

Prediction of Metabolic Transformations using Cross Venn-ABERS Predictors

Staffan Arvidsson[†], Ola Spjuth[†], Lars Carlsson[‡], Paolo Toccaceli^{*}

[†]Department of Pharmaceutical Biosciences, Uppsala University [‡]Quantitative Biology, Discovery Sciences, Innovative Medicines & Early Development, AstraZeneca ^{*}Department of Computer Science, Royal Holloway, University of London

Contact: Staffan Arvidsson <u>staffan.arvidsson@farmbio.uu.se</u> Pharmaceutical Biosciences Uppsala University





About me

- PhD student supervisor Ola Spjuth
- Uppsala University
- Department of Pharmaceutical Biosciences
- Free time: Cyclist







- 1. Defining the problem setting
- 2. Describing reaction types using SMIRKS
- 3. Preprocessing
- 4. ML context
- 5. Results
- 6. Conclusions



1. Defining the problem setting

- 2. Defining SMIRKS as reaction types
- 3. Preprocessing
- 4. ML context
- 5. Results
- 6. Conclusions





Problem setting – Metabolism





Problem setting – Drug Metabolism





Predicting drug metabolism in drug discovery process

- Aid drug discovery process
 - Duration and intensity of pharmacological action
 - Metabolites might be biologically active
 - Provide early warnings \rightarrow decrease costs
- Use known biotransformations
 - MDL Metabolite Database¹ 73599 recorded biotransformations
 - Use Machine Learning to predict metabolism for new drug candidates

¹Elsevier MDL. Mdl metabolite database, 2005.

URL http://accelrys.com/products/ collaborative-science/databases/bioactivity-databases/biovia-metabolite.html.



1. Defining the problem setting

2. Describing reaction types using SMIRKS

- 3. Preprocessing
- 4. ML context
- 5. Results
- 6. Conclusions



Describing reactions - SMIRKS



SMIRKS is a language used for describing chemical reaction transformations in a generic way¹



[Reactant]>>[Product]

Where Reactant and Product are represented as SMILES or reaction graphs

¹ Chemical Information Systems Inc. DAYLIGHT. Smirks - a reaction transform language, 2008. URL http://www.daylight.com/dayhtml/doc/theory/theory.smirks.html.





- 1. Defining the problem setting
- 2. Describing reaction types using SMIRKS

3. Preprocessing

- 4. ML context
- 5. Results
- 6. Conclusions





Derived datasets

Dataset name	SMIRKS	Biotrans- formations	Class 0
Alkyl hydroxylation	[\$([C:1])]>>[C:1][OH]	17793	68%
Aromatic hydroxylation	[\$([c:1])]>>[c:1][OH]	14691	85%
Carboxylation	[\$([CH3:1])]>>[C:1](=O)O	12580	64%
Oxidation of tertiary amine	[\$([N;X3:1])]>>[N+:1][O-]	11040	97% 🖛
Aromatization	[\$([*;R;!a:1])]>>[a:1]	9518	25%



- 1. Defining the problem setting
- 2. Describing reaction types using SMIRKS
- 3. Preprocessing

4. ML context

- 5. Results
- 6. Conclusions



Machine learning context

- Observations z = (x, y)
 - x: Features derived using Signatures Descriptors¹ of height 1 to 3, total number of features was 112710
 - y ∈ {0, 1}, 0="no SMIRKS match", 1="SMIRKS match"
- Cross Venn-ABERS predictors
 - Using 10 Inductive Venn-ABERS predictors
 - Datasets taken fold-wise i.e. 10% in calibration set and 90% in proper training set

¹Jean-Loup Faulon, et al. The signature molecular descriptor using extended valence sequences in qsar and qspr studies. Journal of chemical information and computer sciences, 43(3):707–720, 2003.



Machine learning context cont.

- Evaluation using outer *k*-fold cross validation, *k*=10
- Scoring algorithm: SVM in Python Scikit-learn¹
 - RBF kernel
 - -C = 50, $\gamma = 0.002$ previously found optimal values²
- Study both (p_0, p_1) interval and single prediction derived using: $p = GM(p_1)/(GM(1 - p_0) + GM(p_1))$

 ¹F. Pedregosa, et al. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12:2825–2830, 2011.
²J. Alvarsson, et al. Benchmarking study of parameter variation when using signature fingerprints together with support vector machines. Journal of chemical information and modeling, 54(11):3211–3217, 2014.



- 1. Defining the problem setting
- 2. Describing reaction types using SMIRKS
- 3. Preprocessing
- 4. ML context
- 5. Results
- 6. Conclusions



Calibration results







Performance metrics

Dataset name	Log Loss	AUC
Alkyl hydroxylation	0.538	0.753
Aromatic hydroxylation	0.348	0.793
Carboxylation	0.410	0.881
Oxidation of tertiary amine	0.093	0.904
Aromatization	0.173	0.964





- 1. Defining the problem setting
- 2. Describing reaction types using SMIRKS
- 3. Preprocessing
- 4. ML context
- 5. Results
- 6. Conclusions



Conclusions

- Solve issue with skewed class distributions
 - Incorporate stratified sampling
 - Re-sampling of minority class
 - Mondrian classification in Conformal Prediction might be a better approach
- We achieve very good predictors in most of the datasets
 - Narrow prediction intervals for all datasets even though not always informative predictions



Questions?

