# Graceful handling of large imbalanced datasets using Conformal Prediction

**Ulf Norinder and Fredrik Svensson** 

 $ulfn@dsv.su.se\ ,\ f.svensson@ucl.ac.uk$ 





Proceedings: 3rd Skövde Workshop on Information Fusion Topics, pp 59-62, 2009. Utilizing Information on Uncertainty for *In Silico* Modeling using Random Forests Henrik Boström

AstraZeneca R&D Södertälje

Henrik Bostrom Dept. of Computer and Systems Sciences Stockholm University and Informatics Research Centre University of Skövde

Introducing Uncertainty in Predictive Modeling—Friend or Foe? Ulf Norinder<sup>\*,†,‡,⊥</sup> and Henrik Boström<sup>§</sup>

J. Chem. Inf. Model. 2012, 52, 2815-2822

Representing descriptors derived from multiple conformations as uncertain features for machine learning

Ulf Norinder • Henrik Boström

J Mol Model (2013) 19:2679-2685





"What we ideally would like to know is in fact that a particular prediction is derived from an <u>area of property space</u> from which <u>reliable predictions are to be expected</u>"

> Introducing Conformal Prediction in Predictive Modeling. A Transparent and Flexible Alternative to Applicability Domain Determination

Ulf Norinder,\*\*<sup>†</sup> Lars Carlsson,<sup>‡</sup> Scott Boyer,<sup>‡,⊥</sup> and Martin Eklund<sup>‡,§</sup>

J. Chem. Inf. Model. 2014, 54, 1596-1603





Proceedings: 3rd Skövde Workshop on Information Fusion Topics, pp 59-62, 2009.

Utilizing Information on Uncertainty for In Silico Modeling using Random Forests

> Henrik Boström Dept. of Computer and Systems Sciences Stockholm University and Informatics Research Centre University of Skövde

Ulf Norinder AstraZeneca R&D Södertälje

Proceedings: Stream KDD '10 First International Workshop on Novel Data Stream Pattern Mining Techniques, Washington, D.C. ,2010

#### **Conformal Prediction for Distribution-Independent** Anomaly Detection in Streaming Vessel Data

Rikard Laxhammar University of Skövde Skövde, Sweden & Security and Defense Solutions, Saab AB Järfälla, Sweden rikard.laxhammar@his.se

Göran Falkman University of Skövde Skövde, Sweden goran.falkman@his.se

Introducing Uncertainty in Predictive Modeling-Friend or Foe? Ulf Norinder\*, <sup>†,‡,⊥</sup> and Henrik Boström<sup>§</sup>

J. Chem. Inf. Model. 2012, 52, 2815-2822

Representing descriptors derived from multiple conformations as uncertain features for machine learning

Ulf Norinder · Henrik Boström

J Mol Model (2013) 19:2679-2685







Introducing Uncertainty in Predictive Modeling—Friend or Foe? Ulf Norinder\*\*,  $^{\uparrow,\uparrow,\perp}$  and Henrik Boström  $^{\$}$ 

J. Chem. Inf. Model. 2012, 52, 2815-2822

Representing descriptors derived from multiple conformations as uncertain features for machine learning

Ulf Norinder · Henrik Boström

J Mol Model (2013) 19:2679-2685





Conformal Prediction What is it good for ...?





### **Confidence Predictors**

Conformal Predictors





#### **Confidence Predictors**

#### Conformal Predictors

• Venn (-Abers) Predictors

#### Accurate Hit Estimation for Iterative Screening Using Venn–ABERS Predictors

Ruben Buendia\*<sup>†</sup> (b, Thierry Kogej<sup>‡</sup>, Ola Engkvist<sup>‡</sup>, Lars Carlsson<sup>‡§⊥</sup>, Henrik Linusson<sup>†</sup>, Ulf Johansson<sup>†</sup>, Paolo Toccaceli<sup>§</sup>, and Ernst Ahlberg<sup>\*∎</sup>

<sup>†</sup> Department of Information Technology, University of Borås, SE-501 90 Borås, Sweden

<sup>‡</sup> Discovery Sciences, AstraZeneca IMED Biotech Unit, SE-431 83 Mölndal, Sweden

<sup>§</sup> Department of Computer Science, Royal Holloway, University of London, Egham, Surrey TW20 0EX, United Kingdom <sup>II</sup> Data Science and Al, Drug Safety & Metabolism, AstraZeneca IMED Biotech Unit, SE-431 83 Mölndal, Sweden

J. Chem. Inf. Model., 2019, 59 (3), pp 1230–1237 DOI: 10.1021/acs.jcim.8b00724 Publication Date (Web): February 6, 2019 Copyright © 2019 American Chemical Society

\*E-mail: Ruben.Buendia@astrazeneca.com., \*E-mail: Ernst.AhlbergHelgee@astrazeneca.com. This article is part of the Machine Learning in Drug Discovery special issue. Cite this: J. Chem. Inf. Model. 2019, 59, 3, 1230-1237

Calibrated probabilities: p0 and p1

representing the minimum and maximum probability limits





# **Why Conformal Prediction?**

- Win situation
- Statistical guarantees (on validity)





#### If {Exchangeability} then {conformal predictors are <u>always</u> valid}

#### **Mathematical proof**

Mach Learn (2013) 92:349-376 DOI 10.1007/s10994-013-5355-6

Conditional validity of inductive conformal predictors

Vladimir Vovk

Vovk V, Gammerman A, Shafer G (2005) **Algorithmic learning in a random world**, Springer, New York

If 20 % prediction errors on **validity** acceptable ---> CP will give, <u>at most</u>, 20 % errors!!





# Conformal Prediction validity



**Binary classification** 

In conformal prediction:

If a classification contains the correct class it is correct

*both* = always correct, *empty* = always erroneous

Validity = % of correct classifications (for each class) Efficiency = % of single label classifications (right or wrong)





# Conformal Prediction Why Conformal Prediction?

- Win situation
- Statistical guarantees (on validity)
- CP is instance-based
- The risk is known up-front for the decision taken
- Applicability domain closely linked to model development CP strictly defines the level of similarity (conformity) needed No ambiguity anymore
- Gracefully handles (severely) imbalanced datasets Ratios of 1:100 – 1:1000

No need for over- or undersampling

• CP is a framework (almost any ML algorithm will work)





How does this work?



Determination, J. Chem. Inf. Model., 2014, 54, 596-1603



#### CP is a framework (almost any ML algorithm will work)

- ML algorithm must provide a ranking
- Use current models, descriptors, algorithms
- Add calibration set -

New examples in time







- (non-) similarity function  $\rightarrow$  (non-) conformity function in CP
- Compares new compounds to old (calibration) compounds
  - Defined by the user
  - Probability from the RF trees
  - Distance to decision plane in SVM
  - o (Random numbers)





- (non-) similarity function  $\rightarrow$  (non-) conformity function in CP
- Compares new compounds to old (calibration) compounds
  - Defined by the user
  - Probability from the RF trees
  - Distance to decision plane in SVM

#### **Ranking problem**

CP p-value

 $|\{i = 1, \dots, n : \beta_i \leq \beta_{new}\}| / (n+1) \geq \varepsilon$ 

 $\beta_i$  = probability for the calibration compound i

 $\beta_{new}$  = probability for the new test compound

n = number of calibration set compounds

 $\mathcal{E}$  = significance level (% acceptable errors)

The number of calibration set compounds with probabilities  $\leq$  probability for the new compound divided by (n+1) must be  $\geq \varepsilon$  to be assigned a class label





How does this work?



Determination, J. Chem. Inf. Model., 2014, 54, 596-1603



**Example: Predicting Toxicity** 

Imbalanced dataset (toxic minority class)

A binary RF classifier (100 trees) gives the output:

New compound to predict (is toxic)

0

OH

=0 ″0H

 $NH_2$ 

ΌΗ

OH

ÔН

0

ö

32 trees: toxic68 trees: non-toxic





#### **Example: Predicting Toxicity**

Calibration set, 6 toxic, 7 non-toxic compounds N trees predicting correct class



New compound to predict (is toxic)

32 trees: toxic 68 trees: non-toxic







Mondrian Conformal Prediction



#### **Example: Predicting Toxicity**

Calibration set, 6 toxic, 7 non-toxic compounds N trees predicting correct class







#### **Example: Predicting Toxicity**

Based on the similarity to the known examples in the calibration set:





**Example: Predicting Toxicity** 



Using 80% confidence level (0.2 significance level):

3/8 = 0.43 > 0.2 therefore the compound is assigned to the toxic class

New compound to predict (is toxic)

32 trees: toxic 68 trees: non-toxic



0/8 = 0.0 < 0.2 therefore the compound is not assigned to the non-toxic class



#### **Example: Predicting Toxicity**

Calibration set, 6 toxic, 7 non-toxic compounds N trees predicting correct class





New compound to predict (is toxic)

32 trees: toxic 68 trees: non-toxic



Several p-values (for each class): Use median pvalue



Aggregated Mondrian Conformal Prediction Mondrian Cross-Conformal Prediction



Several pairs of proper train and calibration sets

#### **Mondrian Cross-Conformal Prediction**



University

DRUG DISCOVERY INSTITUTE

# **Binary Mondrian Conformal Prediction p-values**



# **Binary Mondrian Conformal Prediction p-values**



In conformal prediction:

If a classification contains the correct class it is correct **both** = always correct, **empty** = always erroneous

Validity = % of correct classifications (for each class) Efficiency = % of single label classifications (right or wrong)





# PubChem Cytotox Assays

- Results from 16 high throughput cell viability (tox) screens from PubChem
- On average 0.8% toxic compounds

AID	Tested compounds	Toxic compounds	%active	ratio non-tox/tox
624418	386 360	524	0.14	736.3
504648	367 995	600	0.16	612.3
602141	359 040	1302	0.36	274.8
620	86 701	364	0.42	237.2
847	41 152	194	0.47	211.1
903	52 783	338	0.64	155.2
2275	29 938	193	0.64	154.1
588856	404 016	3018	0.75	132.9
1825	290 605	2259	0.78	127.6
2717	299 957	3181	1.06	93.3
648	86 121	924	1.07	92.2
719	84 841	937	1.10	89.5
1486	217 851	2408	1.11	89.5
463	56 465	706	1.25	79.0
430	62 627	1121	1.79	54.9
598	85 162	5139	6.03	15.6





# PubChem Cytotox Assays

- Results from 16 high throughput cell viability (tox) screens from PubChem
- On average 0.8% toxic compounds
- RDKit descriptors
- RF, 500 trees, ensemble of 100 models
- 80 % training set, 20 % external test set



Modelling compound cytotoxicity using conformal prediction and PubChem HTS data†

Fredrik Svensson,<sup>a</sup> Ulf Norinder<sup>b,c</sup> and Andreas Bender\*<sup>a</sup>







Validity of the predictions (test sets) at the 80% confidence level. Models are valid for both classes.







Accuracy of the single label predictions (test sets) at the 80% confidence level.



The accuracy is similar for both the active and the inactive class.



# PubChem & Hansen Datasets

- Four dataset of different sizes and class imbalances
- 10 % randomly selected training sets
- Signature descriptors of heights 0–2 for chemical structure characterization
- Support vector machines (SVM) C-SVC, RBF kernel, parameters C =50, gamma = 0.002
- Ensemble of 100 SVM models

Binary classification of imbalanced datasets using conformal prediction

Ulf Norinder\*, Scott Boyer

Journal of Molecular Graphics and Modelling 72 (2017) 256–265







# Size and imbalance differs considerably between the datasets.





#### Fraction predicted active and inactive compounds.



#### #compounds in **both** class & **empty** class



@acceptable significance level:

Results from new data  $\rightarrow$ 

- Many predictions in *empty* class → outside AD of current model → measure and update model
- Many predictions in *both* class

→ inside AD of current model → lack of information → add new information (descriptors), develop better model (classifier, algorithm)





#### @acceptable significance level (decided by the user)



If a classification contains the correct class it is correct **both** = always correct, **empty** = always erroneous





Efficiency = % of single label classifications (right or wrong)



#### Not over-optimistic models

#### Validity minority class



#### Validity majority class



#### Signif. level 0.2





# Acknowledgements

Dr. Andreas Bender, Cambridge Univ
Dr. Lars Carlsson, Royal Holloway, Univ London
Dr. Martin Eklund, Karolinska Institutet
Dr. Ola Spjuth, Uppsala Univ
Dr. Scott Boyer, former Swetox
Dr. Natalia Aniceto, Cambridge Univ
The Francis Crick Institute



