221 CDK training datasets. The complete feature set was used for MP2D predictions, for
222 CDK predictions the intersection between training and pyrrolizidine alkaloid features
223 was used.

224 **Availability**

- 225 • Source code for this manuscript (GPL3): [https://git.in-silico.ch/lazar/tree/?h=](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
226 [mutagenicity-paper](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
- 227 • Crossvalidation experiments (GPL3): [https://git.in-silico.ch/lazar/tree/models/](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
228 [?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
- 229 • Pyrrolizidine alkaloid predictions (GPL3): [https://git.in-silico.ch/lazar/tree/](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
230 [predictions/?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
- 231 • Public web interface: <https://lazar.in-silico.ch>

232 **Tensorflow models**

233 **Feature Preprocessing**

234 For preprocessing of the CDK features we used a quantile transformation to a uniform
235 distribution. MP2D features were not preprocessed.

236 **Random forests (*RF*)**

237 For the random forest classifier we used the parameters `n_estimators=1000` and
238 `max_leaf_nodes=200`. For the other parameters we used the scikit-learn default values.

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

240 For the logistic regression we used an ensemble of five trained models. For each model
241 we used a batch size of 64 and trained for 50 epochs. As an optimizer ADAM was chosen.
242 For the other parameters we used the tensorflow default values.

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

244 For the logistic regression we used as parameters the scikit-learn default values.

245 **Neural Nets (*NN*)**

246 For the neural network we used an ensemble of five trained models. For each model we
247 used a batch size of 64 and trained for 50 epochs. As an optimizer ADAM was chosen.
248 The neural network had 4 hidden layers with 64 nodes each and a ReLu activation
249 function. For the other parameters we used the tensorflow default values.

250 **Support vector machines (*SVM*)**

251 We used the SVM implemented in scikit-learn. We used the parameters kernel='rbf',
252 gamma='scale'. For the other parameters we used the scikit-learn default values.

253 **Validation**

254 10-fold cross-validation was used for all Tensorflow models.

255 **Pyrrolizidine alkaloid predictions**

256 For the prediction of pyrrolizidine alkaloids we trained the model described above on
257 the training data. For training and prediction only the features were used that were in
258 the intersection of features from the training data and the pyrrolizidine alkaloids.

259 **Availability**

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

261 *Crossvalidation:*

- 262 • MolPrint2D fingerprints: [https://git.in-silico.ch/mutagenicity-paper/tree/](https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/mp2d/tensorflow)
263 crossvalidations/mp2d/tensorflow

- CDK descriptors: <https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/cdk/tensorflow>
- 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 | 89 | 85 | 78 | 83 | 83 | 82 | 83 |
| True negative rate | 78 | 78 | 82 | 84 | 85 | 85 | 86 |
| Positive predictive value | 83 | 80 | 81 | 84 | 84 | 84 | 85 |
| Negative predictive value | 86 | 84 | 80 | 84 | 84 | 83 | 84 |
| Nr. predictions | 5864 | 7782 | 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 | 85 | 82 | 84 | 79 | 80 | 85 | 82 |
| True positive rate | 87 | 84 | 81 | 81 | 80 | 85 | 82 |
| True negative rate | 82 | 80 | 86 | 78 | 80 | 85 | 82 |
| Positive predictive value | 85 | 81 | 85 | 79 | 80 | 85 | 82 |
| Negative predictive value | 85 | 82 | 82 | 80 | 80 | 85 | 82 |
| Nr. predictions | 4872 | 7353 | 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/>.

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%, Piegorsch and Zeiger (1991)). Sensitivities and specificities are balanced in all of these models.

Pyrrolizidine alkaloid mutagenicity predictions

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

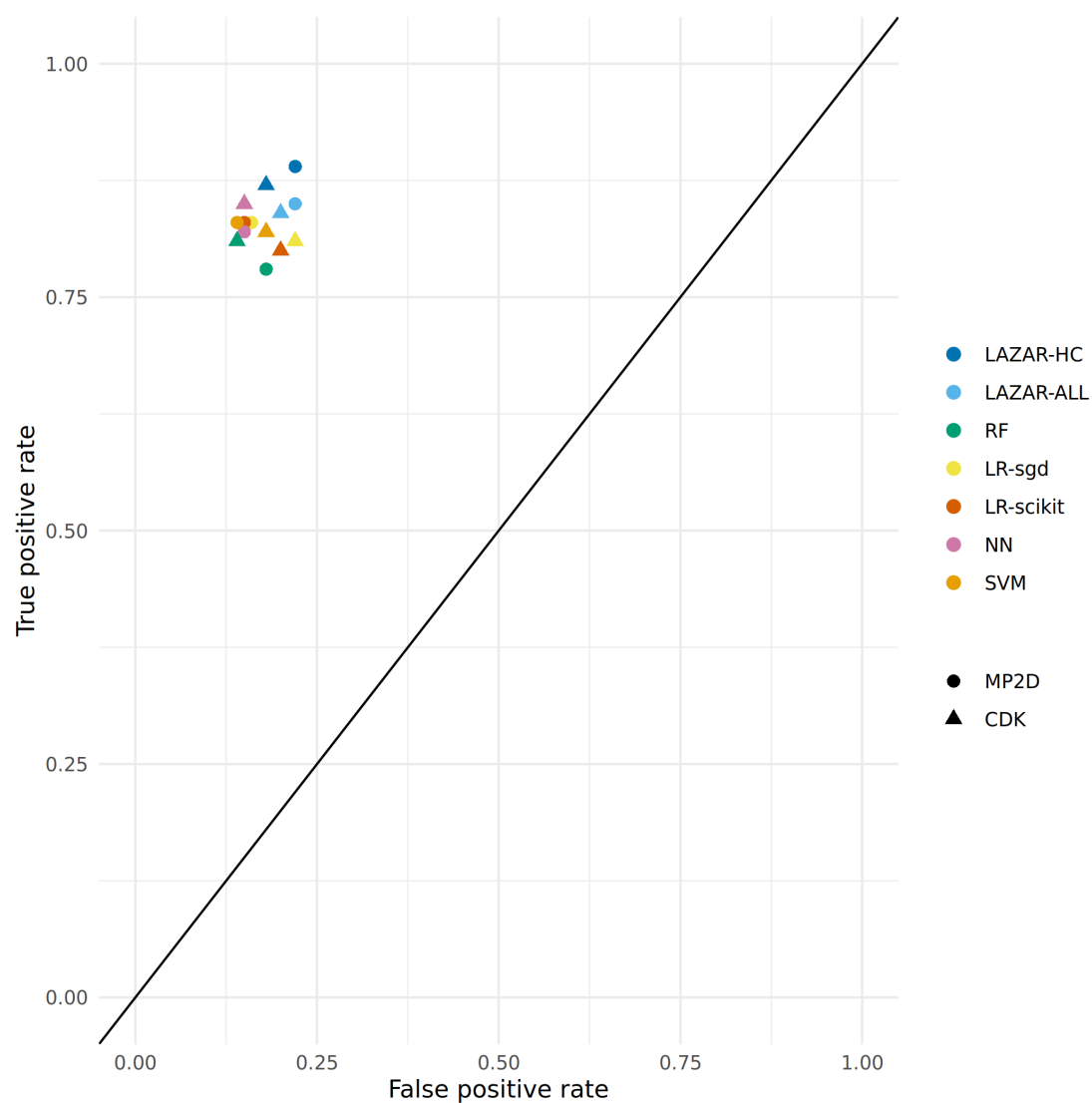


Figure 1: ROC plot of crossvalidation results (lazar-HC: lazar with high confidence, lazar-all: all lazar predictions, RF: random forests, LR-sgd: logistic regression (stochastic gradient descent), LR-scikit: logistic regression (scikit), NN: neural networks, SVM: support vector machines).

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

For the visualisation of the position of pyrrolizidine alkaloids in respect to the training data set we have applied t-distributed stochastic neighbor embedding (t-SNE, Maaten and Hinton (2008)) for MolPrint2D and CDK descriptors. t-SNE maps each high-dimensional object (chemical) to a two-dimensional point, maintaining the high-dimensional distances of the objects. Similar objects are represented by nearby points and dissimilar objects are represented by distant points. t-SNE coordinates were calculated with the R `Rtsne` package using the default settings (perplexity = 30, theta = 0.5, max_iter = 1000).

Figure 2 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training data in MP2D space (Tanimoto/Jaccard similarity), which resembles basically the structural diversity of the investigated compounds.

Figure 3 shows the t-SNE of pyrrolizidine alkaloids (PA) and the mutagenicity training data in CDK space (Euclidean similarity), which resembles basically the physical-chemical properties of the investigated compounds.

Figure 4 and Figure 5 depict two example pyrrolizidine alkaloid mutagenicity predictions in the context of training data. t-SNE visualisations of all investigated models can be downloaded from <https://git.in-silico.ch/mutagenicity-paper/figures>.

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

Table 3: Summary of pyrrolizidine alkaloid predictions

| Model | MP2D Mutagenic | Nr. predictions | CDK Mutagenic | Nr. predictions |
|-----------|----------------|-----------------|---------------|-----------------|
| lazar-all | 20% (111) | 93% (560) | 39% (193) | 83% (500) |
| lazar-HC | 25% (76) | 50% (301) | 45% (111) | 41% (246) |

| Model | MP2D Mutagenic | Nr. predictions | CDK Mutagenic | Nr. predictions |
|-----------|----------------|-----------------|---------------|-----------------|
| RF | 5% (28) | 100% (602) | 2% (10) | 100% (602) |
| LR-sgd | 21% (127) | 100% (602) | 16% (97) | 100% (602) |
| LR-scikit | 20% (118) | 100% (602) | 15% (88) | 100% (602) |
| NN | 21% (124) | 100% (602) | 25% (150) | 100% (602) |
| SVM | 14% (82) | 100% (602) | 3% (19) | 100% (602) |

Figure 6 displays the proportion of positive mutagenicity predictions from all models for the different pyrrolizidine alkaloid groups. Tensorflow models predicted all 602 pyrrolizidine alkaloids, **lazar** MP2D models predicted 560 compounds (301 with high confidence) and **lazar** CDK models 500 compounds (246 with high confidence).

For the **lazar**-HC model, only 50/41% of the PA dataset were within the stricter similarity thresholds of 0.5/0.9 (MP2D/CDK). Reduction of the similarity threshold to 0.2/0.5 in the **lazar**-all model increased the amount of predictable PAs to 93/83%. As the other ML models do not consider applicability domains, all PAs were predicted.

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, Figure 6).

Over all models, the mean value of mutagenic predicted PAs was highest for otonecines (65%, 407/623), followed by macrocyclic diesters (31%, 1042/3356), dehydropyrrolizidines (27%, 74/268), tertiary PAs (19%, 1201/6307) and retronecines (15%, 762/5054).

When excluding the aforementioned three deviating models, the rank order stays the same, but the percentage of mutagenic PAs is higher.

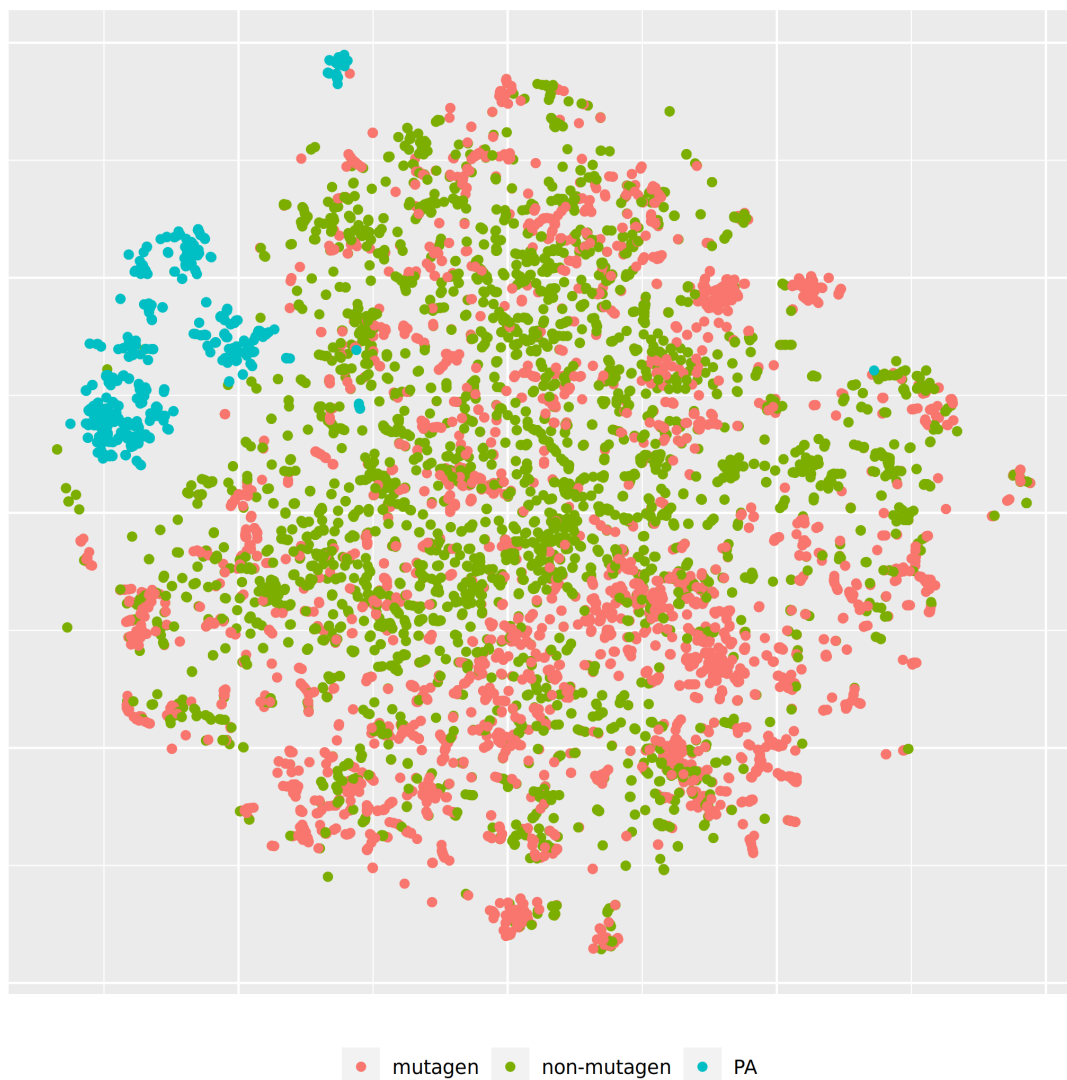


Figure 2: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA) in MP2D space

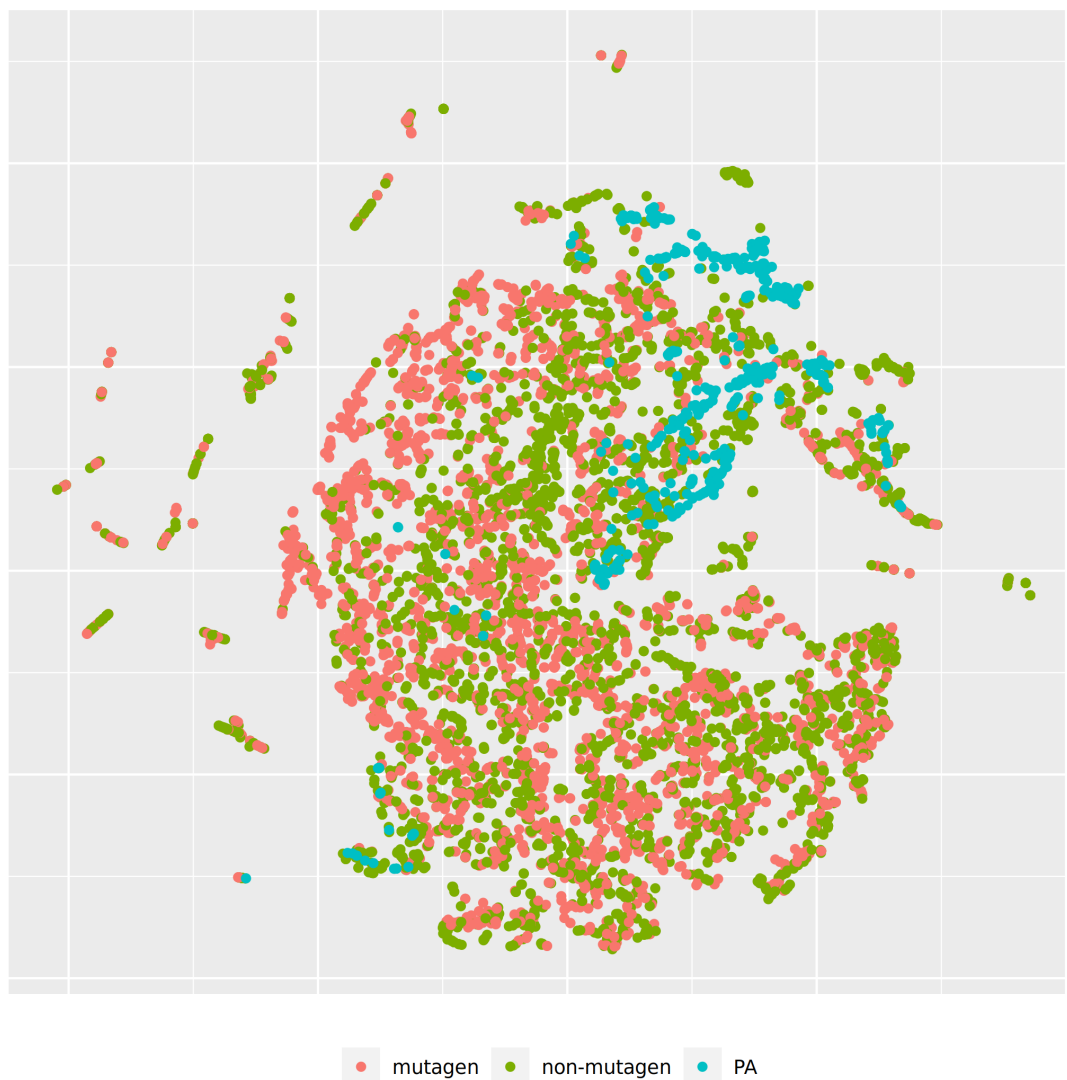


Figure 3: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA) in CDK space

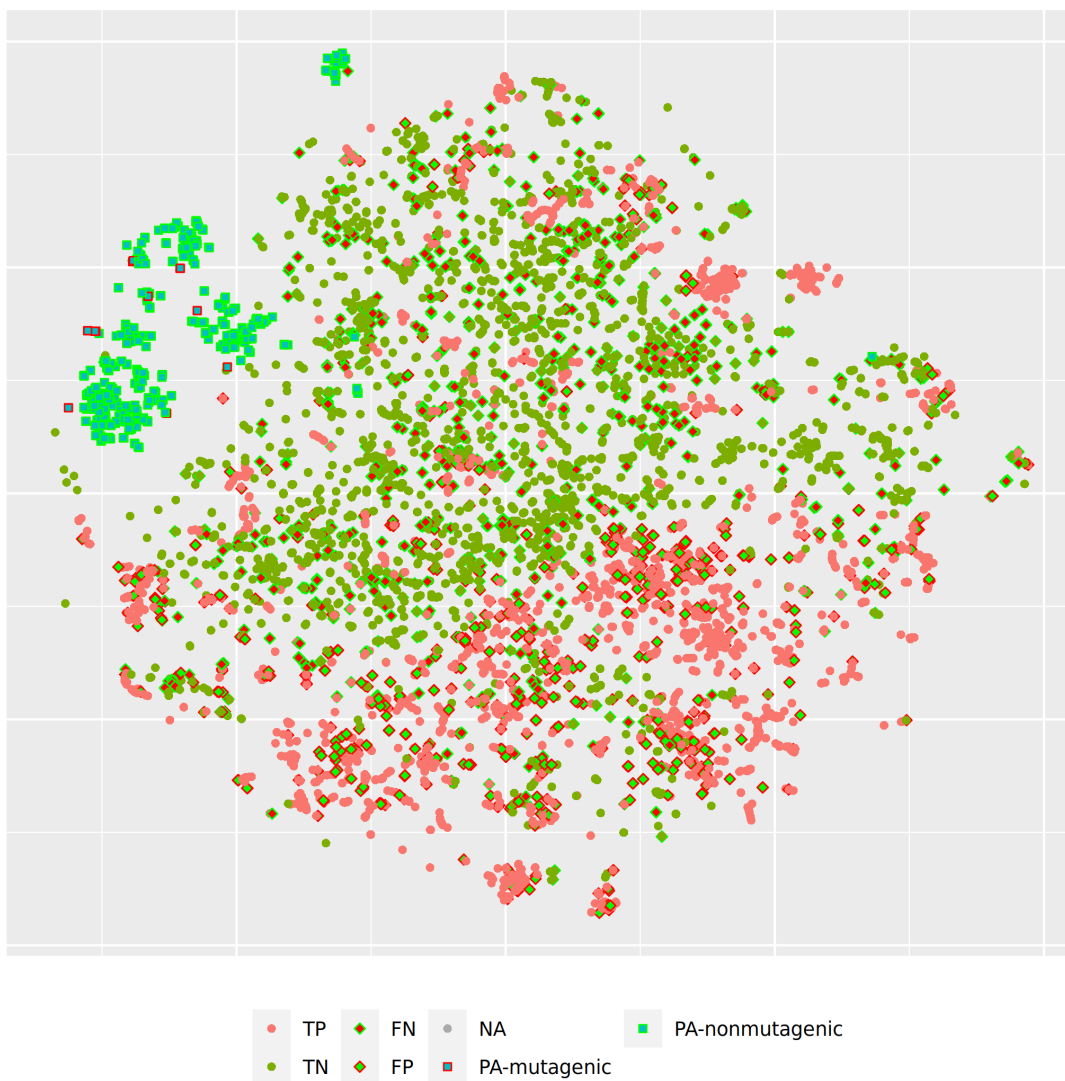


Figure 4: t-SNE visualisation of MP2D random forest predictions

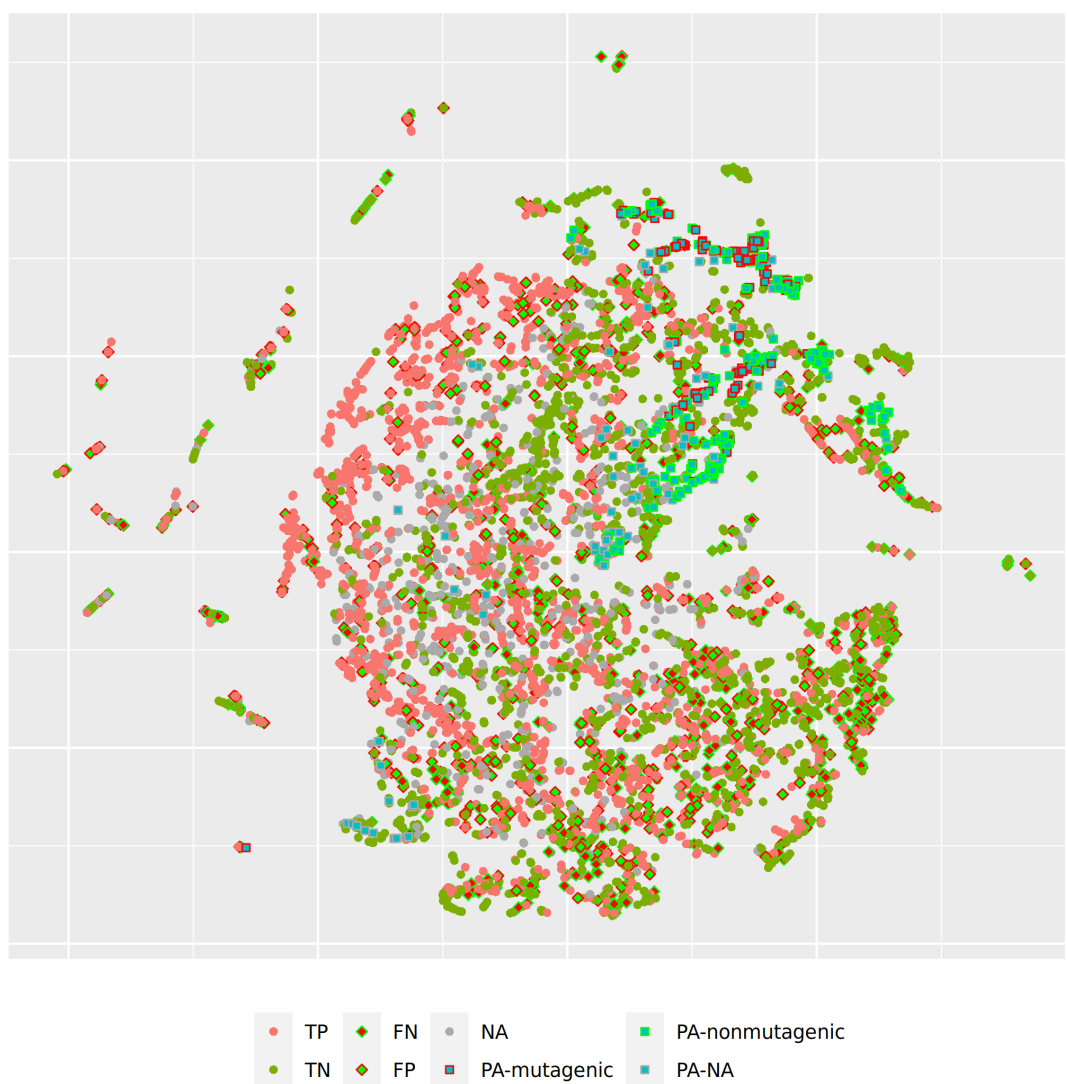


Figure 5: t-SNE visualisation of all CDK lazar predictions

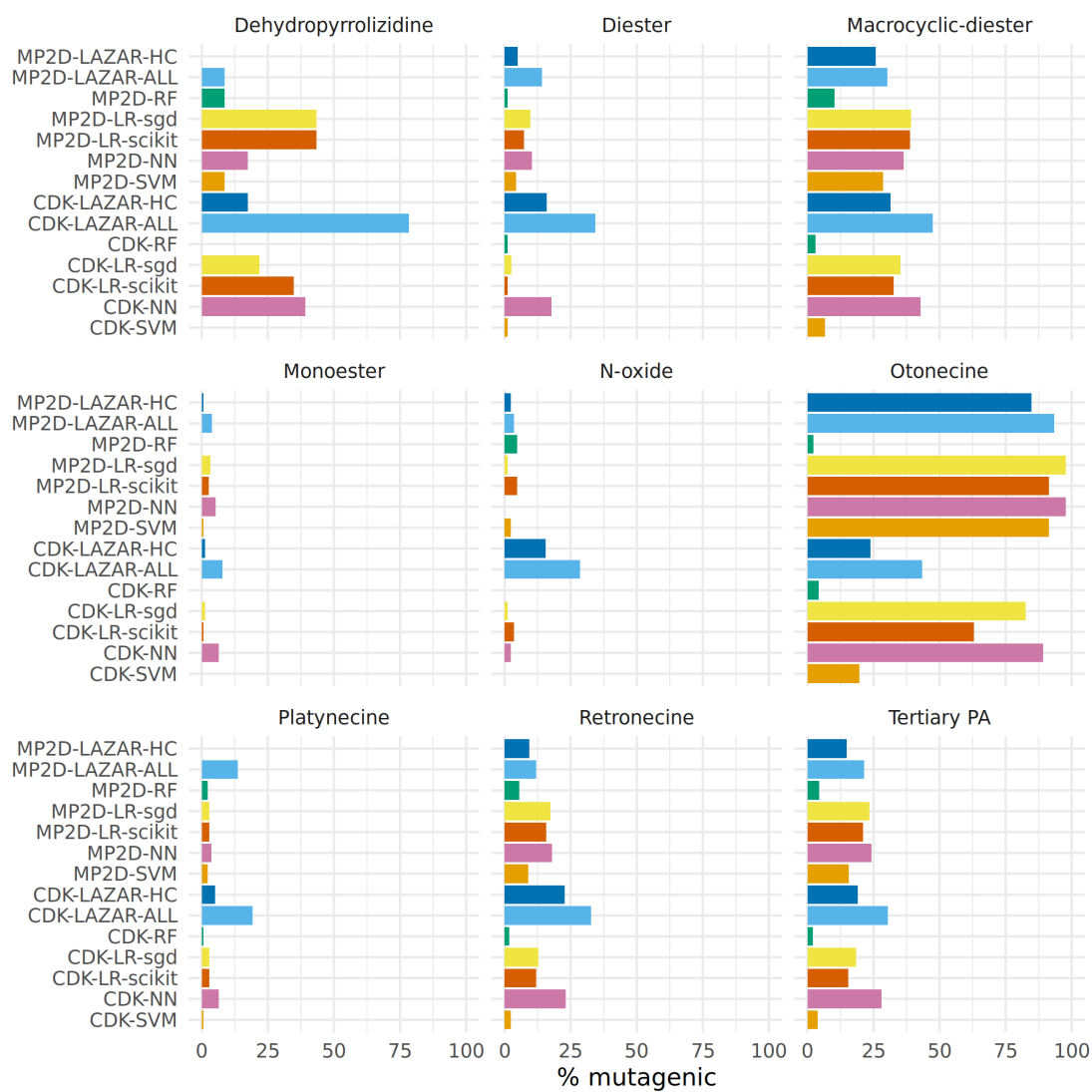


Figure 6: Summary of pyrrolizidine alkaloid predictions

328 The following rank order for mutagenic probability can be deduced from the results of
329 all models taken together:

330 Necine base: Platynecine < Retronecine « Otonecine

331 Necic acid: Monoester < Diester « Macrocyclic diester

332 Modification of necine base: N-oxide < Tertiary PA < Dehydropyrrolizidine

333 Discussion

334 Data

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

341 Algorithms

342 **lazar** is formally a *k-nearest-neighbor* algorithm that searches for similar structures
343 for a given compound and calculates the prediction based on the experimental data for
344 these structures. The QSAR literature calls such models frequently *local models*, because
345 models are generated specifically for each query compound. The investigated tensorflow
346 models are in contrast *global models*, i.e. a single model is used to make predictions for
347 all compounds. It has been postulated in the past, that local models are more accurate,
348 because they can account better for mechanisms that affect only a subset of the training
349 data.

350 Table 1, Table 2 and Figure 1 show that the crossvalidation accuracies of all models are
351 comparable to the experimental variability of the *Salmonella typhimurium* mutagenicity
352 bioassay (80-85% according to Piegorsch and Zeiger (1991)). All of these models have
353 balanced sensitivity (true positive rate) and specificity (true negative rate) and provide
354 highly significant concordance with experimental data (as determined by McNemar’s
355 Test). This is a clear indication that *in silico* predictions can be as reliable as the
356 bioassays. Given that the variability of experimental data is similar to model variability
357 it is impossible to decide which model gives the most accurate predictions, as models
358 with higher accuracies might just approximate experimental errors better than more
359 robust models.

360 Our results do not support the assumption that local models are superior to global
361 models for classification purposes. For regression models (lowest observed effect level)
362 we have found however that local models may outperform global models (Helma et al.
363 (2018)) with accuracies similar to experimental variability.

364 As all investigated algorithms give similar accuracies the selection will depend more on
365 practical considerations than on intrinsic properties. Nearest neighbor algorithms like
366 **lazar** have the practical advantage that the rationales for individual predictions can be
367 presented in a straightforward manner that is understandable without a background in
368 statistics or machine learning (Figure 7). This allows a critical examination of individual
369 predictions and prevents blind trust in models that are intransparent to users with a
370 toxicological background.

371 **Descriptors**

372 This study uses two types of descriptors for the characterisation of chemical structures:
373 *MolPrint2D* fingerprints (MP2D, Bender et al. (2004)) use atom environments (i.e.
374 connected atom types for all atoms in a molecule) as molecular representation, which

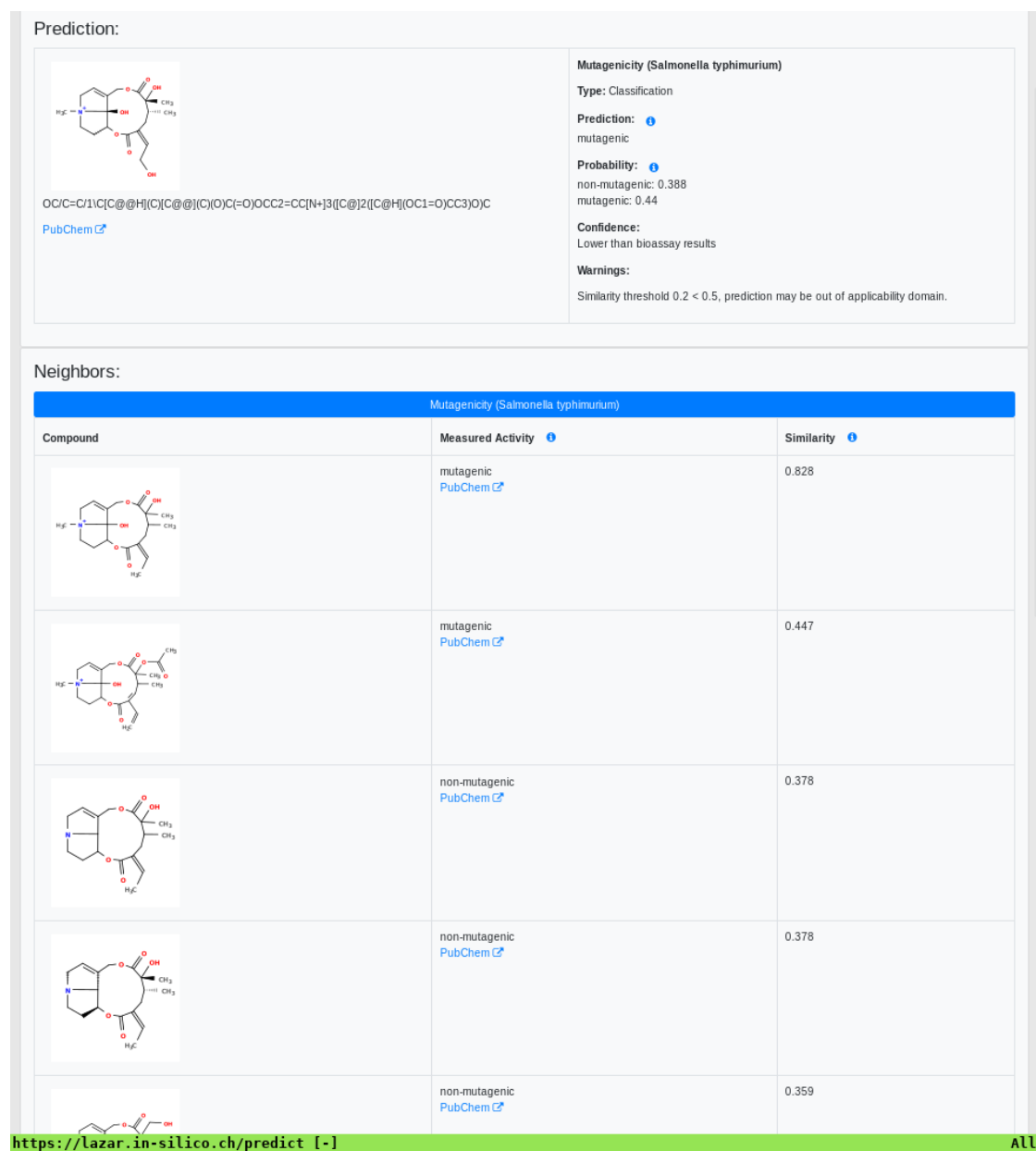


Figure 7: *lazar* screenshot of 12,21-Dihydroxy-4-methyl-4,8-secosenecinonan-8,11,16-trione mutagenicity prediction

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, Mayfield, and Alvarsson (2017)) descriptors were calculated with the PaDEL graphical interface (Yap (2011)). They include 1D and 2D topological descriptors as well as physical-chemical properties.

All investigated algorithms obtained models within the experimental variability for both types of descriptors (Table 1, Table 2, Figure 1).

Given that similar predictive accuracies are obtainable from both types of descriptors the choice depends once 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 which can cause substantial computational overhead.

400 **Pyrrolizidine alkaloid mutagenicity predictions**

401 **Algorithms and descriptors**

402 Figure 6 shows a clear differentiation between the different pyrrolizidine alkaloid groups.
403 Nevertheless differences between predictions from different algorithms and descriptors
404 (Table 3) were not expected based on crossvalidation results.

405 In order to investigate, if any of the investigated models show systematic errors in the
406 vicinity of pyrrolizidine-alkaloids we have performed a detailed t-SNE analysis of all
407 models (see Figure 4 and Figure 5 for two examples, all visualisations can be found at
408 <https://git.in-silico.ch/mutagenicity-paper/figures>.

409 None of the models showed obvious deviations from their expected behaviour, so the
410 reason for the disagreement between some of the models remains unclear at the moment.
411 It is however possible that some systematic errors are covered up by converting high
412 dimensional spaces to two coordinates and are thus invisible in t-SNE visualisations.

413 **Necic acid**

414 The rank order of the necic acid is comparable in all models. PAs from the monoester
415 type had the lowest genotoxic potential, followed by PAs from the open-ring diester
416 type. PAs with macrocyclic diesters had the highest genotoxic potential. The result fits
417 well with current state of knowledge: in general, PAs, which have a macrocyclic diesters
418 as necic acid, are considered to be more toxic than those with an open-ring diester or
419 monoester (EFSA (2011), Fu et al. (2004), Ruan2014b). This was also confirmed by
420 more recent studies, confirming that macrocyclic- and open-diester are more genotoxic
421 *in vitro* than monoesters (Hadi et al. (2021); Allemang et al. (2018), Louisse et al.
422 (2019)).

423 Necine base

424 In the rank order of necine base PAs, platynecine is the least mutagenic, followed by
425 retronecine, and otonecine. Saturated PAs of the platynecine-type are generally accepted
426 to be less or non-toxic and have been shown in *in vitro* experiments to form no DNA-
427 adducts (Xia et al. (2013)). In literature, otonecine-type PAs were shown to be more
428 toxic than those of the retronecine-type (Li et al. (2013)).

429 Modifications of necine base

430 The group-specific results reflect the expected relationship between the groups: the low
431 mutagenic potential of *N*-oxides and the high potential of dehydropyrrolizidines (DHP)
432 (Chen, Mei, and Fu (2010)). However, *N*-oxides may be *in vivo* converted back to their
433 parent toxic/tumorigenic parent PA (Yan et al. (2008)), on the other hand they are
434 highly water soluble and generally considered as detoxification products, which are *in*
435 *vivo* quickly renally eliminated (Chen, Mei, and Fu (2010)).

436 DHP are regarded as the toxic principle in the metabolism of PAs, and are known to
437 produce protein- and DNA-adducts (Chen, Mei, and Fu (2010)). None of our investigated
438 models did meet this expectation and all of them predicted the majority of DHP as non-
439 mutagenic. However, the following issues need to be considered. On the one hand, all
440 DHP were outside of the stricter applicability domain of MP2D **lazar**. This indicates
441 that they are structurally very different than the training data and might be out of the
442 applicability domain of all models based on this training set. In addition, DHP has two
443 unsaturated double bonds in its necine base, making it highly reactive. DHP and other
444 comparable molecules have a very short lifespan *in vivo*, and usually cannot be used in
445 *in vitro* experiments.

446 Overall the low number of positive mutagenicity predictions was unexpected. PAs are
447 generally considered to be genotoxic, and the mode of action is also known. Therefore,

the fact that some models predict the majority of PAs as not mutagenic seems contradictory. To understand this result, the experimental basis of the training dataset has to be considered. The training dataset is based on the *Salmonella typhimurium* mutagenicity bioassay (Ames test). There are some studies, which show mutagenicity of PAs in the Ames test (Chen, Mei, and Fu (2010)). Also, Rubiolo et al. (1992) examined several different PAs and several different extracts of PA-containing plants in the AMES test. They found that the Ames test was indeed able to detect mutagenicity of PAs, but in general, appeared to have a low sensitivity. The pre-incubation phase for metabolic activation of PAs by microsomal enzymes was the sensitivity-limiting step. This could very well mean that the low sensitivity of the Ames test for PAs is also reflected in the investigated models.

A *in vitro* screen of cellular PA effects (metabolic activation and mutagenic effects) in human and rodent hepatocytes (HepG2 and H-4-II-E) showed that results may also critically depend on the cellular model and cell culture conditions and may underestimate the effects of PAs (Forsch et al. (2018)).

In summary, we found marked differences in the predicted genotoxic potential between the PA groups: most toxic appeared the otonecines and macrocyclic diesters, least toxic the platynecines and the mono- and diesters. These results are comparable with *in vitro* measurements in hepatic HepaRG cells (Louisse et al. (2019)), where relative potencies (RP) were determined: for otonecines and cyclic diesters $RP = 1$, for open diesters $RP = 0.1$ and for monoesters $RP = 0.01$.

Due to a lack of differential data, European authorities based their risk assessment in a worst-case approach on lasiocarpine, for which sufficient data on genotoxicity and carcinogenicity were available (HMPC (2014), EMA (2020)). Our data further support a tiered risk assessment based on *in silico* and experimental data on the relative potency of individual PAs as already suggested by other authors (Merz and Schrenk (2016), Rutz

et al. (2020), Louisse et al. (2019)).

Conclusions

A new public *Salmonella* mutagenicity training dataset with 8309 experimental results was created and used to train **lazar** and Tensorflow models with MolPrint2D and CDK descriptors. All investigated algorithm and descriptor combinations showed accuracies comparable to the interlaboratory variability of the Ames test.

Pyrrolizidine alkaloid predictions showed a clear separation between different classes of PAs which were generally in accordance with the current toxicological knowledge about these compounds. Some of the models showed however a substantially lower number of mutagenicity predictions, despite similar crossvalidation results and we were unable to identify the reasons for this discrepancy within this investigation.

Thus the practical question how to choose model predictions in the absence of experimental data remains open. Tensorflow predictions do not include applicability domain estimations and the rationales for predictions cannot be traced by toxicologists. Transparent models like **lazar** may have an advantage in this context, because they present rationales for predictions (similar compounds with experimental data) which can be accepted or rejected by toxicologists and provide validated applicability domain estimations.

Our data show that large difference exist with regard to genotoxic potential between different pyrrolizidine subgroups. These results may allow to adjust risk assessment of pyrrolizidine contamination.

References

- Allemand, Ashley, Catherine Mahony, Cathy Lester, and Stefan Pfuhler. 2018. "Relative Potency of Fifteen Pyrrolizidine Alkaloids to Induce DNA Damage as Measured by Micronucleus Induction in HepaRG Human Liver Cells." *Food and Chemical Toxicology* 121: 72–81. <https://doi.org/https://doi.org/10.1016/j.fct.2018.08.003>.
- Bender, Andreas, Hamse Y. Mussa, Robert C. Glen, and Stephan Reiling. 2004. "Molecular Similarity Searching Using Atom Environments, Information-Based Feature Selection, and a Naïve Bayesian Classifier." *Journal of Chemical Information and Computer Sciences* 44 (1): 170–78. <https://doi.org/10.1021/ci034207y>.
- Chen, T., N. Mei, and P. P. Fu. 2010. "Genotoxicity of Pyrrolizidine Alkaloids." *J. Appl. Toxicol.*, 183–96. <https://doi.org/https://doi.org/10.1002/jat.1504>.
- (ECHA), European Chemicals Agency. 2017. "Guidance on Information Requirements and Chemical Safety Assessment, Chapter R.7a: Endpoint Specific Guidance." <https://doi.org/10.2823/337352>.
- EFSA. 2011. "Scientific Opinion on Pyrrolizidine Alkaloids in Food and Feed." *EFSA Journal*, no. 9: 1–134. <https://doi.org/https://doi.org/10.2903/j.efsa.2011.2406>.
- EFSA. 2016. "Guidance on the Establishment of the Residue Definition for Dietary Assessment: EFSA Panel on Plant Protect Products and Their Residues (PPR)." *EFSA Journal*, no. 14: 1–12. <https://doi.org/https://doi.org/10.2903/j.efsa.2016.4549>.
- EMA. 2020. "Public Statement on the Use of Herbal Medicinal Products Containing Toxic, Unsaturated Pyrrolizidine Alkaloids (Pas) Including Recommendations Regarding Contamination of Herbal Medicinal Products with Pyrrolizidine Alkaloids. European Medicines Agency, Committee on Herbal Medicinal Products (Hmpc), Ema/Hmpc/893108/2011 Rev.1."

519 Forsch, K., V. Schöning, L. Disch, B. Siewert, M. Unger, and J. Drewe. 2018. “Devel-
520 opment of an in Vitro Screening Method of Acute Cytotoxicity of the Pyrrolizidine
521 Alkaloid Lasiocarpine in Human and Rodent Hepatic Cell Lines by Increasing Sus-
522 ceptibility.” *Journal of Ethnopharmacology*, no. 217: 134–39. <https://doi.org/https://doi.org/10.1016/j.jep.2018.02.018>.
523

524 Fu, P. P., Q. Xia, G. Lin, and M. W. Chou. 2004. “Pyrrolizidine Alkaloids–Genotoxicity,
525 Metabolism Enzymes, Metabolic Activation, and Mechanisms.” *Drug Metab. Rev.*, no.
526 36: 1–55. <https://doi.org/https://doi.org/https://doi.org/10.1081/dmr-120028426>.

527 Hadi, Naji Said Aboud, Ezgi Eyluel Bankoglu, Lea Schott, Eva Leopoldsberger, Vanessa
528 Ramge, Olaf Kelber, Hartwig Sievers, and Helga Stopper. 2021. “Genotoxicity of Se-
529 lected Pyrrolizidine Alkaloids in Human Hepatoma Cell Lines HepG2 and Huh6.” *Mu-
530 tation Research/Genetic Toxicology and Environmental Mutagenesis* 861-862: 503305.
531 <https://doi.org/https://doi.org/10.1016/j.mrgentox.2020.503305>.

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

536 Hartmann, T., and L. Witte. 1995. “Chemistry, Biology and Chemoecology of the
537 Pyrrolizidine Alkaloids.” In *Alkaloids: Chemical and Biological Perspectives*, edited by
538 S. W. Pelletier, 155–233. London, New York: Pergamon.

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

543 HMPC. 2014. “Public Statement on the Use of Herbal Medicinal Products 5 Containing

544 Toxic, Unsaturated Pyrrolizidine Alkaloids (Pas), European Medicines Agency, Commit-
545 tee on Herbal Medicinal Products (Hmpc) Ema/Hmpc/8931082011.”

546 International Council for Harmonisation of Technical Requirements for Pharmaceuticals
547 for Human Use (ICH). 2017. “Assessment and Control of DNA Reactive (Mutagenic)
548 Impurities in Pharmaceuticals to Limit Potential Carcinogenic Risk M7(R1).”

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

551 Langel, D., D. Ober, and P. B. Pelsner. 2011. “The Evolution of Pyrrolizidine Alkaloid
552 Biosynthesis and Diversity in the Senecioneae.” *Phytochemistry Reviews*, no. 10: 3–74.

553 Li, Yan Hong, Winnie Lai Ting Kan, Na Li, and Ge Lin. 2013. “Assessment of
554 Pyrrolizidine Alkaloid-Induced Toxicity in an in Vitro Screening Model.” *Journal of*
555 *Ethnopharmacology* 150 (2): 560–67. [https://doi.org/https://doi.org/10.1016/j.jep.2013.](https://doi.org/https://doi.org/10.1016/j.jep.2013.09.010)
556 09.010.

557 Louisse, Jochem, Deborah Rijkers, Geert Stoopen, Wendy Jansen Holleboom, Mona
558 Delagrangé, Elise Molthof, Patrick P. J. Mulder, Ron L. A. P. Hoogenboom, Marc Au-
559 debert, and Ad A. C. M. Peijnenburg. 2019. “Determination of Genotoxic Potencies of
560 Pyrrolizidine Alkaloids in HepaRG Cells Using the H2AX Assay.” *Food and Chemical*
561 *Toxicology* 131: 110532. [https://doi.org/https://doi.org/10.1016/j.fct.2019.05.040.](https://doi.org/https://doi.org/10.1016/j.fct.2019.05.040)

562 Maaten, L. J. P. van der, and G. E. Hinton. 2008. “Visualizing Data Using t-SNE.”
563 *Journal of Machine Learning Research*, no. 9: 2579–2605.

564 Mattocks, AR. 1986. *Chemistry and Toxicology of Pyrrolizidine Alkaloids*. Academic
565 Press.

566 Merz, Karl-Heinz, and Dieter Schrenk. 2016. “Interim Relative Potency Factors for the
567 Toxicological Risk Assessment of Pyrrolizidine Alkaloids in Food and Herbal Medicines.”
568 *Toxicology Letters* 263: 44–57. <https://doi.org/https://doi.org/10.1016/j.toxlet.2016.05.>

569 002.

570 O’Boyle, Noel, Michael Banck, Craig James, Chris Morley, Tim Vandermeersch, and
571 Geoffrey Hutchison. 2011. “Open Babel: An open chemical toolbox.” *J. Cheminf.* 3 (1):
572 33. <https://doi.org/doi:10.1186/1758-2946-3-33>.

573 Piegorsch, W. W., and E. Zeiger. 1991. “Measuring Intra-Assay Agreement for the
574 Ames Salmonella Assay.” In *Statistical Methods in Toxicology, Lecture Notes in Medical*
575 *Informatics*, edited by L. Hotorn, 35–41. Springer-Verlag.

576 Rubiolo, P., L. Pieters, M. Calomme, C. Bicchi, A. Vlietinck, and D. Vanden
577 Berghe. 1992. “Mutagenicity of Pyrrolizidine Alkaloids in the Salmonella Ty-
578 phimurium/Mammalian Microsome System.” *Mutation Research*, no. 281: 143–47.
579 [https://doi.org/https://doi.org/https://doi.org/10.1016/0165-7992\(92\)90050-r](https://doi.org/https://doi.org/https://doi.org/10.1016/0165-7992(92)90050-r).

580 Rutz, L., L. Gao, J. H. Küpper, and others. 2020. “Structure-Dependent
581 Genotoxic Potencies of Selected Pyrrolizidine Alkaloids in Metabolically Com-
582 petent Hepg2 Cells.” *Arch. Toxicol.*, no. 94: 4159–72. <https://doi.org/https://doi.org/10.1007/s00204-020-02895-z>.

584 Schöning, Verena, Felix Hamann, Mark Peinl, and Jürgen Drewe. 2017. “Editor’s
585 Highlight: Identification of Any Structure-Specific Hepatotoxic Potential of Different
586 Pyrrolizidine Alkaloids Using Random Forests and Artificial Neural Networks.” *Toxicol.*
587 *Sci.*, no. 160: 361–70. [https://doi.org/https://doi.org/https://doi.org/10.1093/toxsci/](https://doi.org/https://doi.org/https://doi.org/10.1093/toxsci/kfx187)
588 [kfx187](https://doi.org/https://doi.org/10.1093/toxsci/kfx187).

589 Weininger, David, Arthur Weininger, and Joseph L. Weininger. 1989. “SMILES. 2.
590 Algorithm for Generation of Unique Smiles Notation.” *J. Chem. Inf. Comput. Sci.*, no.
591 29: 97–101. <https://doi.org/https://doi.org/10.1021/ci00062a008>.

592 Willighagen, E. L., J. W. Mayfield, and J. et al. Alvarsson. 2017. “The Chemistry
593 Development Kit (Cdk) V2.0: Atom Typing, Depiction, Molecular Formulas, and Sub-

594 structure Searching.” *J. Cheminform.*, no. 9(33). <https://doi.org/https://doi.org/10.1186/s13321-017-0220-4>.
595

596 Xia, Q., Y. Zhao, L. S. Von Tungeln, D. R. Doerge, G. Lin, G. Cai, and P. P. Fu. 2013.
597 “Pyrrolizidine Alkaloid-Derived DNA Adducts as a Common Biological Biomarker of
598 Pyrrolizidine Alkaloid-Induced Tumorigenicity.” *Chem Res. Toxicol.*, no. 26: 1384–96.
599 <https://doi.org/https://doi.org/https://doi.org/10.1021/tx400241c>.

600 Yan, J., Q. Xia, M. W. Chou, and P. P. Fu. 2008. “Metabolic Activation of
601 Retronecine and Retronecine N-oxide - Formation of DHP-Derived DNA Adducts.”
602 *Toxicol. Ind. Health*, no. 24(3): 181–8. <https://doi.org/https://doi.org/https://doi.org/10.1177/0748233708093727>.
603

604 Yap, C. W. 2011. “PaDEL-descriptor: An Open Source Software to Calculate Molecular
605 Descriptors and Fingerprints.” *Journal of Computational Chemistry*, no. 32: 1466–74.
606 <https://doi.org/https://doi.org/10.1002/jcc.21707>.