

## Stacked Generalization with Applicability Domain Outperforms simple QSAR on in vitro Toxicological Data

Ingrid Grenet, Kevin Merlo, Jean Paul Comet, Romain Tertiaux, David Rouquié, and Frederic Dayan

*J. Chem. Inf. Model.*, **Just Accepted Manuscript** • DOI: 10.1021/acs.jcim.8b00553 • Publication Date (Web): 08 Feb 2019

Downloaded from <http://pubs.acs.org> on February 12, 2019

### Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

# Stacked Generalization with Applicability Domain Outperforms simple QSAR on *in vitro* Toxicological Data

Ingrid Grenet,<sup>†,‡,§</sup> Kevin Merlo,<sup>\*,¶,§</sup> Jean-Paul Comet,<sup>†</sup> Romain Tertiaux,<sup>¶</sup> David  
Rouquié,<sup>‡</sup> and Frédéric Dayan<sup>¶</sup>

*University Côte d'Azur, I3S laboratory, UMR CNRS 7271, CS 40121, 06903 Sophia  
Antipolis Cedex France, Bayer SAS, 06903 Sophia Antipolis Cedex France, and Dassault  
Systèmes SE, 06906 Sophia Antipolis Biot France*

E-mail: Kevin.MERLO@3ds.com

## Abstract

The development of *in silico* tools able to predict bioactivity and toxicity of chemical substances is a powerful solution envisioned to assess toxicity as early as possible. To enable the development of such tools, the ToxCast program has generated and made publicly available *in vitro* bioactivity data for thousands of compounds. The goal of the present study is to characterize and explore the data from ToxCast in terms of Machine Learning capacity. For this, a large scale analysis on the entire database has been performed to build models to predict bioactivities measured in *in vitro* assays. Simple classical QSAR algorithms (ANN, SVM, LDA, Random Forest and Bayesian) were first

---

\*To whom correspondence should be addressed

<sup>†</sup>I3S laboratory

<sup>‡</sup>Bayer SAS

<sup>¶</sup>Dassault Systemes

<sup>§</sup>Equal contributor

1  
2  
3 applied on the data, and the results of these algorithms suggested that they do not seem  
4 to be well suited for datasets with a high proportion of inactive compounds. The study  
5 then showed for the first time that the use of an ensemble method named "Stacked  
6 generalization" could improve the model performance on this type of data. Indeed,  
7 for 61% of 483 models, the Stacked method led to models with higher performance.  
8 Moreover, the combination of this ensemble method with an applicability domain filter  
9 allows one to assess the reliability of the predictions for further compound prioritization.  
10 In particular we showed that for 50% of the models, the ROC score is better if we do  
11 not consider the compounds that are not within the applicability domain.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## 23 Introduction

24  
25  
26 Chemical risk evaluation is aiming at defining the safe conditions of use of chemicals for  
27 both human health and the environment. To do so, the various hazards of the compounds  
28 are determined in a series of *in vitro* assays and *in vivo* studies using rodent and non-rodent  
29 laboratory animals. However, these *in vivo* studies are expensive with regard to time, money  
30 and animals, and are not adapted for the evaluation of thousands of chemicals. Therefore,  
31 alternative solutions were envisioned to assess, as early as possible, the potential toxicity of  
32 compounds and to prioritize them for further testing. The Tox21 partnership between the  
33 US Environmental Protection Agency (US EPA), the National Institutes of Health and the  
34 National Toxicology Program proposed to address in an eponymous project these issues by  
35 using high-throughput screening (HTS) *in vitro* assays to determine bioactivities of selected  
36 chemical substances.<sup>1</sup> This project was expected to help in characterizing pathways of tox-  
37 icity based on chemical bioactivity profiles, in prioritizing compounds for further targeted  
38 toxicity testing, and in developing predictive models for toxicity.<sup>2</sup> The ToxCast program led  
39 by the EPA in response to the Tox21 suggestions is one of the initiatives which generated  
40 public data regarding bioactivity of more than 8000 compounds tested in up to 900 HTS  
41 assays.<sup>3</sup> The main goals of the ToxCast project were to identify *in vitro* assays that were  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 relevant to predict *in vivo* toxicity, to develop predictive models based on these assays, and  
4 to use these models to prioritize compounds for further screening and testing.<sup>4</sup> Since the be-  
5 ginning of ToxCast in 2007, several computational methods such as Machine Learning (ML)  
6 and other statistical analysis have been proposed in order to establish relationships between  
7 *in vitro* data and *in vivo* effects. For example, Judson<sup>5</sup> used a set of *in vitro* assays from  
8 ToxCast to develop classifiers of *in vivo* toxicity using different ML methods. He showed  
9 that some methods are more appropriate than others for classifying the available *in vivo*  
10 toxicity data, which was largely derived from marketed compounds which tend to be of low  
11 toxicity and thus highly imbalanced. Also, Martin<sup>6</sup> and Sipes<sup>7</sup> used Linear Discriminant  
12 Analysis algorithm to predict *in vivo* reproductive toxicology and developmental toxicology  
13 respectively, based on specific *in vitro* assays. Other approaches tried to link bioactivity  
14 to *in vivo* outcomes using statistical methods such as correlation,<sup>4</sup> unsupervised multivari-  
15 ate analysis<sup>8</sup> or linear additive model.<sup>9,10</sup> Finally, *in vitro* assays have been combined with  
16 chemical structures to predict *in vivo* toxicity in ML models.<sup>11,12</sup> Nonetheless, since HTS  
17 programs can also be costly and time-consuming, there is a real interest in developing *in*  
18 *silico* methods, in particular, Quantitative Structure-Activity Relationship (QSAR) models  
19 predicting compound activity based on the chemical structure information,<sup>13</sup> to limit the use  
20 of *in vitro* assays.<sup>14,15</sup> Such QSAR models are used by regulatory agencies for risk and safety  
21 assessment<sup>16</sup> and the Organization for Economic Cooperation and Development (OECD)  
22 proposed five principles to be met by a QSAR model to be used for regulatory purposes.<sup>17</sup>  
23 These principles require: (1) a defined endpoint, (2) an unambiguous algorithm, (3) a defined  
24 domain of applicability, (4) appropriate measures for model evaluation and (5) a mechanis-  
25 tic interpretation, if possible. Some QSAR models have already been developed to predict  
26 bioactivity measured in HTS assays for specific endpoints.<sup>18,19</sup> The Tox21 challenge in 2014  
27 specifically highlighted the need for such *in silico* tools by looking for good models which  
28 predict the activation of nuclear receptor and stress response pathways.<sup>14</sup> In this paper we  
29 explore how ML can be applied to ToxCast data, and what the optimum uses of those data  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 are in order to build the most accurate classifiers based on chemical structural information.  
4  
5 Since we did not want to focus on a specific endpoint, we performed a large scale analysis  
6  
7 using all the available data which, to our knowledge, has never been done before. Indeed,  
8  
9 several QSAR models were built to predict bioactivity measured in numerous assays using  
10  
11 different types of molecular descriptors and different learning algorithms. As classical QSAR  
12  
13 methods are sometimes not well suited for toxicity data, we evaluated if the combination  
14  
15 of different algorithms proposed by ensemble techniques could lead to more robust models.  
16  
17 Here we focus on the Stacked generalization technique<sup>20</sup> applied to toxicological data and  
18  
19 show that it results in an improvement of the model performance. Also, in order to assess  
20  
21 the relevance of the predictions the use of applicability domain (AD), as proposed by the  
22  
23 OECD principle 3 was evaluated. To our knowledge, this is the first time that Stacked  
24  
25 generalization method combined with applicability domain is applied to toxicological data.  
26  
27  
28

## 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

## Methods

### ToxCast data

Both chemical structure and bioactivity data used in this work come from the ToxCast database and are publicly available on the US EPA website<sup>1</sup>. Overall, we retrieved 8599 unique substances, with information for each substance corresponding to a chemical name, molecular formula, CAS registry number (CASRN), and Simplified Molecular Input Line Entry Systems (SMILES) code.

Bioactivity data corresponding to the results of 1192 HTS *in vitro* assays was provided as a hit matrix with values of 0, 1, -1 and NA for each pair of compound/assay respectively meaning “inactive”, “active”, “undetermined” and “non assigned”. A compound is determined as “active” in an assay if the half maximal activity concentration (AC50) could be measured.

---

<sup>1</sup><https://www.epa.gov/chemical-research/toxicity-forecaster-toxcasttm-data>

## Chemical descriptors computation

To build QSAR models, we computed chemical descriptors for all the compounds using their SMILES representation and two different software: PaDEL-Descriptor<sup>21</sup> which is open source, and Pipeline Pilot<sup>22</sup> developed by Dassault Systèmes BIOVIA. We computed 1D and 2D descriptors that respectively encode: (1) chemical composition such as in particular fragment counts and molecular formulae, and (2) topology determined from the molecular graph (number of rotatable bounds, number of rings).

With PaDEL-Descriptor, we computed 1444 descriptors; in order to keep all of them and study their relative importance for ML, we chose to remove compounds for which at least one descriptor value could not be computed. With Pipeline Pilot, 164 descriptors were computed. This constitutes the first step of the data processing workflow, see Figure 1-A. The full lists of descriptors are available in supplementary materials.

## From ToxCast raw data to processed datasets

We built one dataset for each ToxCast bioassay, characterized by molecular descriptors as input and the assay result as output. A total of 2384 datasets were generated (1192 with Pipeline Pilot descriptors and 1192 with PaDEL descriptors), corresponding to the 1192 ToxCast *in vitro* assays. For each dataset, we removed the compounds for which the assay value was not available (NA or -1) in order to build classifiers that can predict the positive or negative results of the assays Figure 1-A. As a consequence, 200 datasets containing only NA or -1 values were discarded. Finally, we ended up with 2184 datasets (1092 for each types of descriptors). Then, in order to discard irrelevant descriptors, we removed the ones that had a variance close to 0, and to limit redundancy in the dataset, we removed one of two correlated descriptors using a threshold of 0.8. Moreover, we kept only datasets that include at least 10 members of each class (i.e. 10 “active” compounds and 10 “inactive” ones) and that have at least 5 times more observations than descriptors. Finally, we ended up with 515 datasets with Pipeline Pilot descriptors and 414 datasets with PaDEL descriptors,

and, respectively denoted PLP datasets and PaDEL datasets (all the datasets are available in the supplementary material).

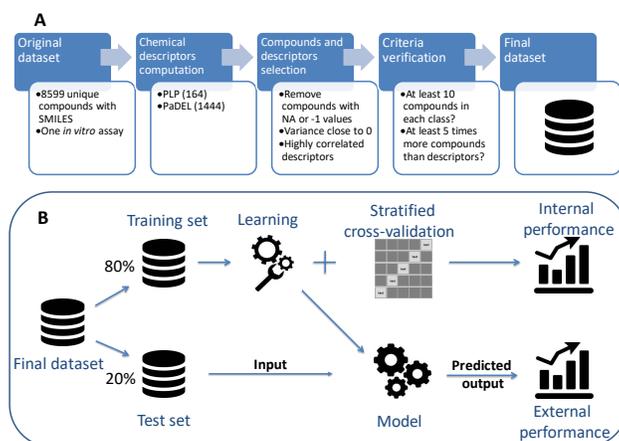


Figure 1: **Workflow of the data processing and learning procedure for one *in vitro* assay.** **A – Data processing:** For each *in vitro* assay, the original dataset contains the list of tested compounds with their structure from which chemical descriptors are computed. The compounds for which the reported value for the assay is NA or -1 are removed and the chemical descriptors are then selected according to their variance and correlation. Only datasets containing at least 10 compounds of the two classes (positive and negative) and at least 5 times more compounds than descriptors are kept. They correspond to the final datasets. **B – Learning procedure:** the dataset obtained in A is split into training set (80%) and test set (20%). The training set is used to learn the model using one of 5 different algorithms and a stratified 5-fold cross-validation is performed to get the internal performance of the model. The test set is then used to compute external performance of the model.

## Simple QSAR classifiers

The learning procedure applied to each dataset is described in Figure 1-B.

We built ML models for all the datasets obtained previously using five supervised ML algorithms. Three algorithms (single hidden layer Artificial Neural Network (ANN), Support Vector Machine (SVM) and Linear Discriminant Analysis (LDA)) were provided by the R packages NNET e1071, and MASS.<sup>23–25</sup> Two other algorithms (Forest multi-tree recursive partitioning classifier,<sup>26</sup> denoted hereafter as RPForest, and two-class Laplacian-modified Bayesian classifier,<sup>27</sup> denoted Bayesian) were implemented in Pipeline Pilot. For the 5 algorithms, we used the default parameters of the Pipeline Pilot machine learning collection.

We performed a stratified 5-fold internal cross-validation for each ML algorithm. To guarantee that each fold contained at least one member of each class of molecules, we first separated active and inactive molecules and the five folds were generated by randomly picking 20% from each class. For external validation, each dataset was randomly split into a training set (80%) and a test set (20%). The process was repeated 5 times and the average of the performance metrics was computed. Note that we chose to use random splitting by common use but other alternatives exist such as temporal split<sup>28</sup> or chemical clustering split.<sup>29</sup>

We used four metrics to evaluate the classification performance of the models, where TP (resp. TN) is the number of true positive (resp. negative), and FP (resp. FN) is the number of false positive (resp. negative): Sensitivity ( $TP/(TP + FN)$ ), Specificity ( $TN/(TN + FP)$ ), Balanced Accuracy, *BA* for short ( $(Sensitivity + Specificity)/2$ ) and ROC score (area under Receiver Operating Characteristic curve) where ROC curve is the plot of *Sensitivity* against  $(1 - Specificity)$ .

In order to perform 2-class predictions, the predicted continuous number that was returned by the algorithm was transformed into a binary one according to a threshold. There are several ways to determine this threshold and it has been shown that the traditional default method (threshold = 0.5) was unreliable for most of the datasets.<sup>30</sup> One of the best approach is to maximize the percentage of correctly classified observations. We chose an equivalent approach which is based on the minimization of the balanced error rate:  $\frac{1}{2} \left( \frac{TP}{TP+FN} + \frac{TN}{TN+FP} \right)$ .

## Stacked Generalization

The Stacked generalization technique is an ensemble method which consists of the training of a learning algorithm that combines the predictions of several other algorithms. We randomly split each dataset into one training set (Train 1) and two test sets (Test 1 and Test 2) in the following proportions: 60%, 24% and 16% (see ovals in Figure 2). On Train 1, we trained a Bayesian and an RPFforest algorithm to build models B1 and RP1 and we tested them on Test 1 to obtain predictions P1-b and P1-rp (see the blue workflow in Figure 2). We then

used these predictions P1-b and P1-rp as input descriptors with the outputs of Test 1 to train a so-called Stacked model. We chose a naïve Bayesian algorithm to build the final Stacked models (orange workflow of Figure 2) due to its performance and its ease of use in particular compared to Random Forest algorithm. We performed a 5-fold cross-validation repeated three times to find the best cutoff for the Stacked model. In some cases, it was impossible to compute a 5-fold cross-validation because of a small number of active compounds in Test 1: in these cases, the whole dataset was removed. Finally, we evaluated external predictive performance of the Stacked model on Test 2 (orange workflow of Figure 2). We applied this workflow for the 515 PLP datasets and generate 483 Stacked models. In order to compare the performance of the Stacked model with the simple QSAR classifiers, we merged Train 1 and Test 1 to train simple Bayesian and Random Forest learners and build models B2 and RF2 (see red workflow in Figure 2), and computed their external predictive performance on Test 2.

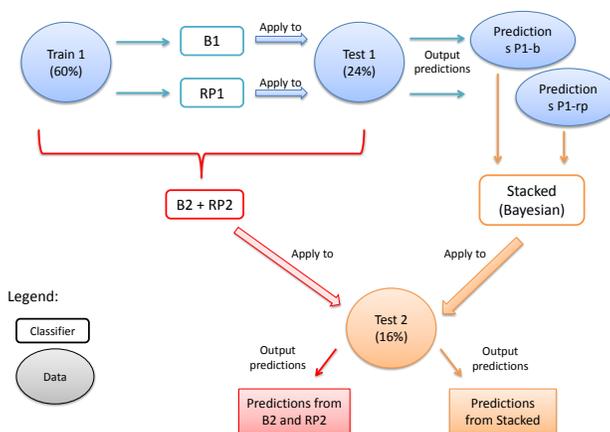


Figure 2: Principle of the Stacked generalization method.

**Blue Workflow:** A Bayesian B1 and a RPFforest RP1 models are built using Train 1 (60% of the dataset). The models are then used to make predictions P1-b and P1-rp on Test 1.

**Orange workflow:** the predictions P1-b and P1-rp are used to train a Stacked Bayesian model (Stacked).

**Red workflow:** By merging Train 1 and Test 1 to train a simple Bayesian model B2 and simple RPFforest model RP2, we are able to compare performance of B2, RP2 and the Stacked model on Test 2.

## Applicability Domain

We used two different approaches to compute the AD.

In the first approach, we performed a principal components analysis (PCA) directly applied on the chemical descriptors (1D and 2D) of the compounds used in the training set in order to reduce the space to only a few principal components (PCs). Basically, the PCs are computed to explain a minimum of 80% of the variance or a minimum of 10 components if 80% of the variance is explained with fewer components. Therefore, the number of principal components depends on the datasets. This allowed computation, for each PCA descriptor, of a range of acceptable values, as defined by the minimal and maximal PCA values observed in the training set. Then for a new compound, if the value of at least one of its PCA descriptors was out of the previously computed range, the compound was flagged “out of domain”. For each model we computed the average performance obtained on an external test set using either all the compounds or only the compounds that are in the AD (note that when the test set had only one active compound in the AD, the dataset was removed, leaving only 487 datasets).

In the second approach, we computed the average Euclidean distance from each molecule of the test set to its three closest compounds from the training set.<sup>31</sup> After sorting the compounds in ascending order of the average distance, we first cut the test set into blocks of 50 compounds. We then counted the number of good predictions in each of these disjoint blocks.

## Results and discussion

### Datasets are imbalanced

The total number of compounds in the PLP datasets ranged between 1115 and 7810 with average and median values of 3054 and 3362 compounds respectively (Figure 3-A). Overall

58% of the datasets (300/515) had less than 10% active compounds (positive in the *in vitro* assay) (Figure 3-B). Regarding PaDEL datasets, their size varied from 1391 to 7516, and there was no dataset with more than 50% active compounds (data not shown). Moreover, 66% of the datasets (275/414) had less than 10% active compounds. Based on a binary classification of the assay values, the datasets were highly imbalanced in favor of inactive compounds (negative in the *in vitro* assays).

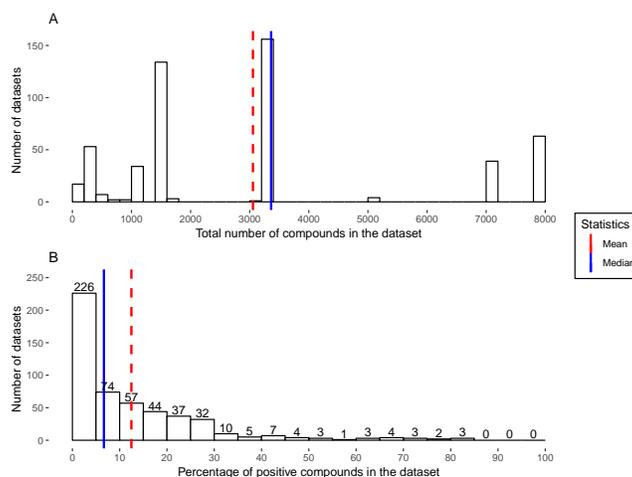


Figure 3: **A) Distribution of the 515 PLP datasets according to the number of compounds in the datasets. B) Distribution of the 515 PLP datasets according to the percentage of active compounds in the datasets. 58% of the datasets (300/515) have less than 10% of active compounds. Similar observations are made for the PaDEL datasets (data not shown).**

## Simple QSAR models

### Performance of 5 simple supervised learning methods

We trained five ML algorithms on the two types of datasets (PLP and PaDEL) and the mean of each performance metric (*Sensitivity*, *Specificity*, *ROCscore* and *BA*) was computed over the 515 and 414 datasets. Results are presented in Figure 4. For the two types of descriptors and the five algorithms, the four metric means were between 0.6 and 0.73 (except ANN-PaDEL *Sensitivity* which was at 0.52) and standard deviations were ranging between 0.06 and 0.19. Overall these results indicate that all algorithms performed similarly on this

type of data and have performance which are below our expectations.

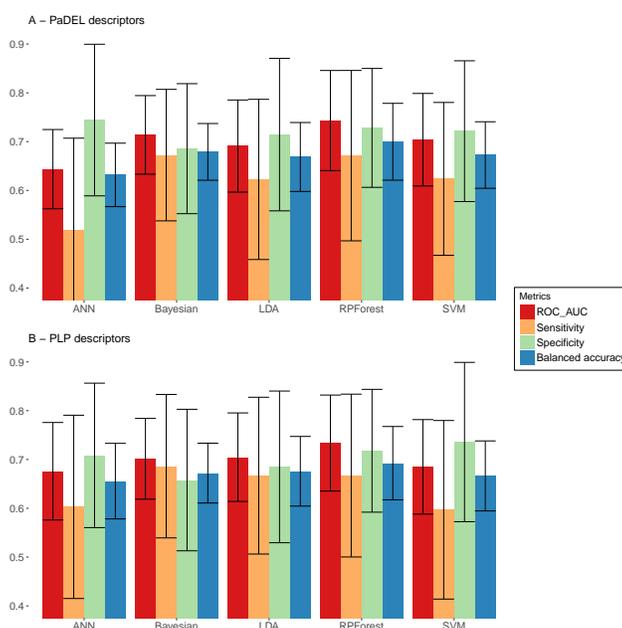


Figure 4: Comparison of performance metrics (*ROC* score, *Sensitivity*, *Specificity* and *BA*) after internal cross-validation for the 5 ML algorithms. A- Models using PaDEL descriptors and all datasets (414). B- Models using PLP descriptors and all datasets (515). There is no difference between the 5 algorithms used, both for PaDEL and PLP.

As we already showed that the data are imbalanced, we considered if this characteristic can explain these results. For all algorithms, plots of *BA* obtained for each dataset according to the percentage of active compounds display a “funnel” shape with *BA* variability decreasing when datasets were more balanced (Figure 5). In particular, we observed lower *BA* variability for datasets containing at least 10% of compounds in the minority class. For datasets with low percentage of positives, most of the *BA* variability *BA* is due to the variability of the *Sensitivity*, which depends on the number of positive compounds in the datasets: for a same percentage of positive compounds, the larger the number of positives, the higher the *Sensitivity*. We obtained similar results with PaDEL datasets (data not shown).

Figure 6 presents the variance of *ROC* score obtained with RPFforest algorithm trained

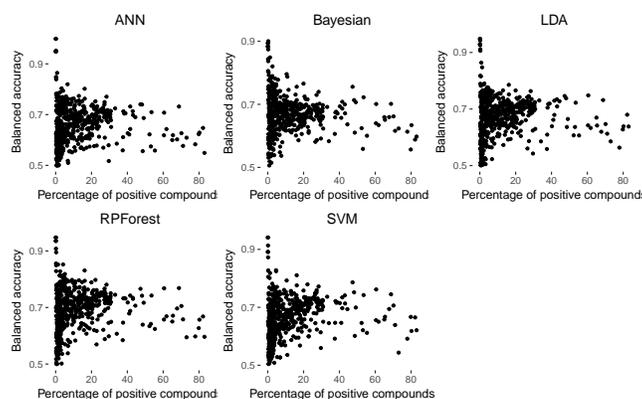


Figure 5: **Balanced accuracy according to percentage of positive compounds in PLP datasets for the 5 ML algorithms.** *BA* is getting stable when percentage of positive compounds in datasets increases. For datasets with low percentage of positives, most of the *BA* variability is due to the variability of the Sensitivity, which depends on the number of positive compounds in the datasets.

on PLP datasets according to the percentage of active compounds in the datasets. For each different range of percentage of positive compounds, we computed the variance of *ROC score* over all the datasets having a percentage of positive compounds in that range. Figure 6 shows that the variance of the *ROC score* tends to decrease when datasets are more balanced. We also see a cut-off: when datasets contain at least 10% of positive compounds, the associated variance is always below 0.0065. Similar results were obtained with PaDEL datasets (data not shown).

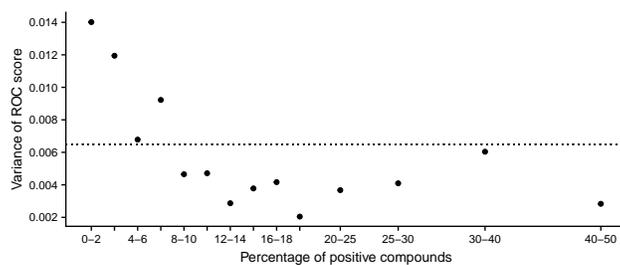


Figure 6: **Variance of *ROC score* according to the percentage of positive compounds after internal cross-validation for Random Forest models based on PLP descriptors.** Variance is lower than 0.0065 when percentage of positive compounds is greater than 10%.

These results suggest that the imbalanced nature of datasets has a negative impact

1  
2  
3 on model performance and that these classical learning methods are not suitable for highly  
4 imbalanced datasets. This finding is in agreement with previous work from different domains  
5 showing that most classifier algorithms assume a relatively balanced distribution of the  
6 data.<sup>32-35</sup> Moreover, the models with very few positive compounds will be characterized  
7 by a limited applicability domain regarding this class of compounds. Consequently, a new  
8 positive compound will have a low chance to be within the applicability domain and its  
9 associated prediction will be of low confidence. In order to be able to follow principle 3 of  
10 the OECD recommendations about the use of QSAR modelling for regulatory purposes and  
11 based on the fact that our results are in line with previous published work, the datasets  
12 containing less than 10% of compounds that belong to the minority class were excluded in  
13 the following analysis.  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

### 25 26 **Performance on more balanced datasets**

27  
28  
29 When focusing on datasets containing at least 10% compounds in the minority class, we  
30 ended up with 139 PaDEL and 215 PLP datasets. Figure 7 shows the means of each per-  
31 formance metric over these datasets. The means of all metrics were between 0.63 and 0.75  
32 and here again the results suggested that all the algorithms present similar performance  
33 (*BA* around 0.68). Interestingly, standard deviations were lower than previously: *BA* stan-  
34 dard deviation decreased in average from 0.07 to 0.04, *Sensitivity* from 0.16 to 0.10 and  
35 *Specificity* from 0.14 to 0.09. In order to quantify the advantage of using more balanced  
36 datasets, we performed a Student test: Table 1 shows the p-values of the t.test that com-  
37 pared the mean of the 4 metrics between the datasets containing strictly less than 10% of  
38 compounds of the minority class and the ones containing at least 10% of compounds, for the  
39 5 algorithms and the 2 types of descriptors used. Most of the metrics means are significantly  
40 different, meaning that the use of more balanced datasets has an impact on the models per-  
41 formance. In particular, the *Sensitivity* is always significantly increased suggesting that the  
42 use of more balanced datasets, and therefore more positive compounds in that case, helps in  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

the detection of true positives. Overall, the same tendencies were observed with both types of descriptors but PLP datasets contain fewer descriptors which confers several advantages. They are generally easier to understand as they are related to well-known physico-chemical properties such as molecular weight or solubility. Also, other advantages of building ML models with fewer descriptors are decreased model complexity, reduced chances of overfitting,<sup>36,37</sup> and decreased computational time. Therefore, we made the decision to focus only on PLP datasets for the rest of the study.

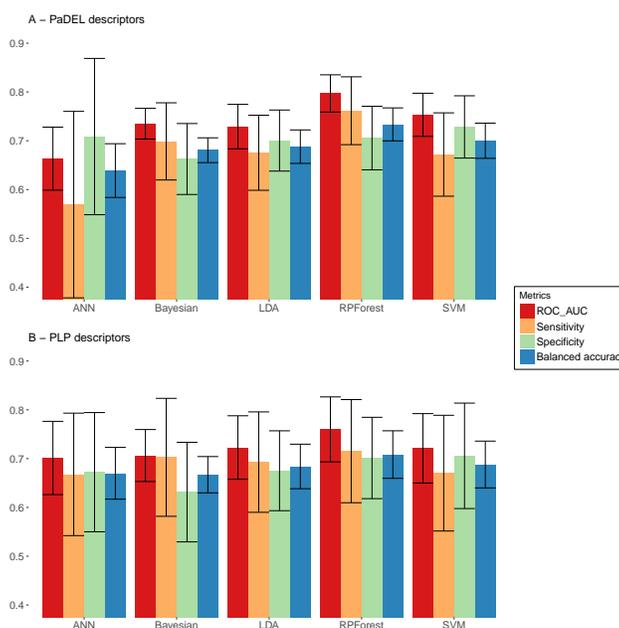


Figure 7: Comparison of performance metrics (*ROCscore*, *Sensitivity*, *Specificity* and *BA*) after internal cross-validation for the 5 ML algorithms. A- Models using PaDEL descriptors and datasets with at least 10 percent of compounds in the minority class (139). B- Models using PLP descriptors and datasets with at least 10% of compounds in the minority class (215). There is no difference between the 5 algorithms used, both for PaDEL and PLP. Standard deviations get smaller when we keep datasets with at least 10% of compounds in the minority class.

## External validation

Figure 8 presents the average *BA*, *Sensitivity* and *Specificity* obtained after external validation for the 5 ML algorithms on all 215 PLP datasets. We observed that, except for the Bayesian algorithm, *Sensitivity* was very low (under 0.4) and *Specificity* was high (greater

Table 1: p-values of Student test performed on the 4 metrics for the 5 algorithms and 2 types of descriptors, between the datasets that contain strictly less than 10% of active compounds and the datasets that contain at least 10% of active compounds. The p-values lower than 0.05 are in bold. Most of the metrics means are significantly different meaning that the use of more balanced datasets has an impact on the models performance.

Method	Descriptor type	ROC AUC	Balanced Accuracy	Sensitivity	Specificity
ANN	PADEL	<b><math>1.10 \times 10^{-4}</math></b>	$9.06 \times 10^{-2}$	<b><math>1.58 \times 10^{-4}</math></b>	<b><math>1.15 \times 10^{-3}</math></b>
Bayesian	PADEL	<b><math>6.32 \times 10^{-7}</math></b>	$5.96 \times 10^{-1}$	<b><math>5.95 \times 10^{-4}</math></b>	<b><math>2.00 \times 10^{-3}</math></b>
LDA	PADEL	<b><math>1.29 \times 10^{-13}</math></b>	<b><math>4.85 \times 10^{-7}</math></b>	<b><math>2.98 \times 10^{-9}</math></b>	$9.03 \times 10^{-2}$
RPForest	PADEL	<b><math>2.40 \times 10^{-23}</math></b>	<b><math>1.04 \times 10^{-15}</math></b>	<b><math>1.70 \times 10^{-22}</math></b>	<b><math>8.48 \times 10^{-4}</math></b>
SVM	PADEL	<b><math>8.76 \times 10^{-22}</math></b>	<b><math>3.55 \times 10^{-13}</math></b>	<b><math>3.85 \times 10^{-8}</math></b>	$3.73 \times 10^{-1}$
ANN	PLP	<b><math>1.66 \times 10^{-7}</math></b>	<b><math>1.20 \times 10^{-4}</math></b>	<b><math>3.30 \times 10^{-13}</math></b>	<b><math>7.59 \times 10^{-7}</math></b>
Bayesian	PLP	$2.11 \times 10^{-1}$	$7.82 \times 10^{-2}$	<b><math>2.55 \times 10^{-2}</math></b>	<b><math>1.47 \times 10^{-4}</math></b>
LDA	PLP	<b><math>3.79 \times 10^{-5}</math></b>	<b><math>1.69 \times 10^{-2}</math></b>	<b><math>6.80 \times 10^{-4}</math></b>	$1.79 \times 10^{-1}$
RPForest	PLP	<b><math>2.15 \times 10^{-8}</math></b>	<b><math>9.39 \times 10^{-6}</math></b>	<b><math>8.98 \times 10^{-10}</math></b>	<b><math>5.69 \times 10^{-3}</math></b>
SVM	PLP	<b><math>6.71 \times 10^{-15}</math></b>	<b><math>2.31 \times 10^{-10}</math></b>	<b><math>1.46 \times 10^{-17}</math></b>	<b><math>1.09 \times 10^{-4}</math></b>

than 0.8) which led to a *BA* around 0.6. Moreover, standard deviations of data obtained with ANN, LDA, RPForest and SVM algorithms were large, and we were unable to draw conclusions on the usefulness of this method based on these results. These external performance make us think that the 4 ML methods are not able to build good models, even when focusing only on datasets having at least 10% of compounds in the minority class. We hypothesized that they were too sensitive to imbalanced datasets as Mazurowski<sup>38</sup> showed for neural networks. Indeed, Mazurowski performed a large scale analysis on simulated imbalanced data and studied the impact of this characteristic on two neural network methods (classical backpropagation and particle swarm optimization). He concluded that even a small imbalance in the training set leads to a deterioration of the performance.

On the contrary the Bayesian algorithm seems to be more suited to imbalanced datasets since *BA*, *Sensitivity* and *Specificity* are greater than 0.6 with smaller standard deviations.

Overall, our results suggest that a unique algorithm alone was not able to generate models with sufficient performance.

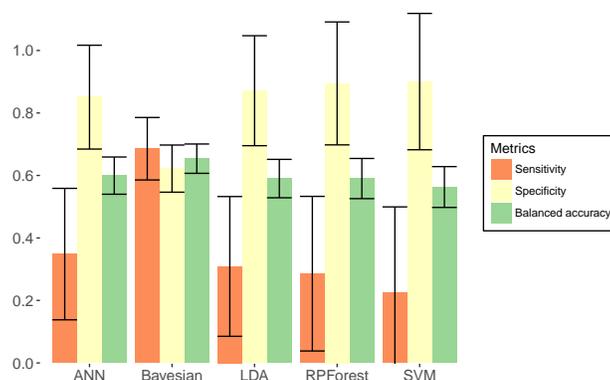


Figure 8: Comparison of performance metrics (sensitivity, Specificity and balanced accuracy) after external validation for the 5 algorithms on PLP datasets with at least 10% of compounds in the minority class (215 datasets). ANN, LDA, RPForest and SVM are not able to build good models whereas Bayesian algorithm seems to be more suitable.

## An ensemble technique: the Stacked generalization

The next step was to test if combining different types of algorithms would improve the results. This approach has been widely studied and led to different approaches so-called "ensemble techniques".<sup>39</sup> Intuitively, each learning algorithm makes a hypothesis to predict as well as possible a particular output. When the choice of these algorithms (and the underlying hypothesis) is not obvious, taking an "ensemble" of models trained on a same dataset allows the combination of multiple hypotheses in a unique model leading to higher performance.<sup>39,40</sup> Here, we used the Stacked generalization technique<sup>20,41</sup> using two methods based on very different internal representations of the training sets (instances and trees): Bayesian and RPForest algorithms which have been shown to lead to the best *Sensitivity* (Bayesian) and high *Specificity* (RPForest). The final stacked model was built using a naïve Bayes classifier. Note that we also built stacked models using 3 to 5 base models (by iteratively adding ANN, SVM and LDA to RPForest and Bayesian) on a subset of datasets (data not shown). Since the performance did not change significantly according to the number of models used, we therefore decided to keep only two base models for a matter of computing time.

### Stacked generalization vs. simple QSAR classifiers

We compared the predictive performance of the 483 simple Bayesian and RPFforest models (B2 and RP2) with the Stacked ones. Table 2 shows, for the three types of methods, the number of models that reach a certain value of *ROCscore*. First, less RP2 models are able to reach *ROCscore* greater than 0.6 compared to the B2 and Stacked ones and only 30 among the 483 have *ROCscore* above 0.8. Furthermore, if an equivalent number of B2 and Stacked models reach 0.6 and 0.7 values of *ROCscore*, when looking at higher and better performance (above 0.75 and 0.8), the Stacked method becomes clearly better than the Bayesian one. Finally, when we look at the method that lead to the highest *ROCscore* for each model, the Stacked is the winner for 61% models (294/483), followed by the Bayesian one with 30% of models (147/483) and the RPFforest with only 9% (42/483) (data not shown). These results confirm that, also in our application and based on the *ROCscore* values, the Stacked generalization method is able to build more models with good performance than simple QSAR algorithms.

Table 2: **Comparison of simple Bayesian B2 and Random Forest RP2 models with Stacked generalization models on the 483 PLP datasets.** The comparison of the number of models that reach a certain value of *ROCscore* shows that the Stacked generalization method is able to build more models with higher *ROCscore* than simple QSAR methods.

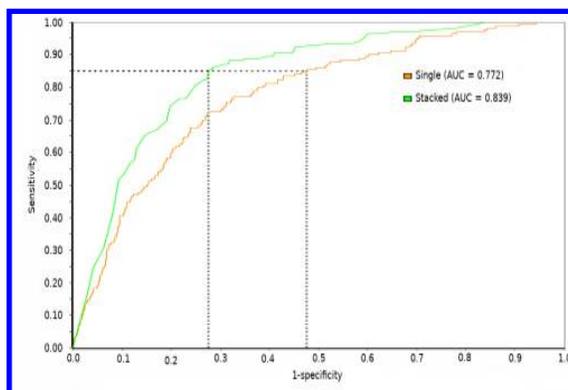
Method	<i>ROC</i> $\geq$ 0.60	<i>ROC</i> $\geq$ 0.70	<i>ROC</i> $\geq$ 0.75	<i>ROC</i> $\geq$ 0.80
Stacked	417	319	253	144
Bayesian	416	321	223	90
RPFforest	356	205	89	30

Figure 9 shows the ROC curves of B2 and stacked models for one particular assay (TOX21\_ERa\_BLA\_Antagonist\_ratio) measuring the expression of the Estrogen Receptor gene<sup>2</sup> and for which model performance were among the best. The associated dataset contains 7810 molecules, 13% of which were active in the *in vitro* assay. We observed that the ROC curve of the Stacked model is always above that of the B2 model. For this particular

<sup>2</sup>see <https://actor.epa.gov/dashboard>.

1  
2  
3 assay, the *ROC scores* were 0.84 and 0.77 for Stacked and B2 models respectively. Further-  
4  
5 more, since the ROC curve displays the *Sensitivities* and their corresponding *Specificities*  
6  
7 obtained for all threshold values between 0 and 1, we can choose a specific threshold according  
8  
9 to a desired *Sensitivity* or *Specificity*. As an example, a *Sensitivity* of 85% corresponds to  
10  
11 a *Specificity* of about 73% with the Stacked model and only of 52% with the Bayesian one,  
12  
13 see dotted lines in Figure 9. This again illustrates the difference of performance between the  
14  
15 two models and the ability of the Stacked model to detect more inactive compounds than  
16  
17 the Bayesian one, for the same number of active detected. Naturally, one can move this  
18  
19 threshold depending on the necessary stringency of the model output.  
20

21 The same analysis on other assays led to the same observations and conclusions (data not  
22  
23 shown).  
24



38 Figure 9: ROC curves obtained for models trained on the dataset of the assay  
39 TOX21\_ERa\_BLA\_Antagonist\_ratio with the two types of methods (Stacked generalization and  
40 simple Bayesian). ROC curve of the Stacked model is always above the one of simple Bayesian model. For a given *Sensitivity*  
41 of 85%, the Bayesian model detects 52% of the inactive molecules whereas the Stacked model detects 73% of them.  
42  
43  
44  
45  
46  
47

## 48 Focus on *in vitro* bioactivity assays proved to be linked to *in vivo* toxicity

49  
50 Since the bioactivity assays can be seen as an intermediate step towards the evaluation of the  
51  
52 *in vivo* toxicity, several works have relevantly focused on the link between ToxCast *in vitro*  
53  
54 assays and toxicity outcomes observed *in vivo*. In particular, in 2015 Liu and co-workers  
55  
56 built ML models that predict *in vivo* chronic toxicity observed in liver<sup>12</sup> based on either  
57  
58  
59  
60

1  
2  
3 chemical descriptors, bioactivity descriptors (i.e ToxCast *in vitro* assays) or a combination  
4 of both. This study has been extended in 2017 to 19 other organs.<sup>42</sup> They extracted in both  
5 studies the 36 (resp. 50) *in vitro* assays most frequently used in their models and which were  
6 supposed to be the most correlated with *in vivo* liver (resp. 19 other organs) toxicity.  
7  
8

9  
10 Since we built QSAR models for the majority of the ToxCast *in vitro* assays, we proposed  
11 here to focus on the ones that predict these assays. More precisely, among the 36 (resp.  
12 50) *in vitro* assays highlighted by Liu, we were able to build 25 (resp. 38) corresponding  
13 QSAR models. Because 11 models were common to both sets (25 and 38), we finally got 52  
14 QSAR models, using the simple algorithms and the stacked method. Table 3 summarizes the  
15 best *ROCscore* we obtained for these 52 QSAR models and the corresponding method used  
16 (simple Bayesian, RPFforest or Stacked generalization). For 71% of the assays (37/52), the  
17 Stacked generalization was the method leading to the best ROC score. Also, for 62% of the  
18 assays (32/52) the *ROCscore* was greater than 0.75 meaning that we were able to build good  
19 QSAR models to predict some of the *in vitro* assays highlighted by Liu. Altogether, these  
20 results show that the Stacked generalization method allows one to build QSAR classifier  
21 models that predict *in vitro* assays which have been previously shown to be associated to  
22 *in vivo* toxicity outcomes. This suggest that we could think about replacing all or part of  
23 these *in vitro* assays by *in silico* predictions and use these predictions as input of Liu's ML  
24 models for example.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

### 43 **Applicability domain information to reinforce the confidence of pre-** 44 **dictions** 45

46  
47 The OECD QSAR Validation principles recommend that a model should be used within  
48 its applicability domain (AD).<sup>17</sup> Several definitions of the AD have been proposed and the  
49 Setubal Workshop report<sup>43</sup> proposed the following one: "The AD of a (Q)SAR is the physico-  
50 chemical, structural, or biological space, knowledge or information on which the training set  
51 of the model has been developed, and for which it is applicable to make predictions for  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table 3: Most frequently used assays in Liu's models and the best *ROC* scores we obtained with either Stacked generalization or simple methods (Bayesian and RPFforest). Assays are sorted by decreasing *ROC* score. Stacked generalization method has the best *ROC* score for 71% of the assays and this score is greater than 0.75 for 62%.

Assay name	Liver	Others	ROC score	Bayesian(B)/Stacked(S) RPFforest (RP)
BSK KF3CT SRB down		×	0,855	S
TOX21 TR LUC GH3 Antagonist	×		0,842	S
APR HepG2 CellLoss 24h dn		×	0,842	B
TOX21 ERa BLA Antagonist ratio	×		0,839	S
ATG VDRE CIS up	×	×	0,836	S
ATG SREBP CIS up		×	0,834	S
ATG PBREM CIS up		×	0,834	S
ATG MRE CIS up		×	0,834	S
ATG PXR TRANS up		×	0,833	S
APR HepG2 MitoticArrest 72h up	×	×	0,830	S
TOX21 ERa LUC BG1 Antagonist		×	0,828	S
TOX21 PPARd BLA agonist ratio		×	0,827	S
TOX21 Aromatase Inhibition	×		0,821	S
ATG TGFb CIS up		×	0,819	S
ATG RARa TRANS up	×		0,815	S
APR HepG2 StressKinase 1h up		×	0,815	B
ATG RORE CIS up		×	0,812	B
BSK 3C Vis down		×	0,811	S
ATG PXRE CIS up	×	×	0,810	S
BSK BE3C SRB down		×	0,809	S
ATG NRF2 ARE CIS up	×		0,793	S
ATG PPRE CIS up	×	×	0,791	S
NVS GPCR hOpiate mu	×		0,791	S
ATG RXRb TRANS up	×		0,787	S
ATG C EBP CIS up		×	0,785	S
ATG LXRb TRANS up		×	0,780	B
ATG BRE CIS up	×	×	0,780	S
TOX21 PPARg BLA antagonist ratio		×	0,775	S
ATG Oct MLP CIS up	×		0,771	S
APR HepG2 MicrotubuleCSK 72h up		×	0,760	B
APR HepG2 MitoMembPot 1h dn		×	0,759	B
ATG NFI CIS up		×	0,754	B
ATG NF kB CIS up		×	0,749	B
NVS MP rPBR	×	×	0,749	RP
NVS ADME hCYP2C19	×	×	0,733	S
APR HepG2 CellCycleArrest 24h up		×	0,733	B
NVS NR hAR	×		0,727	B
ATG ERE CIS up		×	0,717	S
NVS ADME hCYP1A2	×		0,717	B
NVS NR mERa	×		0,708	S
ATG IR1 CIS up		×	0,708	S
APR HepG2 CellCycleArrest 72h dn		×	0,705	S
NVS NR hPXR	×		0,691	B
NVS NR hCAR Antagonist		×	0,689	S
NVS MP hPBR	×	×	0,682	S
TOX21 ERa LUC BG1 Agonist	×	×	0,676	S
NVS TR hNET	×		0,668	S
APR HepG2 NuclearSize 72h up	×		0,650	B
NVS NR hER	×	×	0,649	S
TOX21 TR LUC GH3 Agonist		×	0,646	S
APR HepG2 MitoMass 24h up		×	0,599	B
APR HepG2 MitoMass 72h up	×	×	0,537	S

new compounds. [...] Ideally, the (Q)SAR should only be used to make predictions within that domain by interpolation not extrapolation". Basically, the AD enable to estimate the similarity between training set and test sets.<sup>31,44</sup> In practice, there are different approaches to estimate if a molecule is within the AD or out of it. We chose to consider two of these approaches, first the position of the new molecule in the space described by descriptors and second the distance between the new molecule and the closest molecules of the training set. Table 4 presents the results of the first approach and shows the percentage of assays for

1  
2  
3 which performance metrics were higher when considering the test set with the molecules  
4 in AD only compared to the entire test set. For more than 50% of the assays, *ROCscore*  
5 and *BA* were higher when "out of AD" compounds were excluded. *Specificity* was higher  
6 for 88% of the assays but only 16% showed higher *Sensitivity* when using only "in AD"  
7 compounds. This could be explained by the low number of active compounds in the test  
8 sets leading to equivalent performance regardless of the molecules taken into consideration.  
9 Moreover, the number of compounds "out of AD" was very low for all test sets (between 5  
10 and 10) which led to similar test sets and also explains the results of Table 4. It could be  
11 interesting to test the importance of this AD approach for test sets with more "out of AD"  
12 compounds or to calculate the *BA* of compounds that are "out of AD".  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24  
25 **Table 4: Percentage of assays with higher**  
26 **performance when using only "in AD" com-**  
27 **ounds than when using all compounds.**

Metric	% of assays
ROC score	53.2
BA	53.0
Sensitivity	16.0
Specificity	88.3

28  
29  
30  
31  
32  
33  
34  
35  
36 We applied the second approach to the dataset of the assay TOX21\_ERa\_BLA\_Antagonist\_ratio  
37 which was reported by Liu to be linked to chronic liver toxicity.<sup>12</sup> The test set was composed  
38 of 1251 molecules, corresponding to 24 disjoint blocks of 50 molecules and one block of 51.  
39 Figure 10 displays the results of the AD analysis obtained for these 25 blocks for the sim-  
40 ple Bayesian and Stacked models for which *ROCscores* are 0.77 and 0.84 respectively. We  
41 observed that the percentage of good predictions decreased (from 100% to less than 50%)  
42 when the average distance to the compounds in the training set increased. This shows that  
43 we can be more confident in the predictions when the compounds are close to the ones of  
44 the training set. Moreover, the percentage of good predictions with the Stacked model was  
45 greater than with the Bayesian one (except for blocks 20 to 23). We also performed the  
46 analysis for the assay TOX21\_AR\_BLA\_Antagonist\_ratio and found comparable results  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

(data not shown).

This analysis shows that the use of the Stacked generalization method to predict *in vitro* assays from the chemical structure, in association with a tool to measure whether a new compound belongs to the AD, can lead to good and reliable predictions when this new compound is actually inside the AD.

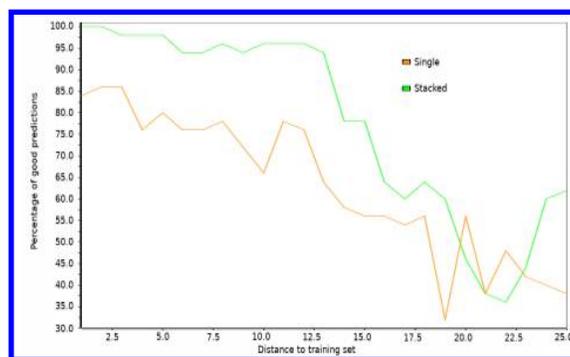


Figure 10: Number of good predictions made by the Stacked model and the simple Bayesian in all blocks of 50 predictions according to the average distance between the molecules of the test set and the three closest compounds of the training set. The dataset used corresponds to the assay TOX21\_Era\_BLA\_Antagonist\_ratio. Percentage of good predictions decreases when the average Euclidean distance to the training set increases. Percentage of good predictions of Stacked model is almost always greater than that of the Bayesian model.

## Conclusion

In this work, we performed a large scale analysis in order to characterize the Toxcast data which is generally very imbalanced with only a few active compounds for each *in vitro* assay. We first used classical learning methods, using two types of descriptors (PLP and PaDEL), in order to build models aiming at the prediction of results of ToxCast *in vitro* assays from chemical structures. The results indicated that all algorithms performed similarly and appeared to be below our expectations for this type of data. We then built Stacked models and, as recently shown by Madasamy,<sup>45</sup> demonstrated that this technique is more appropriate for this type of imbalanced data. In particular, we were able to build models

1  
2  
3 with good performance for *in vitro* assays that have previously been shown to be related to  
4 specific *in vivo* toxicity outcomes.<sup>12,42</sup>  
5  
6

7 Moreover, we demonstrated that the AD information is an important parameter to evaluate  
8 the reliability of predictions and can help to support decision making and prioritization.  
9  
10 Indeed, it could be a complementary and stringent filter allowing the selection of compounds  
11 for further *in vivo* testing. For example, molecules within the AD and predicted as toxic  
12 would be the first to be further characterized *in vivo*. For Toxcast data, combining Stacked  
13 generalization methods with an AD filter led to better classifier than did the classical learning  
14 method for Toxcast data.  
15  
16  
17  
18  
19  
20

21 Choices made in this study highlight directions for future work. First, instead of build-  
22 ing classifiers, one could use regression algorithms to predict the AC50 values obtained in  
23 each *in vitro* assay and not simply the “active” vs “inactive” labels. Here, by applying the  
24 threshold before the learning, we might have lost information. It could be interesting to see  
25 if more accurate models can be built by applying the threshold after the predictions. Second,  
26 we did not explore all the possible parameters for the different learners; these parameters  
27 can be tuned in order to increase the performance but the operation is time consuming.<sup>46,47</sup>  
28  
29 Regarding the stacked generalization method, we could think about exploring more combi-  
30 nations of customized base models and using other algorithms as meta-learner to build the  
31 stacked model in order to reach higher performance. Third, to complement the 1D and 2D  
32 molecular descriptors used in this paper, we could also use 3D descriptors and fingerprints  
33 which are known to be good for QSAR.<sup>48,49</sup> Then, to face the effect of the highly imbalanced  
34 data, we could implement data augmentation techniques<sup>50</sup> which balance the training set  
35 by adding or removing data and are supposed to increase model performance.<sup>51</sup> Finally, we  
36 could test other ensemble techniques such as bagging, boosting or bucket of models<sup>45</sup> to  
37 compare with the stacked generalization method. We actually already tried to use bagging  
38 with Random Forest and Bayesian algorithms on 15 datasets (having more than 1000 com-  
39 pounds and at least 30% of positive ones) and these first results showed that the bagging did  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 not lead to better performance than the stacked generalization. However, this work needs  
4 to be extended to all the datasets and more methods. More interestingly we would like to  
5 evaluate the ability of predicting *in vivo* toxicity by considering both structure and *in vitro*  
6 information, either by combining them as input descriptors or by chaining up two prediction  
7 steps.  
8  
9  
10  
11  
12

## 13 14 15 **Acknowledgement**

16  
17 The authors thank Gabriel Sarrazin and Moeka Shishido who provided insights and com-  
18 ments that greatly improved the manuscript. The authors would also like to show their  
19 gratitude to the Dassault Systèmes BIOVIA Drug Design Services team members (Celine  
20 Ferre, Tanguy Devos and Patrice Pistone) for their support. Editing of the manuscript by  
21 Elizabeth Shipp and Leo Bleicher was gratefully appreciated by the authors.  
22  
23  
24  
25  
26

27 This work was partly supported by the OSEO BioIntelligence program. I.G hold a doctoral  
28 fellowship from the Association Nationale de la Recherche Technique (ANRT CIFRE PhD  
29 funding).  
30  
31  
32

33 Finally, the authors acknowledge the reviewers for their comments and questions which have  
34 allowed us to improve the quality of the paper.  
35  
36  
37

38  
39 **Notes** The authors declare no conflict of interest.  
40  
41

## 42 43 **Supporting Information Available**

44  
45 The following files are available free of charge.  
46  
47

- 48 • supplementary1\_PaDELdescriptors.xlsx: list of PaDEL descriptors
- 49
- 50 • supplementary2\_PLPdescriptors.xlsx: list of PLP descriptors
- 51
- 52 • supplementary3\_datasets.zip: datasets used for learning
- 53
- 54
- 55

56 This material is available free of charge via the Internet at <http://pubs.acs.org/>.  
57  
58

## References

- (1) Council, N. R. *Toxicity Testing in the 21st Century*; National Academies Press: Washington, D.C., 2007.
- (2) Knudsen, T.; Martin, M.; Chandler, K.; Kleinstreuer, N.; Judson, R.; Sipes, N. Predictive models and computational toxicology. *Methods in molecular biology (Clifton, N.J.)* **2013**, *947*, 343–374.
- (3) Dix, D. J.; Houck, K. A.; Martin, M. T.; Richard, A. M.; Setzer, R. W.; Kavlock, R. J. The ToxCast Program for Prioritizing Toxicity Testing of Environmental Chemicals. *Toxicol. Sci.* **2007**, *95*, 5–12.
- (4) Judson, R. S.; Houck, K. A.; Kavlock, R. J.; Knudsen, T. B.; Martin, M. T.; Mortensen, H. M.; Reif, D. M.; Rotroff, D. M.; Shah, I.; Richard, A. M.; Dix, D. J. In Vitro Screening of Environmental Chemicals for Targeted Testing Prioritization: The ToxCast Project. *Environ. Health Perspect.* **2009**, *118*, 485–492.
- (5) Judson, R.; Elloumi, F.; Setzer, R. W.; Li, Z.; Shah, I. A comparison of machine learning algorithms for chemical toxicity classification using a simulated multi-scale data model. *BMC Bioinf.* **2008**, *9*, 241.
- (6) Martin, M. T.; Knudsen, T. B.; Reif, D. M.; Houck, K. A.; Judson, R. S.; Kavlock, R. J.; Dix, D. J. Predictive Model of Rat Reproductive Toxicity from ToxCast High Throughput Screening<sup>1</sup>. *Biol. Reprod.* **2011**, *85*, 327–339.
- (7) Sipes, N. S.; Martin, M. T.; Reif, D. M.; Kleinstreuer, N. C.; Judson, R. S.; Singh, A. V.; Chandler, K. J.; Dix, D. J.; Kavlock, R. J.; Knudsen, T. B. Predictive Models of Prenatal Developmental Toxicity from ToxCast High-Throughput Screening Data. *Toxicol. Sci.* **2011**, *124*, 109–127.
- (8) Shah, I.; Houck, K.; Judson, R. S.; Kavlock, R. J.; Martin, M. T.; Reif, D. M.;

- 1  
2  
3 Wambaugh, J.; Dix, D. J. Using nuclear receptor activity to stratify hepatocarcino-  
4 gens. *PloS one* **2011**, *6*, 1–11.  
5  
6  
7  
8 (9) Browne, P.; Judson, R. S.; Casey, W. M.; Kleinstreuer, N. C.; Thomas, R. S. Screening  
9 Chemicals for Estrogen Receptor Bioactivity Using a Computational Model. *Environ.*  
10 *Sci. Technol.* **2015**, *49*, 8804–8814.  
11  
12  
13  
14 (10) Kleinstreuer, N. C.; Ceger, P.; Watt, E. D.; Martin, M.; Houck, K.; Browne, P.;  
15 Thomas, R. S.; Casey, W. M.; Dix, D. J.; Allen, D.; Sakamuru, S.; Xia, M.; Huang, R.;  
16 Judson, R. Development and Validation of a Computational Model for Androgen Re-  
17 ceptor Activity. *Chem. Res. Toxicol.* **2017**, *30*, 946–964.  
18  
19  
20  
21  
22 (11) Thomas, R. S.; Black, M. B.; Li, L.; Healy, E.; Chu, T.-M.; Bao, W.; Andersen, M. E.;  
23 Wolfinger, R. D. A comprehensive statistical analysis of predicting in vivo hazard using  
24 high-throughput in vitro screening. *Toxicol. Sci.* **2012**, *128*, 398–417.  
25  
26  
27  
28 (12) Liu, J.; Mansouri, K.; Judson, R. S.; Martin, M. T.; Hong, H.; Chen, M.; Xu, X.;  
29 Thomas, R. S.; Shah, I. Predicting Hepatotoxicity Using ToxCast in Vitro Bioactivity  
30 and Chemical Structure. *Chem. Res. Toxicol.* **2015**, *28*, 738–751.  
31  
32  
33  
34 (13) Hansch, C. Quantitative structure-activity relationships and the unnamed science. *Acc.*  
35 *Chem. Res.* **1993**, *26*, 147–153.  
36  
37  
38 (14) Capuzzi, S. J.; Politi, R.; Isayev, O.; Farag, S.; Tropsha, A. QSAR Modeling of Tox21  
39 Challenge Stress Response and Nuclear Receptor Signaling Toxicity Assays. *Front. En-*  
40 *viron. Sci.* **2016**, *4*, 3.  
41  
42  
43 (15) Gadaleta, D.; Manganelli, S.; Roncaglioni, A.; Toma, C.; Benfenati, E.; Mombelli, E.  
44 QSAR Modeling of ToxCast Assays Relevant to the Molecular Initiating Events of  
45 AOPs Leading to Hepatic Steatosis. *J. Chem. Inf. Model* **2018**, *58*, 1501–1517.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 (16) Roy, K.; Kar, S.; Ambure, P. On a simple approach for determining applicability domain  
4 of QSAR models. *Chemom. Intell. Lab. Syst.* **2015**, *145*, 22–29.  
5  
6  
7  
8 (17) OCDE, *Guidance Document on the Validation of (Quantitative) Structure-Activity Re-*  
9 *lationship [(Q)SAR] Models*; OECD Series on Testing and Assessment; 2014; p 154.  
10  
11  
12 (18) Ng, H. W.; Doughty, S. W.; Luo, H.; Ye, H.; Ge, W.; Tong, W.; Hong, H. Development  
13 and Validation of Decision Forest Model for Estrogen Receptor Binding Prediction of  
14 Chemicals Using Large Data Sets. *Chem. Res. Toxicol.* **2015**, *28*, 2343–2351.  
15  
16  
17 (19) Zang, Q.; Rotroff, D. M.; Judson, R. S. Binary Classification of a Large Collection  
18 of Environmental Chemicals from Estrogen Receptor Assays by Quantitative Struc-  
19 ture–Activity Relationship and Machine Learning Methods. *J. Chem. Inf. Model.* **2013**,  
20 *53*, 3244–3261.  
21  
22  
23 (20) Wolpert, D. H. Stacked generalization. *Neural Netw.* **1992**, *5*, 241–259.  
24  
25  
26 (21) Yap, C. W. PaDEL-descriptor: An open source software to calculate molecular descrip-  
27 tors and fingerprints. *J. Comput. Chem.* **2011**, *32*, 1466–1474.  
28  
29  
30 (22) Dassault Systèmes BIOVIA, *Pipeline Pilot v. 17.2.0.1361*; 2012.  
31  
32  
33 (23) R Core Team, R: A Language and Environment for Statistical Computing. 2013; ISBN  
34 3-900051-07-0.  
35  
36  
37 (24) Meyer, D. Support Vector Machines. The Interface to libsvm in package e1071. 2001.  
38  
39  
40 (25) Venables, W. N.; Ripley, B. D. *Modern Applied Statistics with S*, 4th ed.; Springer:  
41 New York, 2002.  
42  
43  
44 (26) Breiman, L. Random Forests. *Machine Learning* **2001**, *45*, 5–32.  
45  
46  
47 (27) Xia, X.; Maliski, E. G.; Gallant, P.; Rogers, D. Classification of Kinase Inhibitors Using  
48 a Bayesian Model. *J. Med. Chem.* **2004**, *47*, 4463–4470.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 (28) Sheridan, R. P. Time-Split Cross-Validation as a Method for Estimating the Goodness  
4 of Prospective Prediction. *J. Chem. Inf. Model* **2013**, *53*, 783–790.  
5  
6  
7  
8 (29) Martin, T. M.; Harten, P.; Young, D. M.; Muratov, E. N.; Golbraikh, A.; Zhu, H.;  
9 Tropsha, A. Does Rational Selection of Training and Test Sets Improve the Outcome  
10 of QSAR Modeling? *J. Chem. Inf. Model* **2012**, *52*, 2570–2578.  
11  
12  
13  
14 (30) Freeman, E. A.; Moisen, G. G. A comparison of the performance of threshold criteria  
15 for binary classification in terms of predicted prevalence and kappa. *Ecol. Modell.* **2008**,  
16 *217*, 48–58.  
17  
18  
19  
20  
21 (31) Sheridan, R. P.; Feuston, B. P.; Maiorov, V. N.; Kearsley, S. K. Similarity to molecules  
22 in the training set is a good discriminator for prediction accuracy in QSAR. *J. Chem.*  
23 *Inf. Comput. Sci.* **2004**, *44*, 1912–1928.  
24  
25  
26  
27  
28 (32) Chawla, N. V.; Japkowicz, N.; Kotcz, A. Special issue on learning from imbalanced  
29 data sets. *ACM Sigkdd Explorations Newsletter* **2004**, *6*, 1–6.  
30  
31  
32  
33 (33) Fawcett, T.; Provost, F. Adaptive fraud detection. *Data Min. Knowl. Discov* **1997**, *1*,  
34 291–316.  
35  
36  
37  
38 (34) Kubat, M.; Holte, R. C.; Matwin, S. Machine learning for the detection of oil spills in  
39 satellite radar images. *Machine learning* **1998**, *30*, 195–215.  
40  
41  
42  
43 (35) Riddle, P.; Segal, R.; Etzioni, O. Representation design and brute-force induction in a  
44 Boeing manufacturing domain. *Applied Artificial Intelligence* **1994**, *8*, 125–147.  
45  
46  
47  
48 (36) Goodarzi, M.; Dejaegher, B.; Heyden, Y. V. Feature Selection Methods in QSAR Stud-  
49 ies. *J. AOAC Int.* **2012**, *95*, 636–651.  
50  
51  
52  
53 (37) Guyon, I.; Elisseeff, A.; De, A. M. An Introduction to Variable and Feature Selection.  
54 *JMLR* **2003**, *3*, 1157–1182.  
55  
56  
57  
58  
59  
60

- 1  
2  
3 (38) Mazurowski, M. A.; Habas, P. A.; Zurada, J. M.; Lo, J. Y.; Baker, J. A.; Tourassi, G. D.  
4 Training neural network classifiers for medical decision making: The effects of imbal-  
5 anced datasets on classification performance. *Neural Netw.* **2008**, *21*, 427–436.  
6  
7  
8  
9  
10 (39) Yang, P.; Hwa Yang, Y.; B. Zhou, B.; Y. Zomaya, A. A Review of Ensemble Methods  
11 in Bioinformatics. *Curr. Bioinf.* **2010**, *5*, 296–308.  
12  
13  
14 (40) Zhang, L.; Ai, H.; Chen, W.; Yin, Z.; Hu, H.; Zhu, J.; Zhao, J.; Zhao, Q.; Liu, H.  
15 CarcinoPred-EL: Novel models for predicting the carcinogenicity of chemicals using  
16 molecular fingerprints and ensemble learning methods. *Sci. Rep.* **2017**, *7*, 2118.  
17  
18  
19 (41) Güneş, F.; Wolfinger, R.; yi Tan, P. Stacked Ensemble Models for Improved Prediction  
20 Accuracy. SAS Global Forum Proc. 2017.  
21  
22  
23  
24  
25  
26 (42) Liu, J.; Patlewicz, G.; Williams, A. J.; Thomas, R. S.; Shah, I. Predicting Organ  
27 Toxicity Using in Vitro Bioactivity Data and Chemical Structure. *Chem. Res. Toxicol.*  
28 **2017**, *30*, 2046–2059.  
29  
30  
31  
32  
33 (43) Roy, K.; Kar, S.; Das, R. N. *Understanding the basics of QSAR for applications in*  
34 *pharmaceutical sciences and risk assessment*; 2015.  
35  
36  
37  
38 (44) Sheridan, R. P. The Relative Importance of Domain Applicability Metrics for Estim-  
39 ating Prediction Errors in QSAR Varies with Training Set Diversity. *J. Chem. Inf. Model.*  
40 **2015**, *55*, 1098–1107.  
41  
42  
43  
44 (45) Madasamy, K.; Ramaswami, M. Data Imbalance and Classifiers: Impact and Solutions  
45 from a Big Data Perspective. *IJCIR* **2017**, *13*, 2267–2281.  
46  
47  
48  
49 (46) Snoek, J.; Larochelle, H.; Adams, R. P. In *Advances in Neural Information Processing*  
50 *Systems 25*; Pereira, F., Burges, C. J. C., Bottou, L., Weinberger, K. Q., Eds.; 2012;  
51 pp 2951–2959.  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 (47) Bergstra, J.; Bengio, Y. Random search for hyper-parameter optimization. *JMLR* **2012**,  
4 *13*, 281–305.  
5  
6  
7  
8 (48) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**,  
9 *50*, 742–754.  
10  
11  
12  
13 (49) Zang, Q.; Mansouri, K.; Williams, A. J.; Judson, R. S.; Allen, D. G.; Casey, W. M.;  
14 Kleinstreuer, N. C. In Silico Prediction of Physicochemical Properties of Environmental  
15 Chemicals Using Molecular Fingerprints and Machine Learning. *J. Chem. Inf. Model.*  
16 **2017**, *57*, 36–49.  
17  
18  
19  
20  
21 (50) He, H.; Garcia, E. Learning from Imbalanced Data. *IEEE TKDE* **2009**, *21*, 1263–1284.  
22  
23  
24 (51) Cortes-Ciriano, I.; Bender, A. Improved Chemical Structure–Activity Modeling  
25 Through Data Augmentation. *J. Chem. Inf. Model.* **2015**, *55*, 2682–2692.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Graphical TOC Entry

