

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

22 Introduction

23 **TODO:** rationale for investigation

24 As case study we decided to apply these mutagenicity models to Pyrrolizidines alkaloids
25 (PAs) in order to highlight potentials and problems with the applicability of mutagenicity
26 models for compounds with very limited experimental data.

27 Pyrrolizidine alkaloids (PAs) are characteristic metabolites of some plant families,
28 mainly: *Asteraceae*, *Boraginaceae*, *Fabaceae* and *Orchidaceae* (Hartmann and Witte
29 (1995), Langel, Ober, and B. (2011)) and form a powerful defence mechanism against
30 herbivores. PAs are heterocyclic ester alkaloids composed of a necine base (two fused
31 five-membered rings joined by a single nitrogen atom) and a necic acid (one or two
32 carboxylic ester arms), occurring principally in two forms, tertiary base PAs and PA
33 N-oxides. Several *in vitro* studies have shown the mutagenic potential of PAs, which
34 seems highly dependent on structure of necine base and necic acid (Hadi et al. (2021);
35 Allemang et al. (2018), Louisse et al. (2019)). However, due to limited availability of
36 pure substances, only a limited number of PAs have been investigated with regards to
37 their structure-specific mutagenicity. To overcome this bottleneck, the prediction of
38 structure-specific mutagenic potential of PAs with different machine learning models
39 could provide further inside in the mechanisms.

40 Summing up the main objectives of this study were

- 41 • to generate a new mutagenicity training dataset, by combining the most compre-
42 hensive public datasets
- 43 • to compare the performance of MolPrint2D (*MP2D*) fingerprints with Chemistry
44 Development Kit (*CDK*) descriptors
- 45 • to compare the performance of global QSAR models (random forests (*RF*), support
46 vector machines (*SVM*), logistic regression (*LR*), neural nets (*NN*)) with local

- 47 models (`lazar`)
- 48 • to apply these models for the prediction of pyrrolizidine alkaloid mutagenicity

49 **Materials and Methods**

50 **Data**

51 **Mutagenicity training data**

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

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

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

64 Dataset merges were based on unique SMILES (*Simplified Molecular Input Line En-*
65 *try Specification*, Weininger, Weininger, and Weininger (1989)) strings of the compound
66 structures. Duplicated experimental data with the same outcome was merged into a
67 single value, because it is likely that it originated from the same experiment. Contradic-
68 tory results were kept as multiple measurements in the database. The combined training
69 dataset contains 8290 unique structures and 8309 individual measurements.

70 Source code for all data download, extraction and merge operations is pub-
71 licly available from the git repository <https://git.in-silico.ch/mutagenicity-paper>
72 under a GPL3 License. The new combined dataset can be found at <https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/mutagenicity.csv>.
73

74 **Pyrrolizidine alkaloid (PA) dataset**

75 The pyrrolizidine alkaloid dataset was created from five independent, necine base sub-
76 structure searches in PubChem (<https://pubchem.ncbi.nlm.nih.gov/>) and compared to
77 the PAs listed in the EFSA publication EFSA (2011) and the book by Mattocks Mattocks
78 (1986), to ensure, that all major PAs were included. PAs mentioned in these publica-
79 tions which were not found in the downloaded substances were searched individually
80 in PubChem and, if available, downloaded separately. Non-PA substances, duplicates,
81 and isomers were removed from the files, but artificial PAs, even if unlikely to occur in
82 nature, were kept. The resulting PA dataset comprised a total of 602 different PAs.

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

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

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

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

- 91 • N-oxide-type (84 compounds)
- 92 • Tertiary-type (PAs which were neither from the N-oxide- nor DHP-type, 495 com-
93 pounds)

- Dehydropyrrolizidine-type (pyrrolic ester, 23 compounds)

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

- 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/>)

117 mp2d/mutagenicity-fingerprints.csv.gz)

118 *Pyrrolizidine alkaloids:*

- 119 • 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)
120 mp2d/fingerprints.mp2d)
- 121 • 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)
122 mp2d/pa-fingerprints.csv.gz)

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

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

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

132 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)
133 [tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv](https://git.in-silico.ch/mutagenicity-paper/tree/mutagenicity/cdk/mutagenicity-mod-2.new.csv).

134 The same procedure was applied for the pyrrolizidine dataset yielding descriptors for
135 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)
136 [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).

137 Algorithms

138 **lazar**

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

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

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

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

151 Feature preprocessing

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

156 Neighbour identification

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

159 For MolPrint2D fingerprints chemical similarity between two compounds a and b is
160 expressed as the proportion between atom environments common in both structures
161 $A \cap B$ and the total number of atom environments $A \cup B$ (Jaccard/Tanimoto index).

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

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

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

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

- 168 • First a similarity threshold of 0.5 (MP2D/Tanimoto) or 0.9 (CDK/Cosine) is used
169 to collect neighbours, to create a local QSAR model and to make a prediction for
170 the query compound. This are predictions with *high confidence*.
- 171 • If any of these steps fails, the procedure is repeated with a similarity threshold of
172 0.2 (MP2D/Tanimoto) or 0.7 (CDK/Cosine) and the prediction is flagged with a
173 warning that it might be out of the applicability domain of the training data (*low*
174 *confidence*).
- 175 • These Similarity thresholds are the default values chosen by software developers
176 and remained unchanged during the course of these experiments.

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

179 Local QSAR models and predictions

180 Only similar compounds (neighbours) above the threshold are used for local QSAR
181 models. In this investigation, we are using a weighted majority vote from the neigh-
182 bour’s experimental data for mutagenicity classifications. Probabilities for both classes
183 (mutagenic/non-mutagenic) are calculated according to the following formula and the
184 class with the higher probability is used as prediction outcome.

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

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

186 $\sum \text{sim}_{n,c}$ Sum of similarities of neighbours with class c

187 $\sum \text{sim}_n$ Sum of all neighbours

188 Applicability domain

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

197 Validation

198 10-fold cross validation was performed for model evaluation.

199 **Pyrrolizidine alkaloid predictions**

200 For the prediction of pyrrolizidine alkaloids models were generated with the MP2D and
201 CDK training datasets. The complete feature set was used for MP2D predictions, for
202 CDK predictions the intersection between training and pyrrolizidine alkaloid features
203 was used.

204 **Availability**

- 205 • Source code for this manuscript (GPL3): [https://git.in-silico.ch/lazar/tree/?h=](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
206 [mutagenicity-paper](https://git.in-silico.ch/lazar/tree/?h=mutagenicity-paper)
- 207 • Crossvalidation experiments (GPL3): [https://git.in-silico.ch/lazar/tree/models/](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
208 [?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/models/?h=mutagenicity-paper)
- 209 • Pyrrolizidine alkaloid predictions (GPL3): [https://git.in-silico.ch/lazar/tree/](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
210 [predictions/?h=mutagenicity-paper](https://git.in-silico.ch/lazar/tree/predictions/?h=mutagenicity-paper)
- 211 • Public web interface: <https://lazar.in-silico.ch>

212 **Tensorflow models**

213 **Feature Preprocessing**

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

216 **Random forests (*RF*)**

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

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

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

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

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

225 **Neural Nets (*NN*)**

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

230 **Support vector machines (*SVM*)**

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

233 **Validation**

234 10-fold cross-validation was used for all Tensorflow models.

235 **Pyrrolizidine alkaloid predictions**

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

239 Availability

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

241 *Crossvalidation:*

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

244 • CDK descriptors: [https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/](https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/cdk/tensorflow)
245 [cdk/tensorflow](https://git.in-silico.ch/mutagenicity-paper/tree/crossvalidations/cdk/tensorflow)

246 *Pyrrolizidine alkaloids:*

247 • MolPrint2D fingerprints: [https://git.in-silico.ch/mutagenicity-paper/tree/](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/tensorflow)
248 [pyrrolizidine-alkaloids/mp2d/tensorflow](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/mp2d/tensorflow)

249 • CDK descriptors: [https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/tensorflow)
250 [cdk/tensorflow](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/cdk/tensorflow)

251 • CDK desc

252 Results

253 10-fold crossvalidations

254 Crossvalidation results are summarized in the following tables: Table 1 shows results
255 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

	lazar-HC	lazar-all	RF	LR-sgd	LR-scikit	NN	SVM
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* mu-

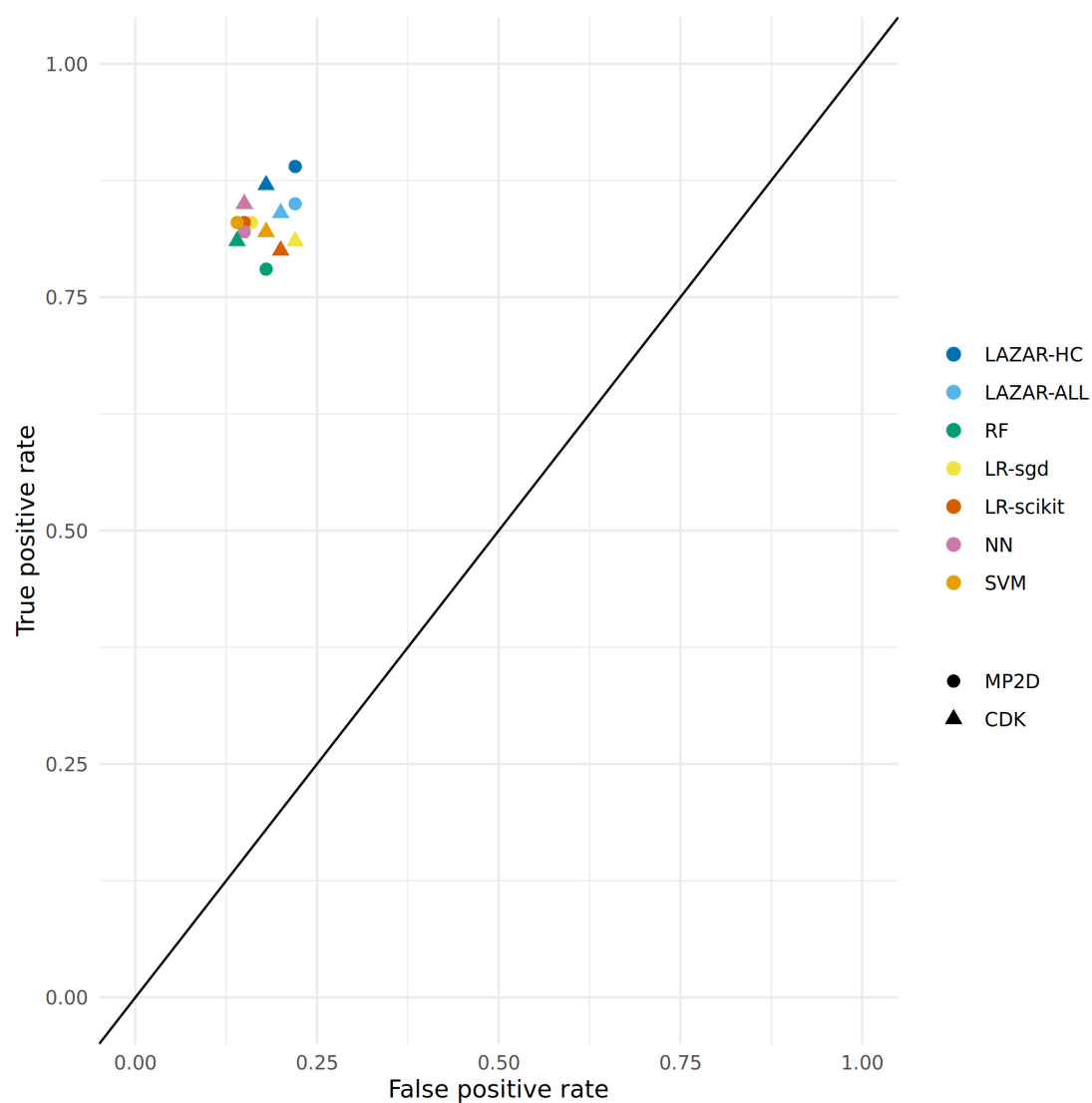


Figure 1: ROC plot of crossvalidation results (lazar-HC: lazar with high confidence, lazar-all: all lazar predictions, RF: random forests, LR-sgd: logistic regression (stochastic gradient descent), LR-scikit: logistic regression (scikit), NN: neural networks, SVM: support vector machines).

263 tagenicity bioassay (80-85%, Benigni and Giuliani (1988)). Sensitivities and specificities
264 are balanced in all of these models.

265 **Pyrrolizidine alkaloid mutagenicity predictions**

266 Mutagenicity predictions of 602 pyrrolizidine alkaloids (PAs) from all investigated
267 models can be downloaded from [https://git.in-silico.ch/mutagenicity-paper/tree/
268 pyrrolizidine-alkaloids/pa-predictions.csv](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/pa-predictions.csv). A visual representation of all PA predictions
269 can be found at [https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/
270 pa-predictions.pdf](https://git.in-silico.ch/mutagenicity-paper/tree/pyrrolizidine-alkaloids/pa-predictions.pdf).

271 For the visualisation of the position of pyrrolizidine alkaloids in respect to the train-
272 ing data set we have applied t-distributed stochastic neighbor embedding (t-SNE,
273 Maaten and Hinton (2008)) for MolPrint2D and CDK descriptors. t-SNE maps
274 each high-dimensional object (chemical) to a two-dimensional point, maintaining the
275 high-dimensional distances of the objects. Similar objects are represented by nearby
276 points and dissimilar objects are represented by distant points. t-SNE coordinates were
277 calculated with the R *Rtsne* package using the default settings (perplexity = 30, theta
278 = 0.5, max_iter = 1000).

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

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

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



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

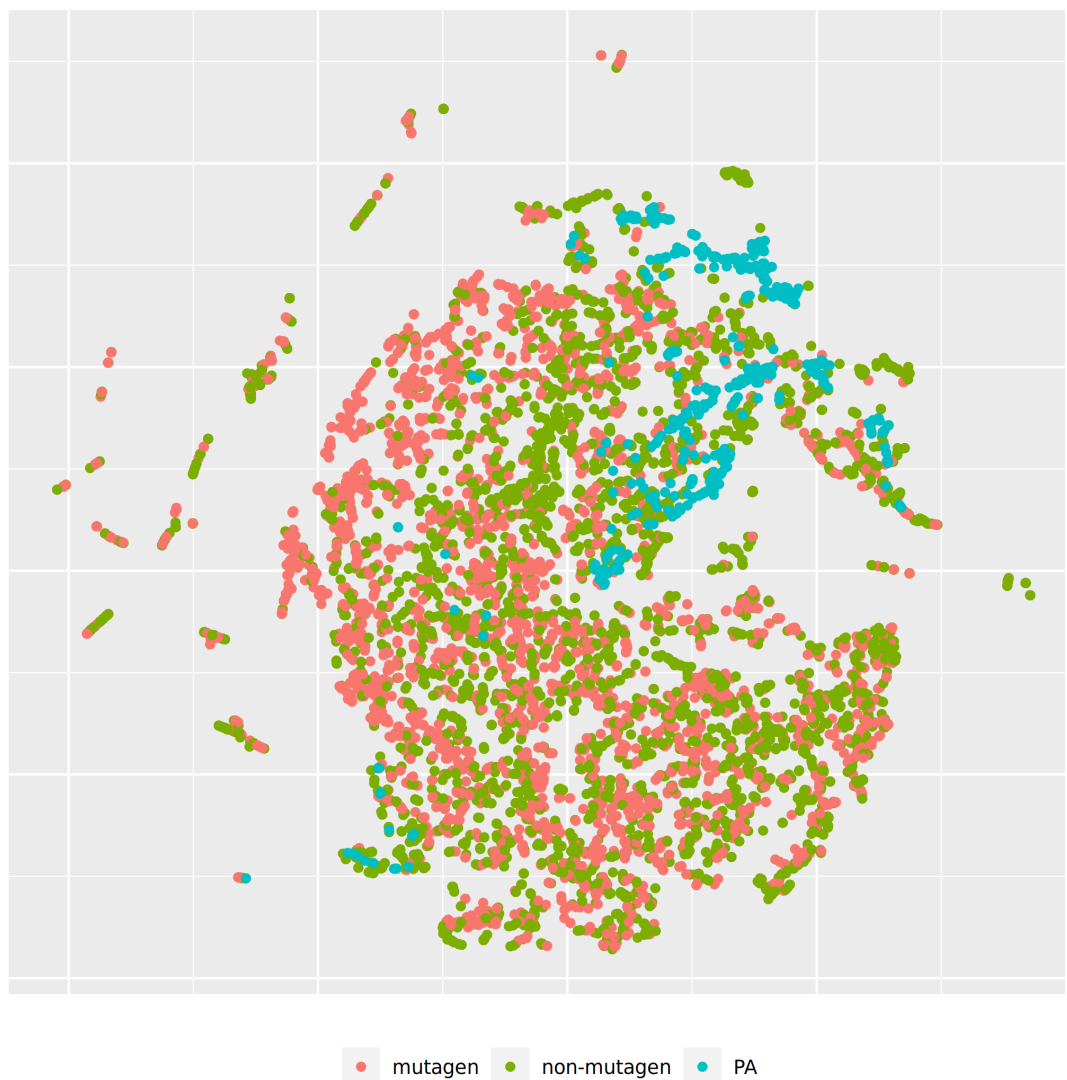


Figure 3: t-SNE visualisation of mutagenicity training data and pyrrolizidine alkaloids (PA) in CDK space

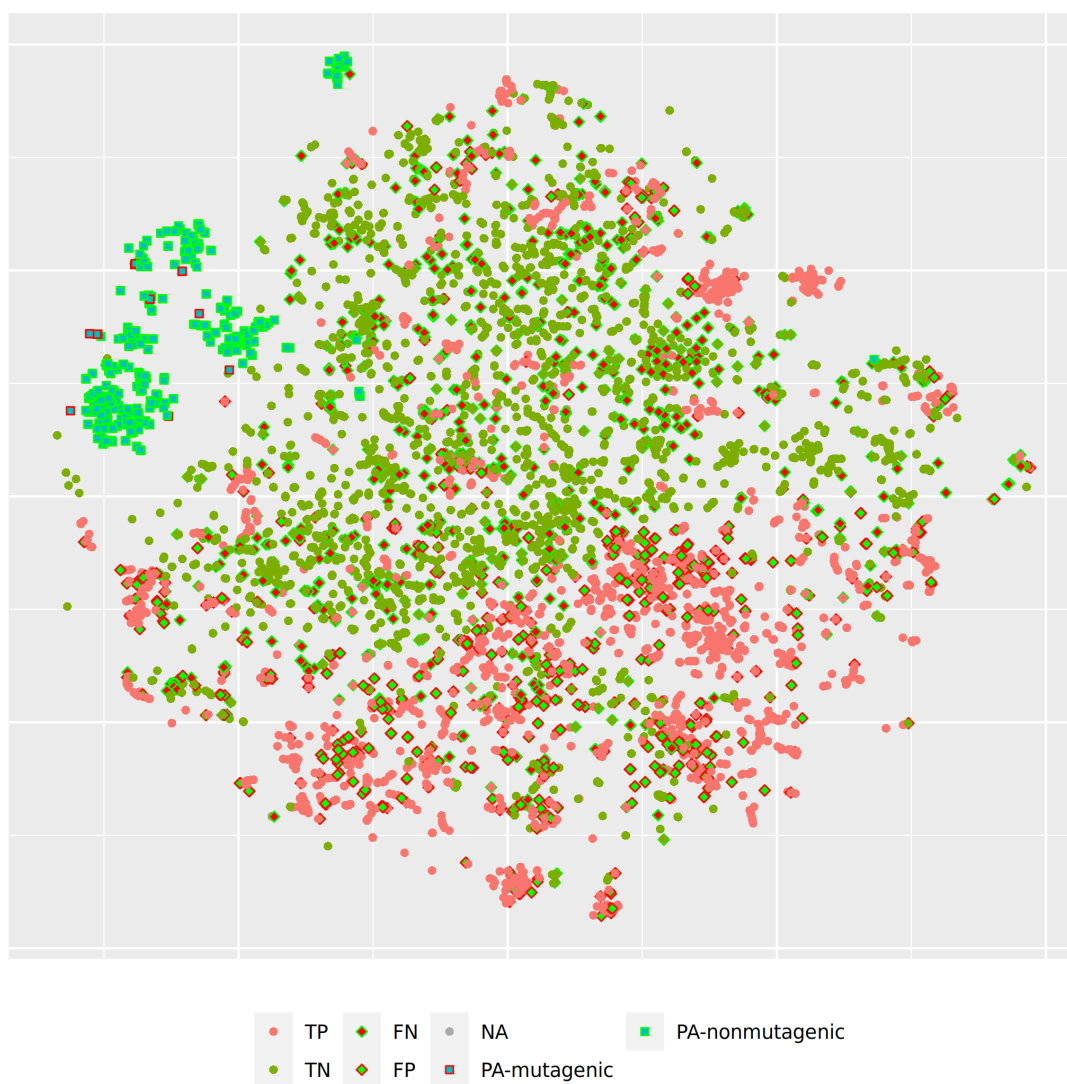


Figure 4: t-SNE visualisation of MP2D random forest predictions

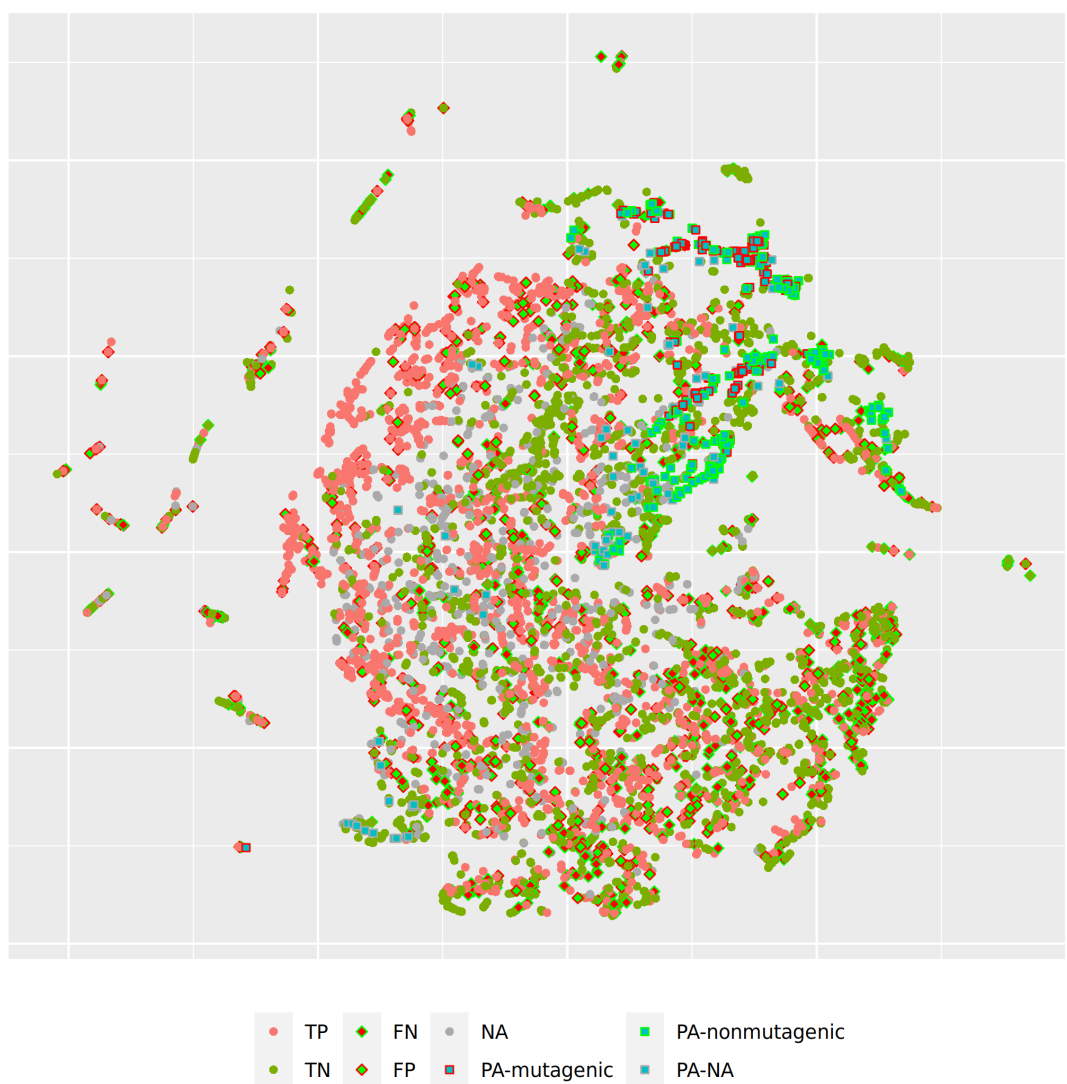


Figure 5: t-SNE visualisation of all CDK lazar predictions

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)
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

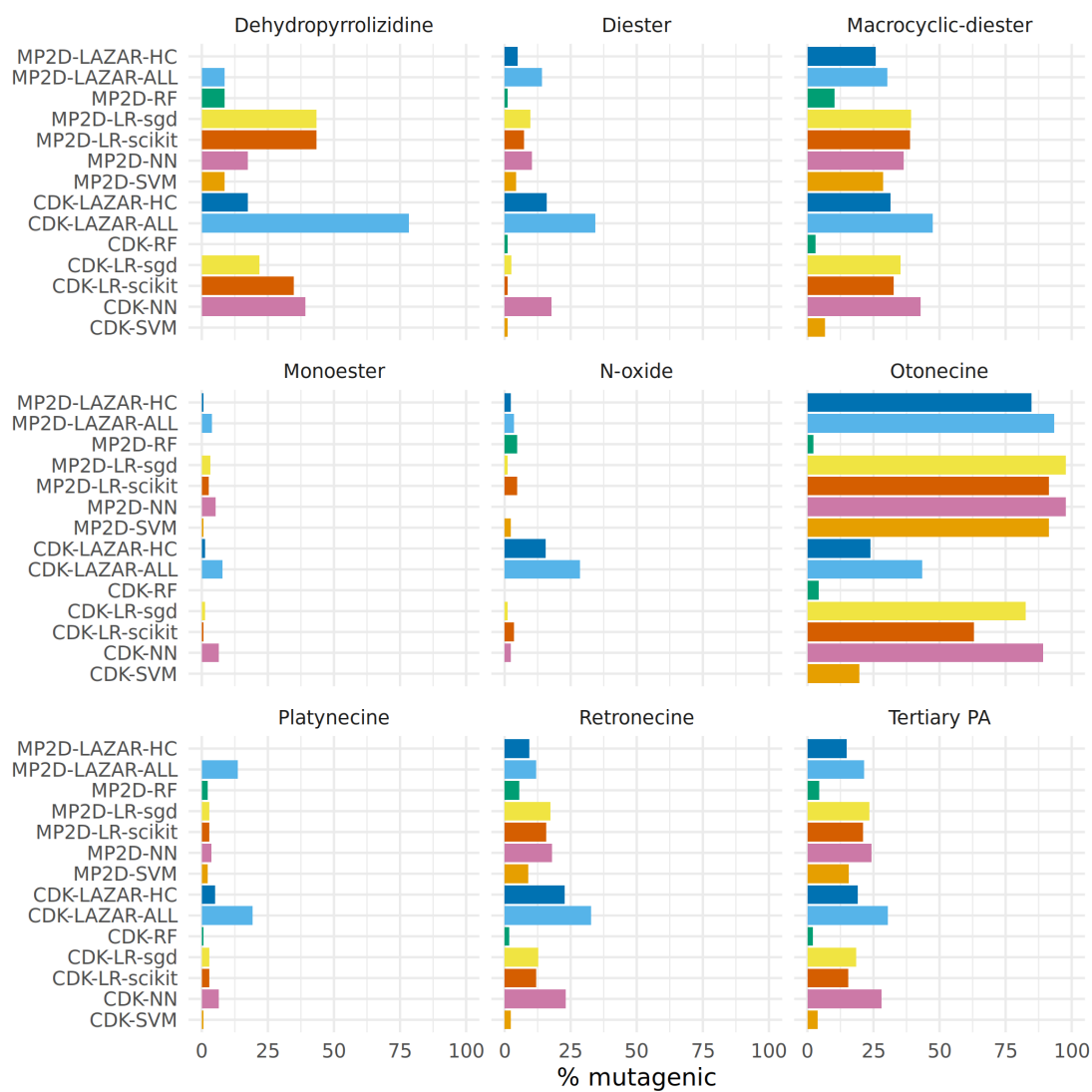


Figure 6: Summary of pyrrolizidine alkaloid predictions

303 (65%, 407/623), followed by Macrocyclic diesters (31%, 1042/3356), Dehydropy-
304 rrolizidine (27%, 74/268), Tertiary PAs (19%, 1201/6307) and Retronecines (15%,
305 762/5054).

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

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

310 Necine base: Platynecine < Retronecine « Otonecine

311 Necic acid: Monoester < Diester « Macrocyclic diester

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

313 Discussion

314 Data

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

321 Algorithms

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

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 the crossvalidation accuracies of all models 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 might just approximate experimental errors better than more robust models.

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.

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

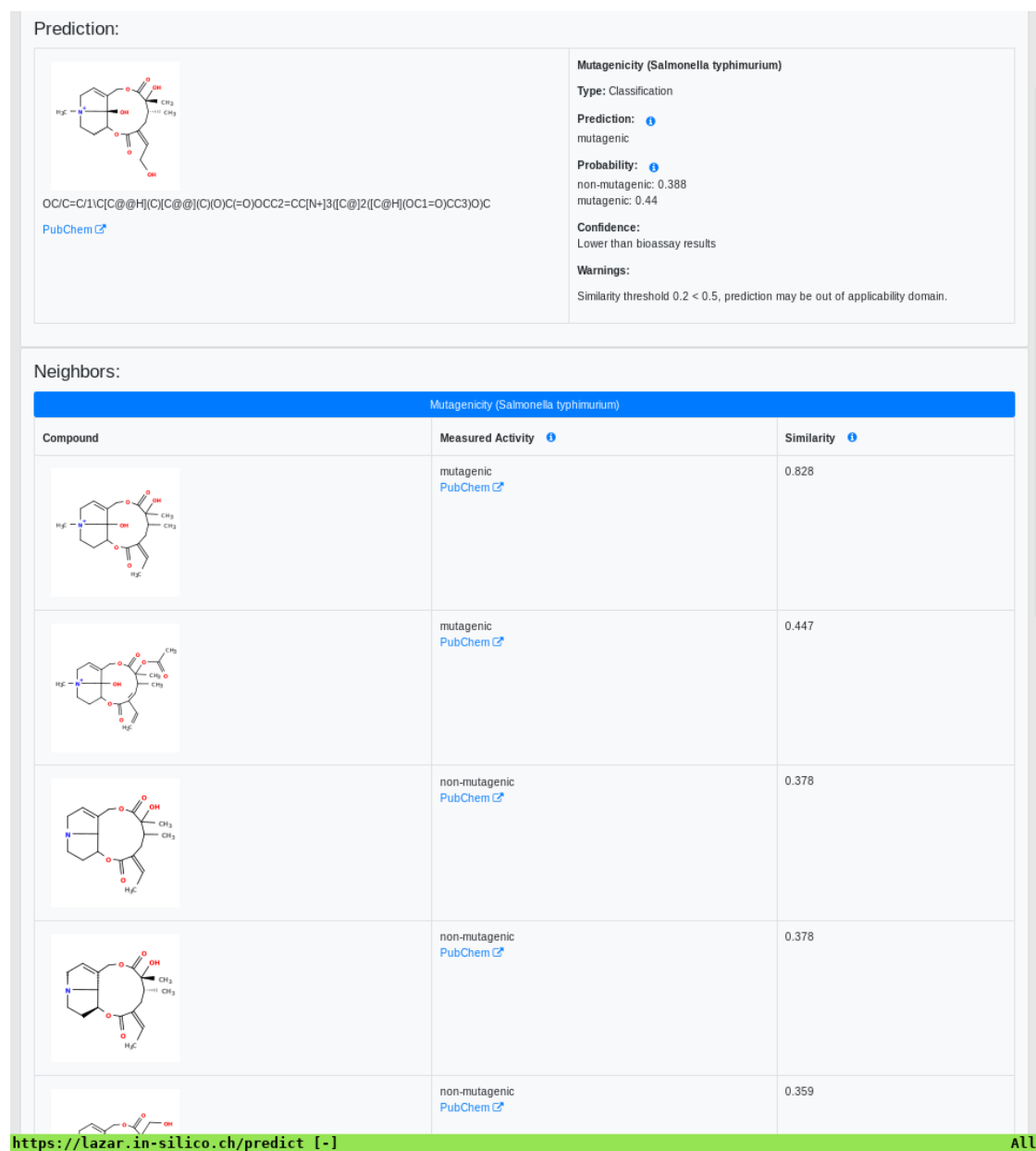


Figure 7: Lazar screenshot of 12,21-Dihydroxy-4-methyl-4,8-secosenecinonan-8,11,16-trione mutagenicity prediction

351 Descriptors

352 This study uses two types of descriptors for the characterisation of chemical structures:

353 *MolPrint2D* fingerprints (MP2D, Bender et al. (2004)) use atom environments (i.e.
354 connected atom types for all atoms in a molecule) as molecular representation, which
355 resembles basically the chemical concept of functional groups. MP2D descriptors are
356 used to determine chemical similarities in the default **lazar** settings, and previous ex-
357 periments have shown, that they give more accurate results than predefined fingerprints
358 (e.g. MACCS, FP2-4).

359 *Chemistry Development Kit* (CDK, Willighagen, Mayfield, and Alvarsson (2017)) descrip-
360 tors were calculated with the PaDEL graphical interface (Yap (2011)). They include 1D
361 and 2D topological descriptors as well as physical-chemical properties.

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

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

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

370 MolPrint2D fragments are generated dynamically from chemical structures and can be
371 used to determine if a compound contains structural features that are absent in training
372 data. This feature can be used to determine applicability domains. CDK descriptors
373 contain in contrast a predefined set of descriptors with unknown toxicological relevance.

374 MolPrint2D fingerprints can be represented very efficiently as sets of features that are
375 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.

Pyrrolizidine alkaloid mutagenicity predictions

Algorithms and descriptors

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

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

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

Necic acid

The rank order of the necic acid is comparable in all models. PAs from the monoester type had the lowest genotoxic potential, followed by PAs from the open-ring diester type. PAs with macrocyclic diesters had the highest genotoxic potential. The result fits well with current state of knowledge: in general, PAs, which have a macrocyclic diesters as necic acid, are considered to be more toxic than those with an open-ring diester or monoester (EFSA (2011), ("Pyrrolizidine Alkaloids–Genotoxicity, Metabolism Enzymes,

400 Metabolic Activation, and Mechanisms” 2004), Ruan2014b). This was also confirmed by
401 more recent studies, confirming that macrocyclic- and open-diesters are more genotoxic
402 *in vitro* than monoesters (Hadi et al. (2021); Allemang et al. (2018), Louisse et al.
403 (2019)).

404 **Necine base**

405 In the rank order of necine base PAs, platynecine is the least mutagenic, followed by
406 retronecine, and otonecine. Saturated PAs of the platynecine-type are generally ac-
407 cepted to be less or non-toxic and have been shown in *in vitro* experiments to form
408 no DNA-adducts ((“Pyrrolizidine Alkaloid-Derived Dna Adducts as a Common Biolog-
409 ical Biomarker of Pyrrolizidine Alkaloid-Induced Tumorigenicity” 2013)). In literature,
410 otonecine-type PAs were shown to be more toxic than those of the retronecine-type
411 ((“Assessment of Pyrrolizidine Alkaloid-Induced Toxicity in an in Vitro Screening Model”
412 2013)).

413 **Modifications of necine base**

414 The group-specific results reflect the expected relationship between the groups: the low
415 mutagenic potential of N-oxides and the high potential of Dehydropyrrolizidines (DHP)
416 (Chen, Mei, and Fu (2010)).

417 Dehydropyrrolizidines are regarded as the toxic principle in the metabolism of PAs, and
418 known to produce protein- and DNA-adducts (Chen, Mei, and Fu (2010)). None of
419 the models did not meet this expectation and predicted the majority of DHP as non-
420 mutagenic. However, the following issues need to be considered. On the one hand, all
421 DHP were outside of the stricter applicability domain of MP2D laser. This indicates
422 that they are structurally very different than the training data and might be out of the
423 applicability domain of all models based on this training set. In addition, DHP has two

424 unsaturated double bonds in its necine base, making it highly reactive. DHP and other
425 comparable molecules have a very short lifespan, and usually cannot be used in *in vitro*
426 experiments.

427 PA N-oxides are easily conjugated for extraction, they are generally considered as detox-
428 ification products, which are *in vivo* quickly renally eliminated (Chen, Mei, and Fu
429 (2010)).

430 Overall the low number of positive mutagenicity predictions was unexpected. PAs are
431 generally considered to be genotoxic, and the mode of action is also known. Therefore,
432 the fact that some models predict the majority of PAs as not mutagenic seems contradic-
433 tory. To understand this result, the experimental basis of the training dataset has to be
434 considered. The training dataset is based on the *Salmonella typhimurium* mutagenicity
435 bioassay (Ames test). There are some studies, which show mutagenicity of PAs in the
436 Ames test (Chen, Mei, and Fu (2010)). Also, Rubiolo et al. (1992) examined several
437 different PAs and several different extracts of PA-containing plants in the AMES test.
438 They found that the Ames test was indeed able to detect mutagenicity of PAs, but in
439 general, appeared to have a low sensitivity. The pre-incubation phase for metabolic
440 activation of PAs by microsomal enzymes was the sensitivity-limiting step. This could
441 very well mean that the low sensitivity of the Ames test for PAs is also reflected in the
442 investigated models.

443 Conclusions

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

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

453 Thus the practical question how to choose model predictions in the absence of experi-
454 mental data remains open. Tensorflow predictions do not include applicability domain
455 estimations and the rationales for predictions cannot be traced by toxicologists. Trans-
456 parent models like **lazar** may have an advantage in this context, because they present
457 rationales for predictions (similar compounds with experimental data) which can be
458 accepted or rejected by toxicologists and provide validated applicability domain estima-
459 tions.

460 References

461 Allemang, Ashley, Catherine Mahony, Cathy Lester, and Stefan Pfuhler. 2018. “Rela-
462 tive Potency of Fifteen Pyrrolizidine Alkaloids to Induce Dna Damage as Measured by
463 Micronucleus Induction in Heparg Human Liver Cells.” *Food and Chemical Toxicology*
464 121: 72–81. <https://doi.org/https://doi.org/10.1016/j.fct.2018.08.003>.

465 “Assessment of Pyrrolizidine Alkaloid-Induced Toxicity in an in Vitro Screening Model.”
466 2013. *J. Ethnopharmacol.*, no. 150: 560–7.

467 Bender, Andreas, Hamse Y. Mussa, Robert C. Glen, and Stephan Reiling. 2004. “Molec-
468 ular Similarity Searching Using Atom Environments, Information-Based Feature Selec-
469 tion, and a Naïve Bayesian Classifier.” *Journal of Chemical Information and Computer*
470 *Sciences* 44 (1): 170–78. <https://doi.org/10.1021/ci034207y>.

471 Benigni, R., and A. Giuliani. 1988. “Computer-assisted Analysis of Interlaboratory

Ames Test Variability.” *Journal of Toxicology and Environmental Health* 25 (1): 135–48.
<https://doi.org/10.1080/15287398809531194>.

Chen, T., N. Mei, and P. P. Fu. 2010. “Genotoxicity of Pyrrolizidine Alkaloids.” *J. Appl. Toxicol.*, 183–96.

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

———. 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.

Hadi, Naji Said Aboud, Ezgi Eyluel Bankoglu, Lea Schott, Eva Leopoldsberger, Vanessa Ramge, Olaf Kelber, Hartwig Sievers, and Helga Stopper. 2021. “Genotoxicity of Selected Pyrrolizidine Alkaloids in Human Hepatoma Cell Lines Hepg2 and Huh6.” *Mutation Research/Genetic Toxicology and Environmental Mutagenesis* 861-862: 503305.
<https://doi.org/https://doi.org/10.1016/j.mrgentox.2020.503305>.

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

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

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

497 Kazius, J., R. McGuire, and R. Bursi. 2005. "Derivation and Validation of Toxicophores
498 for Mutagenicity Prediction." *J Med Chem*, no. 48: 312–20.

499 Langel, D., D. Ober, and Pelser P. B. 2011. "The Evolution of Pyrrolizidine Alkaloid
500 Biosynthesis and Diversity in the Senecioneae," no. 10: 3–74.

501 Louisse, Jochem, Deborah Rijkers, Geert Stoopen, Wendy Jansen Holleboom, Mona
502 Delagrangé, Elise Molthof, Patrick P. J. Mulder, Ron L. A. P. Hoogenboom, Marc Au-
503 debert, and Ad A. C. M. Peijnenburg. 2019. "Determination of Genotoxic Potencies of
504 Pyrrolizidine Alkaloids in Heparg Cells Using the H2AX Assay." *Food and Chemical*
505 *Toxicology* 131: 110532. <https://doi.org/https://doi.org/10.1016/j.fct.2019.05.040>.

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

508 Mattocks, AR. 1986. *Chemistry and Toxicology of Pyrrolizidine Alkaloids*. Academic
509 Press.

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

513 "Pyrrolizidine Alkaloid-Derived Dna Adducts as a Common Biological Biomarker of
514 Pyrrolizidine Alkaloid-Induced Tumorigenicity." 2013. *Chem Res. Toxicol.*, no. 26:
515 1384–96.

516 "Pyrrolizidine Alkaloids–Genotoxicity, Metabolism Enzymes, Metabolic Activation, and
517 Mechanisms." 2004. *Drug Metab. Rev.*, no. 36: 1–55.

518 Rubiolo, P., L. Pieters, M. Calomme, C. Bicchi, A. Vlietinck, and D. Vanden
519 Berghe. 1992. "Mutagenicity of Pyrrolizidine Alkaloids in the Salmonella Ty-
520 phimurium/Mammalian Microsome System." *Mutation Research*, no. 281: 143–47.

521 Schöning, Verena, Felix Hammann, Mark Peinl, and Jürgen Drewe. 2017. "Editor's
 522 Highlight: Identification of Any Structure-Specific Hepatotoxic Potential of Different
 523 Pyrrolizidine Alkaloids Using Random Forests and Artificial Neural Networks." *Toxicol.*
 524 *Sci.*, no. 160: 361–70.

525 Weininger, David, Arthur Weininger, and Joseph L. Weininger. 1989. "SMILES. 2.
 526 Algorithm for Generation of Unique Smiles Notation." *J. Chem. Inf. Comput. Sci.*, no.
 527 29: 97–101. <https://doi.org/https://doi.org/10.1021/ci00062a008>.

528 Willighagen, E. L., J. W. Mayfield, and J. et al. Alvarsson. 2017. "The Chemistry
 529 Development Kit (Cdk) V2.0: Atom Typing, Depiction, Molecular Formulas, and Sub-
 530 structure Searching." *J. Cheminform.*, no. 9(33). [https://doi.org/https://doi.org/10.](https://doi.org/https://doi.org/10.1186/s13321-017-0220-4)
 531 [1186/s13321-017-0220-4](https://doi.org/https://doi.org/10.1186/s13321-017-0220-4).

532 Yap, CW. 2011. "PaDEL-Descriptor: An Open Source Software to Calculate Molecular
 533 Descriptors and Fingerprints." *Journal of Computational Chemistry*, no. 32: 1466–74.