

## Motivation

Biodegradability describes the capacity of substances to be mineralized by free living bacteria. It is a crucial property to estimate a compounds long-term impact on the environment. In order to substitute laborious experimental testing, it is of interest to reliably predict biodegradability. Yet, this endpoint is difficult to model due to inconsistency of available experimental data. Our approach makes use of OCHEM's [2] rich supply of machine learning methods and descriptor sets to build classification models for biodegradability.

## Data Set

The initial dataset has been arranged from several sources:

1. ~1400 measurements extracted from the CHRIP (Chemical Risk Information Platform) and ECHA (European Chemical Agency) database.
2. ~1500 measurements assembled by Cheng at al.[1] from the Japanese NITE database and the BIOWIN data set [3].
3. ~60 measurements of fragrances gathered from various online resources.

## Data Validation

The preliminary dataset has been manually cleaned, excluding compounds:

- listed as oligomers or mixtures
- having ambiguous structures
- with inconsistent biodegradability values
- with duplicates
- causing descriptor calculation errors

The final dataset consists of **1971** compounds out of which **732** are readily biodegradable and **1239** are not readily biodegradable and is freely available in OCHEM.

## Methods

**Model Validation:** Stratified Bagging Validation has been applied. The optimal bagsize was analysed in the range between 64 and 512 for all initial models. 64 bags emerged to yield the best performing models. It was used for validation for all further models.

**Parameter Optimization:** For all machine learning methods, parameter optimization was performed on the initial models using stratified bagging validation with 64 bags.

**Model Performance:** Balanced accuracy is defined as the arithmetic mean of sensitivity and specificity. It was used as a performance measure.

## Initial Models

Models were generated for all combinations of machine learning methods and descriptors.  
best Ø perf. descriptor: Dragon6 (**83.4%**)  
best Ø perf. method: ASNN (**82.4%**)

82.8	81.8	82.5	83.2	82	82.9	81.3	ASNN
81.7	81.2	81.6	83.3	81.4	83.3	81	LibSVM
79.4	78.3	80.9	81.7	77.7	77.9	75.2	FSMLR
75.1	75.9	75.6	75.6	77.7	77.4	74.7	KNN
78.7	73.9	80.5	79.8	69.5	79.6	76.8	PLS
81.5	81.1	81.4	82.9	82.2	83.4	80.7	Weka.J48
83	80.6	83.1	82.6	82	82.8	80.5	Weka.RF
EState, AlogPs	GSFrag	ISIDA	Dragon6	Adriana	CDK	Chemaxon(7.4)	

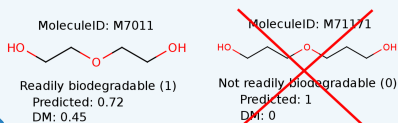
Williams Plot for best model (Weka J48 + CDK) with Bagging-STD as dist. to model.



## Exclusion Models

The distance to model (Bagging-STD) measure was used to exclude 51 compounds that were incorrectly predicted with high confidence in the previously best models.

Balanced Accuracy of Weka-J48 + CDK model:  
initial model (51 outliers left out): 85.7 %  
exclusion model (1920 cpds.): 86.6 %



best Ø perf. descriptor: CDK (**86%**)  
best Ø perf. method: ASNN (**85.8%**)

85.7	85.4	86.4	85.6	ASNN
84.4	84.4	85.8	85.4	LibSVM
84.1	84.4	84.9	<b>86.6</b>	WEKA.J48
85.7	85.2	84.5	86.2	WEKA.RF
AlogPs, EState	ISIDA	Dragon6	CDK	

## Subsets and Combinations

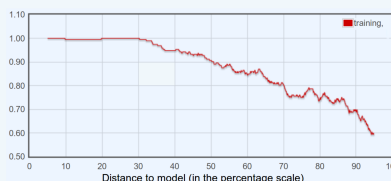
Subsets (3D and non-3D) of the Dragon6 and CDK descriptor sets were analysed, with non-3D descriptors yielding better performing models.

	3D	non-3D
Dragon6	84.2%	86.2%
CDK	82.3%	85.6%

The previously determined best working descriptor sets have been combined:

86.4	<b>87</b>	87	86.6	86.6	86.6	ASNN
85.6	84.8	86	86.7	86.1	86.1	LibSVM
85.1	85.8	85.1	86.8	85.1	85.4	WEKA.RF
85.6	86	85.2	86.2	85.2	86.5	WEKA.J48
EState, AlogPs, Dragon6	EState, AlogPs, CDK	ISIDA, Dragon6	ISIDA, CDK	Dragon6, CDK (3D)	Dragon6, CDK (non-3D)	

Williams plot for best model (ASNN + EState, AlogPs, CDK) with Bagging-STD as distance to model.



## Regression model

For less than half of the dataset (767 compounds) concrete BOD values were available. From initial regression models, 34 outliers were excluded (p-values < 0.001).

Method	Descriptor	Q <sup>2</sup>
ASNN	ISIDA, Dragon	0.66
ASNN	non-3D Dragon6	0.66
ASNN	Dragon6	0.65

## Outlook

Ongoing work is concerned with analysis of the relationship between characteristic compound properties and biodegradation. Special focus is set on outlier analysis and identification of primarily responsible descriptors.

## References

- [1] Cheng, Feixiong et al. In Silico Assessment of Chemical Biodegradability In *Journal of Chemical Information and Modeling*, 52:3, 655-669, 2012
- [2] Sushko, Iuri et al. Online chemical modeling environment (OCHEM): web platform for data storage, model development and publishing of chemical information. In *Journal of Computer-Aided Molecular Design*, 25:6, 533 - 554, 2011
- [3] Tunkel, J. et al. Predicting Ready Biodegradability in the Japanese Ministry of International Trade and Industry In *Test. Environ. Toxicol. Chem.*, 19, 2478-2485, 2000