



# Neural-based approaches to overcome feature selection and applicability domain in drug-related property prediction



María Virginia Sabando\*, Ignacio Ponzoni, Axel J. Soto

Institute for Computer Science and Engineering, UNS – CONICET, Argentina

Department of Computer Science and Engineering, Universidad Nacional del Sur, Argentina

## ARTICLE INFO

### Article history:

Received 31 January 2019

Received in revised form 6 August 2019

Accepted 13 September 2019

Available online 19 September 2019

### Keywords:

Neural networks

QSAR modeling

Model interpretability

Applicability domain

Feature selection

## ABSTRACT

In the fields of pharmaceutical research and biomedical sciences, QSAR modeling is an established approach during drug discovery for prediction of biological activity of drug candidates. Yet, QSAR modeling poses a series of open challenges. First, chemical compounds are represented on a high-dimensional space and thus feature selection is typically applied, although this task entails a challenging combinatorial problem with potential loss of information. Second, the definition of the applicability domain of a QSAR model is a desirable aspect to determine the reliability of predictions on unseen chemicals, which is often difficult to assess due to the extent of the chemical space. Finally, interpretability of these models is also a critical issue for drug designers. The purpose of this work is to thoroughly assess the application of neural-based methods and recent advances deep learning for QSAR modeling. We hypothesize that neural-based methods can overcome the need to perform a descriptor selection phase. We developed three QSAR models based on neural networks for prediction of relevant chemical and biomedical properties that, in the absence of any feature selection step, can outperform the state-of-the-art models for such properties. We also implemented an embedded applicability domain technique based on network output probabilities that proved to be effective; its application improved the predictive performance of the model. Finally, we proposed the use of a *post hoc* feature analysis technique based on an aggregation of network weights, which enabled effective detection of relevant features in the model.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

The integration of computational sciences into the pharmaceutical and biomedical industry has yielded several applications and technological advances during the last decades [1,2]. The pharmaceutical industry is primarily headed towards improving the long and costly process of discovery and development of new drugs, which involves several stages including *in vitro* and *in vivo* wet-lab experiments. Computer-aided rational drug design has allowed accelerating drug candidate identification and prioritization while reducing costs and has helped improve the critical attrition rate in drug discovery projects [3,4]. The purpose of *in silico* drug discovery is to design models for predicting biological activity and physicochemical properties of drug candidate compounds. These models, referred to as Quantitative Structure-Activity Relationship (QSAR), are regression or classification models used in chemical and biological sciences to predict the relationship between features encoding the molecular

structure of compounds and the target property or biological activity under study. QSAR models are extensively used for virtual screening and prediction of categorical properties of drug candidates [5].

The development of QSAR models generally involves dealing with high-dimensional data representations. Drug candidates can be described by a large number of features or descriptors, which encode structural properties of the molecules. Feature selection is usually applied prior to the development of a QSAR model in order to harness high dimensionality [6,7], and mostly because traditional machine learning techniques do not perform properly in this high-dimensional scenario [8]. However, considering the large variety of possible descriptors that can be calculated from compounds, feature selection represents a difficult combinatorial problem that may neglect valuable information.

Another important aspect of QSAR modeling is determining the reliability of predictions on unseen compounds. The Applicability Domain (AD) of a QSAR model is the molecular subspace where predictions performed by the model are expected to be accurate. AD analysis is a significant step in the process of building a reliable QSAR model [9] and the identification of the AD of a QSAR model remains a current matter of research [10,11].

\* Corresponding author at: Department of Computer Science and Engineering, Universidad Nacional del Sur, Argentina.

E-mail address: [virginia.sabando@cs.uns.edu.ar](mailto:virginia.sabando@cs.uns.edu.ar) (M.V. Sabando).

A last important aspect of QSAR modeling is interpretability, as such models are managed by medicinal chemists in the process of searching for drug candidates. Being able to gain insight into the features that are the most relevant for prediction is valuable for experts [12], as such interpretation makes it possible to understand the molecular substructures that play a significant role in the biological activity or property of the chemical compound.

Regarding the techniques that have been used for QSAR modeling, two of them stand out. On the one hand, the use of meta-classifiers and consensus approaches has been widespread in QSAR modeling, and these methods have become the state of the art for predicting several physicochemical properties and biological activities [13,14]. On the other hand, artificial neural networks – a bio-inspired technique [15] – have also been used for QSAR modeling, although their adoption has been criticized due to their lack of generalization and the difficulty in the interpretation of such models in physicochemical terms [16]. Moreover, in recent years QSAR modeling has witnessed the advent of deep learning, which has brought several advantages as well as challenges. Recent advances in Deep Neural Network (DNN) approaches have made neural models less prone to overfitting, and hence more likely to be applied successfully for predicting unseen compounds. In addition, DNN-based models have been found effective for solving large-scale and high-dimensional data analysis problems [17].

The goals of this work are to build QSAR models that incorporate recent advances in deep neural networks for prediction of three relevant properties in biomedical sciences, and to benchmark these models against the state of the art. In addition, we aim to explore the potential of applying confidence estimation to neural-based models as an effective way for AD assessment, as well as to study the possibility to interpret these models in terms of the molecular features used to represent the chemical data. In order to address these goals, we propose the development of neural-based QSAR models for bioactivity prediction of three different properties. On these models, we evaluate the network output probabilities as a means of performing an AD estimation. In addition, we provide a post-hoc analysis of the most relevant features for each property. The first two properties are Cytochrome P450-drug interaction for isoforms 2C9 and 3A4, which are a family of enzymes involved in the oxidation of compounds. These two isoforms are particularly relevant to drug metabolism. It has been proven that inhibition of CYP enzymes leads to adverse side effects of drug–drug interactions [18], and hence the study of CYP interactions has become of major interest in the fields of drug discovery. The third property is Ready Biodegradability (RB). Biodegradation is highly relevant to biomedical sciences, since the presence of certain substances persisting over an extended period of time has been linked to major health risks, such as cancers, neurological dysfunction and hormonal changes [19,20]. Biodegradation properties are also relevant in the design of polymeric materials used for biomedical purposes [21]. Therefore, predicting biodegradability properties on chemicals represents a critical aspect for several biomedical areas.

The contributions of this work can be summarized as follows:

- We applied recent advances in deep neural networks to the development of neural-based QSAR models and obtained higher performance compared to state-of-the-art models, while at the same time overcoming the need for a potentially detrimental feature selection phase.
- We proved the effectiveness of using network output probabilities to perform AD estimation, which represents an advancement over consensus-based AD models that merely provide a binary signal with regard to inclusion or exclusion in the applicability domain.
- We applied a *post hoc* interpretability method based on an analysis of the network weights that has never been applied before in the context of QSAR models. We presented the results by means of a novel visualization based on heat maps, and proved that it effectively allows to gain insight into the interpretability of the proposed neural-based models.
- Our models outperformed the current state of the art for three different properties of high relevance in biomedical sciences.

This paper is organized as follows: in Section 2 we conduct a survey on relevant articles in the area and how they relate to our work. In Section 3, the datasets used for our experiments as well as the proposed methods are detailed. We present the results obtained for our proposed models and discuss their implications in Sections 4 and 5. Finally, conclusions and future lines of work are presented in Section 6.

## 2. Related work

For the past decades the process of rational drug design has relied on computer modeling techniques, and various *in silico* methods have been widely applied with the aim of both speeding up the discovery process and reduce costs [22–24]. Traditional techniques, such as Support Vector Machines (SVM), Decision Trees, Naïve Bayes and k-Nearest Neighbors, have been extensively used for building QSAR models because of their relatively good performance and simplicity [25–27]. Recently, there has been a strong tendency to consensus-based approaches, which consist in assembling different base classifiers to combine their predictions and, as a consequence, increase the prediction capabilities of the model [28–32]. This type of models are typically among the top performing techniques for the prediction of several chemical properties in QSAR modeling, but at the expense of limited interpretability. Besides, they are constrained by their base models, which usually rely on a feature selection step in order to perform at their best [33,34], and they are normally not able to capture complex relationships between descriptors or rule out redundant information [28,35].

### 2.1. Prediction techniques for Cytochrome P450-drug interaction and ready biodegradability

Extensive research work has been carried out for predicting Cytochrome P450-drug interaction, and the majority of these works usually involve a feature selection process [36–38]. Jensen et al. [36] presented two Gaussian kernel weighted k-Nearest Neighbors models. It was the first work to incorporate the use of Extended Connectivity Fingerprints (ECFP) and Functional Class Fingerprints (FCFP) [39] as features for CYP2D6 and CYP3A4 inhibition prediction. Cheng et al. [37] developed consensus-based models for prediction of five different CYP isoforms, using SVM, C4.5 Decision Tree, k-Nearest Neighbors and Naïve Bayes as base classifiers, combined by a backpropagation artificial neural network. They also showed an AD estimation that improves prediction accuracy. More recently, Shah et al. [38] developed a joint QSAR model based on feed-forward multi-layer neural networks for prediction of drug metabolism of isoforms 3A4, 2C9 and 2D6 of Cytochrome P450. Fingerprints were used as input features, and the three biological activities were embedded in a multitask deep neural network. Nembri et al. [40] developed two consensus-based models for prediction of isoforms 3A4 and 2C9 inhibition, and also performed an AD analysis. Both of the reported models were constructed upon two different voting approaches and used variations of k-Nearest Neighbors and a classification tree as base classifiers. Each base classifier was constructed employing either ECFP or a small number of molecular

descriptors obtained from a two-phase feature selection process. The best performing model reported for isoforms 2C9 and 3A4 was one of the voting approaches (namely *Consensus 1*). Since the work by Nembri et al. [40] reports one of the best prediction performances of biological activity for Cytochrome P450 and also due to the provision of all the data necessary for reproducibility, it constitutes the reference research work on CYP inhibition that we use for comparison.

There are also several research studies for prediction of both Biodegradability and Ready Biodegradability of compounds [10, 41–45], where two main approaches stand out: consensus and neural-based models. Consensus models are predominant for the prediction of this property, where feature selection techniques are applied in most cases [41,42,44,46]. Involving neural-based techniques, Goh et al. [45] developed a multimodal architecture for biodegradability prediction combining a Convolutional Neural Network with a fully-connected multi-layer perceptron, and using both domain-specific hand-engineered features and learned representations from raw data. In Mansouri et al. [41], two consensus models were proposed for prediction of Ready Biodegradability of compounds over three different base classifiers: Partial least Squares Discriminant Analysis (PLS-DA), SVM and k-Nearest Neighbors. The proposed consensus models were based on two different voting approaches and an AD analysis was carried out on the developed QSAR models. The best voting approach *Consensus 2* is reported as the best predictive model, which represents the best classification performance compared to other published QSAR models on biodegradation, and thus we chose this work Mansouri et al. [41] as our reference method for comparison.

Deep learning has emerged in the last years as a widely used soft-computing technique for the development of QSAR models and other areas in drug discovery research, and it has established itself as the state-of-the-art prediction technique [3,47]. Although artificial neural networks have already been used for QSAR models in the past [24,48], there is a recent tendency to adopt new strategies for training neural-based models, such as the application of novel techniques for avoiding overfitting and vanishing/exploding gradients during training. Although the application of deep learning in QSAR modeling is still in its beginnings, several research studies have developed deep learning-based models for various drug discovery problems successfully [49–51]. In Ma et al. [49], models based on DNNs achieved higher prediction performance than Random Forest on a group of large and diverse QSAR datasets. Lenselink et al. [50] compared five different techniques over a ChEMBL bioactivity benchmark set and found that DNNs outperformed traditional methods. They also showed that an ensemble of DNNs with additional tuning further improves the performance obtained by more simple DNN-based models. Koutsoukas et al. [51] showed that DNN-based models statistically outperform models based on traditional methods, such as Naïve Bayes, k-Nearest Neighbors, SVM and Random Forest over diverse datasets.

## 2.2. Applicability domain and interpretability of prediction models

The determination of the applicability domain of a QSAR model is a crucial aspect of the modeling process, since it allows to determine the molecular subspace of compounds where the QSAR model is expected to make reliable predictions [10,52]. Most research articles in the area address AD determination using a standalone method, where different strategies and statistical measures are adopted to determine AD boundaries [10]. A good number of them focus on defining different molecular similarity criteria for identifying outliers, which are then excluded from the AD of the model [53–55]. Klingspohn et al. [10] performed a comprehensive study in order to define a taxonomy of AD methods

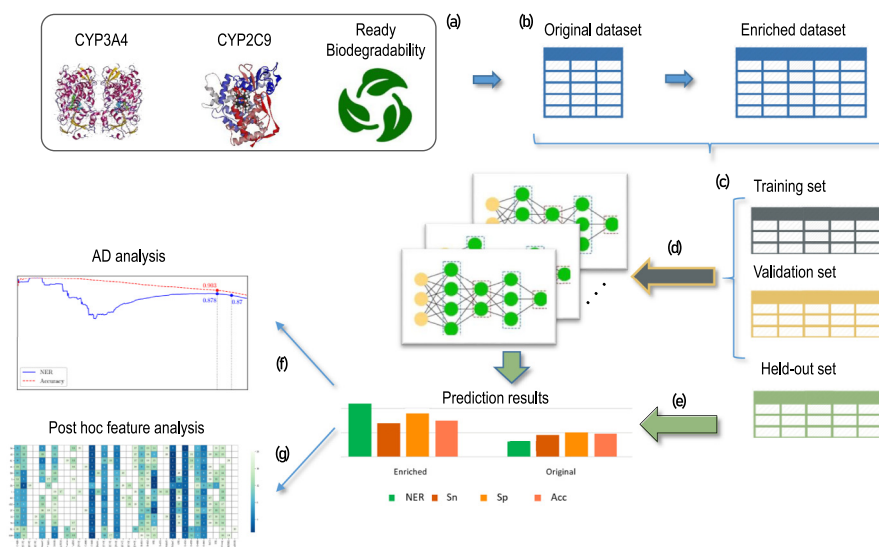
and find the best approaches for estimating the AD of different classification methods. In this article, two main categories of techniques for determining the AD were identified and compared: those based on novelty detection (identification of outliers) and those based on confidence estimation (inferred from the trained classifier). Experiments using six different binary classification techniques on ten datasets were performed. It was concluded that AD measures based on confidence estimation consistently perform better than novelty detection techniques, and thus they are suitable approaches for defining the AD.

Since QSAR models are meant to assist experts during drug discovery, their results should be as interpretable as possible [56]. Consensus-based approaches, in spite of having good predictive performance, tend to lack interpretability since their output result is a combination of different base classifiers. Interpretability of neural network-based models has been studied for several years within the machine learning community [57–59] and it also remains a matter of research in drug discovery. Approaches based on *post hoc* interpretability have been explored recently [60–62]. This type of techniques takes a trained model and makes an attempt to understand predictions in terms of the features used by the model. As opposed to a low-level algorithmic comprehension of the model, which is the most usual approach taken for interpretability analysis, *post hoc* techniques aim to characterize the behavior of the predictive model without attempting to explain its internal representation and operations, but providing a functional understanding of it in terms of its features. In this line of work, Tsang et al. [63] developed a framework to discover statistical interactions between the input features in a feed-forward multi-layer neural network, by direct interpretation of its learned weights. Their method proved to be effective on both synthetic and real-world application datasets, and thus we chose it as our reference paper for *post hoc* feature analysis.

## 2.3. Our proposal

Based on the observed limitations and the current state of the art, we propose the development of neural-based methods to model three physicochemical properties, namely: Cytochrome P450-drug interaction for isoforms CYP2C9 and CYP3A4, and Ready Biodegradability. We compare our approach with the consensus models that, to the best of our knowledge, represent the top-performing models that have been published for the prediction of these properties. We present an embedded applicability domain technique, which is derived from our trained models. This approach would be categorized as prediction confidence estimation according to Klingspohn's taxonomy [10]. We propose the application of a *post hoc* interpretability technique based on an aggregative analysis of the weight contributions of the network, which is based on Tsang et al. [63]. To the best of our knowledge, this method has not yet been employed for interpretability of QSAR models. Additionally, the results of this *post hoc* analysis are summarized using a novel visualization based on heat maps. Finally, our models and results are contrasted to those reported in Nembri et al. [40] and Mansouri et al. [41]. The reason for choosing these latter articles as our baselines for comparison is due to the possibility of reproducing their data and experiments and the high performance attained in the reported results.

The entire workflow of our approach is depicted in Fig. 1. First, we preprocessed the three datasets under study (a). Second, we enriched the Original datasets by adding new molecular descriptors (b) and we split the Enriched and Original datasets into the partitions for training and validating our models (c). Then, we developed our neural-based models by performing an iterative process for hyperparameter tuning and we train the chosen models (d). Next, we evaluated their performance using several



**Fig. 1.** Depiction of the entire workflow of our method: (a) data preprocessing, (b) dataset enrichment, (c) dataset partition, (d) development and training of models, (e) evaluation of models, (f) AD analysis, (g) *post hoc* feature analysis.

metrics and we contrasted these results with different baselines (e). Finally, we developed an AD model based on confidence estimation and applied (f) a *post hoc* feature analysis method, which allows to determine the most influential features to our models (g).

### 3. Materials and methods

In this section, we provide an overall description of all the techniques used in our work, as well as the preparation of the datasets and the selection of the hyperparameters of the model.

We say that we *enrich* a dataset when we extend its set of features by including new molecular descriptors not previously considered in the Original dataset (Fig. 1-b). It is worth noting that the included descriptors are members of the family of descriptors already included in the Original datasets. The reported partitions of the datasets in Train, Validation and Held-out sets<sup>1</sup> were kept the same during the construction of our models (Fig. 1-c). For the sake of completeness, we also trained our *Best\_E* models using 5-fold Cross Validation. The details of this process and its results are summarized in the Supplementary Material.

All of our models are based on feed-forward multi-layer neural networks. The architecture for each model varies depending on the number of input features or molecular descriptors. As a general approach, larger inputs demand more complex architectures, so the Enriched versions of the datasets yielded models with more nodes than those built for the Original versions of the datasets. Our neural-based models were obtained following a two-phase process (Fig. 1-d). The first one consisted in an exploratory phase, where different architectures and optimization strategies were considered. In this phase we developed prototypes and tuned their parameters. The need for an exploratory phase when constructing neural networks has been reported previously in the context of QSAR modeling [64]. The second phase consisted in selecting the best prototype from the first phase. This selection was performed by assessing classification performance on the Validation set. Due to the inherent stochasticity of neural networks, we repeated the training process using a set of fifteen random seeds for each dataset. As a result of this two-phase

process we obtained fifteen models with the same hyperparameters but initialized differently. The average performance of these fifteen models is reported as the *Average* model, whereas we report as the *Best* model the one that performs the best on the Validation set.

After the best model for each Enriched dataset was obtained, we built a new model for the Original version of each dataset – namely *Best\_O* – using the same hyperparameters as for the Enriched models but using less nodes per layer. The construction of two models, one based on the Enriched dataset and another based on the Original dataset, enable us to compare the performance of our proposed strategy over one same set of compounds with and without the application of a feature selection process, and it allows to analyze the potential of our approach on high-dimensional datasets.

#### 3.1. Datasets

The three datasets used in this work, namely CYP2C9, CYP3A4 [40] and RB [41], are publicly available and were selected taking into account their relevance in QSAR modeling in the context of biomedical data analysis. We made datasets CYP2C9, CYP3A4 and RB in their Enriched versions publicly available.<sup>2</sup> Further information on the calculation of molecular descriptors for the construction of the Enriched datasets can be found in the Supplementary Material.

##### 3.1.1. CYP2C9 and CYP3A4

Datasets CYP2C9 and CYP3A4 have a total of 11 940 and 12 118 compounds, and the proportion between active/inactive compounds on their Training and Validation sets is 49/100 for CYP2C9 and 66/100 for CYP3A4. As for the Held-out sets, the ratios are 56/100 in CYP2C9 and 98/100 in CYP3A4. CYP2C9 and CYP3A4 share the same compounds in their Training sets as well as in their Validation sets. In the Original datasets, CYP2C9 includes ten molecular descriptors, whereas CYP3A4 includes eight molecular descriptors. They also include a 1024-bit ECFP for each compound [39]. For the Enriched versions of both datasets we added a total of 2701 molecular descriptors to the CYP2C9 dataset and 2699 molecular descriptors to the dataset CYP3A4, leading to a

<sup>1</sup> Note that the Original datasets refer to these partitions in their papers as *Training*, *Test* and *External Validation*, respectively.

<sup>2</sup> <https://github.com/VirginiaSabando/DNN-QSAR-2019.git>.



total of 2711 and 2707 descriptors, respectively, in addition to the 1024-bit ECFP. We performed the calculation of molecular descriptors and ECFP using Dragon 7 [65].

There were a few compounds on the datasets provided by Nembri et al. [40] whose SMILES codes were not properly formed, and hence we were unable to calculate their molecular descriptors. As a result, one molecule was removed from both Training sets, six molecules were removed from the Held-out set of CYP2C9, and one molecule was removed from the Held-out set of CYP3A4.

### 3.1.2. Ready biodegradability (RB)

Dataset RB comprises 1725 compounds, where the ratio between active/inactive compounds is 51/100 for the Training set, 49/100 for the Validation set and 40/100 for the Held-out set. A total of 41 molecular descriptors were provided for dataset RB in its Original form. We calculated additional molecular descriptors to obtain its Enriched version, which gave a total of 1480 molecular descriptors. We computed these descriptors using Dragon 7 [65].

### 3.2. Model parameterization

In this section, we describe all the model parameters used for the three Enriched datasets. Parameters for the remaining models and their training can be found in the Supplementary Material. The models were built using Tensorflow 1.7 [66].

The input features for our predictive models were molecular descriptors and ECFP. In the cases of CYP2C9 and CYP3A4, we split the ECFP into 1024 separate bits and then considered each one of these bits as a single input feature. All nodes in the input and hidden layers of the models use rectified linear units (ReLU) as activation function [67], and the output layer implements a softmax function, which has two nodes for each class output probability. The networks were trained using standard backpropagation [68] and we chose a minibatch size of 200 instances to train the networks. We use cross-entropy with logits as cost function and Adam optimizer [69] for minimizing it.

For weight initialization we experimented with Xavier initialization [70] and He Normal initialization [71], which are both considered state-of-the-art initialization techniques [72,73]. We also applied Batch Normalization [74] in all layers of our networks for faster training and to avoid exploding/vanishing gradients. In order to avoid overfitting, several regularization techniques were used in each model. We used Dropout [75] with varying dropout rates according to the number of nodes in each layer, and also implemented L2-regularization [76] varying the penalization coefficient  $\lambda$  in each model. We applied early stopping to avoid overtraining the models, hence helping to prevent possible overfitting.

The QSAR model obtained for Enriched dataset CYP2C9 is a feed-forward multi-layer neural network architecture consisting of one input layer of 3735 nodes (for 2711 molecular descriptors plus 1024 bits from ECFP) and three hidden layers of 50, 20 and 5 nodes, respectively. Batch Normalization was applied with a decay value of 0.9 to prevent the weights from growing too large, and we used a learning rate value of 0.00001. We initialized the network weights by using Xavier initialization. A penalization coefficient  $\lambda = 0.0001$  was used for L2-regularization. In order to deal with class imbalance we optimized a weighted cost function, which penalized mispredicted instances from the least popular class by increasing the loss by a factor proportional to the class imbalance observed in the Training set.

The architecture of the QSAR model that we developed for Enriched dataset CYP3A4 is similar to that used for CYP2C9, with the difference that the input layer consists of 3731 nodes (for 2707

molecular descriptors plus 1024 bits from ECFP). As in the case of dataset CYP2C9, all layers implement Batch Normalization with a decay value of 0.9. For CYP3A4, we initialized network weights by applying He Normal initialization, and the same regularization criteria than for CYP2C9 was taken into account for this dataset. For class imbalance mitigation we applied a stratified sampling technique, where an equal number of compounds belonging to each class was drawn to build each of the minibatches during training. The compounds were sampled with replacement from the training set and randomly shuffled before they were fed to the network during the training phase.

Lastly, we developed a QSAR model for Enriched dataset RB based on a less dense feed-forward multi-layer neural network architecture, considering that the input features were fewer than those in the previously described models. The input layer comprises 1480 nodes for molecular descriptors, and the network is also made of three hidden layers of 20, 10 and 5 nodes, respectively. All layers implement Batch Normalization with a decay value of 0.9, as in the case of the previous models. We initialized the network weights by using Xavier initialization, and as regularization techniques we used Dropout and L2-regularization with a penalization coefficient  $\lambda = 0.001$ . The learning rate was set to 0.0001. We used a stratified sampling technique in order to counteract class imbalance, with the same sampling technique as described in the case of CYP3A4.

### 3.3. Applicability domain

The applicability domain (AD) of a QSAR model is the molecular subspace in which the predictions made by the model are expected to be accurate [77,78]. In other words, the definition of an AD allows the expert to determine whether a prediction on a new compound is likely to be reliable or not.

We propose using class probability provided by the output layer of our models to estimate their AD. This leads to AD models which are embedded into the prediction models. The embedded AD models were evaluated as follows. First, we computed class probability values using the network softmax layer for every compound. Then, we sorted these values in decreasing order to elaborate a ranking. Finally, in order to evaluate the goodness of the ranking of confident predictions, we computed Mean Average Precision (MAP) [79], where several metrics (Accuracy, *NER*, etc.) were calculated on the *k*-highest ranked compounds, where *k* is varied from 1 to *n*, and *n* being the total number of compounds. All these different metrics computed for different number of compounds are averaged. This AD approach (Fig. 1-f) allows to evaluate the performance of the models at any desired threshold of membership to the AD.

### 3.4. Post hoc feature analysis

As QSAR models are tools for the benefit of chemists and drug developers alike, it remains an important asset to provide means of interpretability for any proposed models [56,80]. For domain specialists it is useful to know the features that make a particular family of compounds to show some degree of activity regarding a property of interest, since this allows to reduce the search space during drug discovery.

We propose a *post hoc* feature analysis technique as a way of providing interpretability to our neural-based models, so that domain experts can determine the most relevant molecular descriptors in the context of a prediction model (Fig. 1-g). By analyzing the network weight contributions in an aggregative manner, it is possible to gain insight into the descriptors that are more influential on the predicted target value.

The proposed *post hoc* feature analysis technique is described as follows: once training is completed, we calculate a score for every descriptor by taking into account the sequence of contributions from the input to the output nodes. For a given feature, these contributions are calculated by aggregating the weights of the neural model that are connected to this feature. More formally, for any layer, the score of a node  $j$  is computed using

$$S(n_j) = \frac{1}{k} \sum_{i=1}^k |w_{j,i}| S(n_i), \quad (1)$$

which is the average of the  $k$  products between the weights connecting node  $j$  with the  $k$  nodes in the following layer and their scores. Given that this is a recursive definition, by setting the score corresponding to the node of the output layer to 1, we can compute all node scores by starting from the output nodes and going backwards until the scores corresponding to the input nodes are computed. We considered the absolute values of the weights, as a way to analyze quantitative impact of the descriptors on the output, regardless of whether they contribute in a positive or negative manner on the result. The rationale is that input features exhibiting high scores are likely to be more relevant than those showing low scores, as slight changes in their values would have greater impact on the outcome of the network.

#### 4. Results

We performed an evaluation of each of our models by comparing them against *Consensus 1* and *Consensus 2*, which are the top-performing methods ever reported for the three datasets under study [40,41] (Fig. 1-e). To account for a fair comparison, we used the same metrics as reported in Nembri et al. [40], Mansouri et al. [41], i.e., Sensitivity ( $Sn$ ), Specificity ( $Sp$ ) and NER, as well as Accuracy ( $Acc$ ) and MAP, as described in Section 3.3. Sensitivity and Specificity quantify the accuracy in predicting the active and inactive class, respectively, while NER is the arithmetic mean of  $Sn$  and  $Sp$ . Additionally, other performance metrics can be found in the Supplementary Material.

Since our neural-based models are inherently stochastic, fifteen different trials were run to train and test the models, each one using a different random seed. Therefore, we report both the average performance of the developed models, i.e., taking into account all trials, and the best performance in terms of  $NER$ , i.e., the model obtained from the best seed.

Regarding the AD analysis, we report both  $NER$  and the percentage of compounds that are not within the AD – which are referred to as *Not Assigned* compounds ( $\%na$ ) – as it was also done by Nembri et al. [40] and Mansouri et al. [41]. We report  $NER$  when  $\%na$  is fixed to the value of the best reported method in the referenced articles. Similarly, we report  $\%na$  when  $NER$  is fixed to the same values reported in the referenced papers.

##### 4.1. CYP2C9

The results for the Original and Enriched versions of CYP2C9 are presented in Table 1. We also include the results of the best model reported by Nembri et al. [40], i.e., *Consensus 1*. It can be seen that for both Validation and Held-out sets our best Enriched model, namely *Best\_E*, performs better than the best model on the Original dataset, namely *Best\_O*, and both of them outperform *Consensus 1*. In addition, all of our models achieve equivalent or better  $NER$  values than *Consensus 1* when keeping  $\%na$  at the same value.

A more comprehensive evaluation of the performance of our models on the CYP2C9 Held-out set is presented in Figs. 2 and 3.

**Table 1**

Results on the Validation and Held-out sets of CYP2C9. Consistently with the results reported by Nembri et al. [40] the percentage of not assigned compounds ( $\%na$ ) was set to 40% for Validation set, and 45% for Held-out set.

CYP2C9		Validation set				Held-out set			
		NER	Sn	Sp	Acc	NER	Sn	Sp	Acc
Original	Consensus 1	0.89	0.89	0.88	–	0.83	0.85	0.82	–
	Average_O	0.91	0.90	0.92	0.91	0.83	0.79	0.88	0.86
	Best_O	0.92	0.92	0.92	0.92	0.85	0.82	0.87	0.86
Enriched	Average_E	0.92	0.93	0.91	0.92	0.85	0.87	0.83	0.84
	Best_E	<b>0.93</b>	<b>0.94</b>	<b>0.93</b>	<b>0.93</b>	<b>0.87</b>	<b>0.89</b>	<b>0.86</b>	<b>0.87</b>

**Table 2**

Results on the Validation and Held-out sets of CYP3A4. Consistently with the results reported by Nembri et al. [40] the percentage of not assigned compounds ( $\%na$ ) was set to 36% for Validation set, and 42% for Held-out set.

CYP3A4		Validation set				Held-out set			
		NER	Sn	Sp	Acc	NER	Sn	Sp	Acc
Original	Consensus 1	0.88	0.92	0.83	–	0.80	0.89	0.70	–
	Average_O	0.89	0.83	0.94	0.91	0.82	0.76	0.88	0.83
	Best_O	0.91	0.91	0.91	0.91	0.83	0.84	0.82	0.83
Enriched	Average_E	0.92	0.89	0.94	0.92	0.84	0.84	0.85	0.84
	Best_E	<b>0.93</b>	<b>0.91</b>	<b>0.94</b>	<b>0.93</b>	<b>0.85</b>	<b>0.86</b>	<b>0.84</b>	<b>0.85</b>

These figures show different performance measures when different cutoff values for the AD are considered. In Fig. 2 we present the mean of all trials—i.e., *Average* performance, whereas in Fig. 3 the results for the best trial are presented. Both figures correspond to the models trained on the Enriched version of CYP2C9. The horizontal axis represents the number of compounds sorted by class probability, so that the left-most compounds are the most confidently predicted ones. The vertical axis represents different performance measures evaluated over the set.

By looking at these plots it is possible to set any cutoff point in the horizontal axis in order to evaluate performance when the compounds with least certain prediction – those to the right of the cutoff point – are discarded. In particular, two cutoff points are marked as noteworthy; these are the cutoff values where  $\%na$  and  $NER$  match with the ones reported by Nembri et al. [40]. MAP results for the Validation sets of the three datasets can be found in the Supplementary Material.

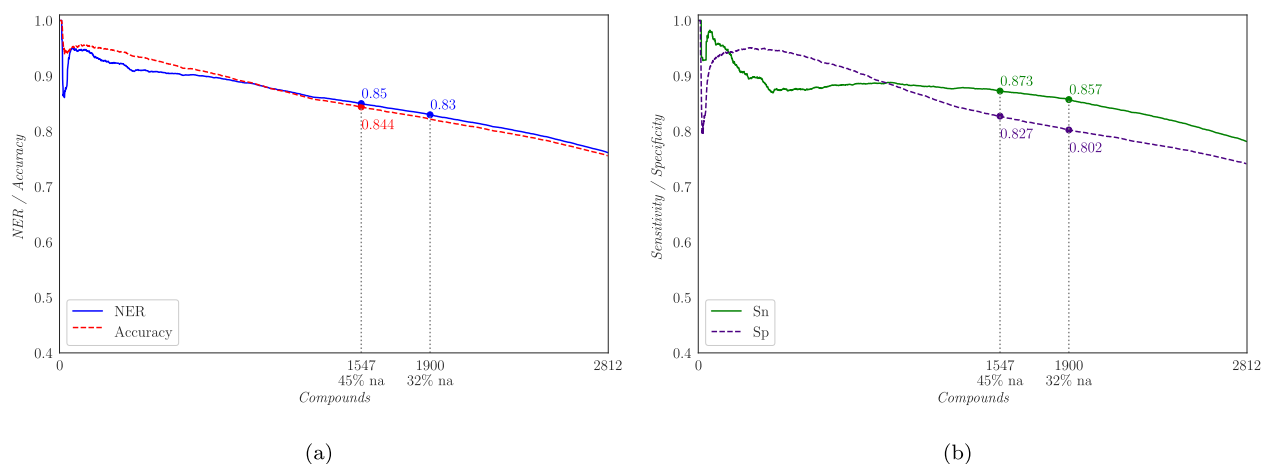
##### 4.2. CYP3A4

We present the results for CYP3A4 in Table 2. For both Validation and Held-out sets our best Enriched model show better performance than the best model trained on the Original version of the dataset, which in turn overcomes the results reported for *Consensus 1*. All of our models obtain higher Non-Error Rate for the same number of discarded compounds than the reference model, yet ours exhibiting balanced values of Sensitivity and Specificity.

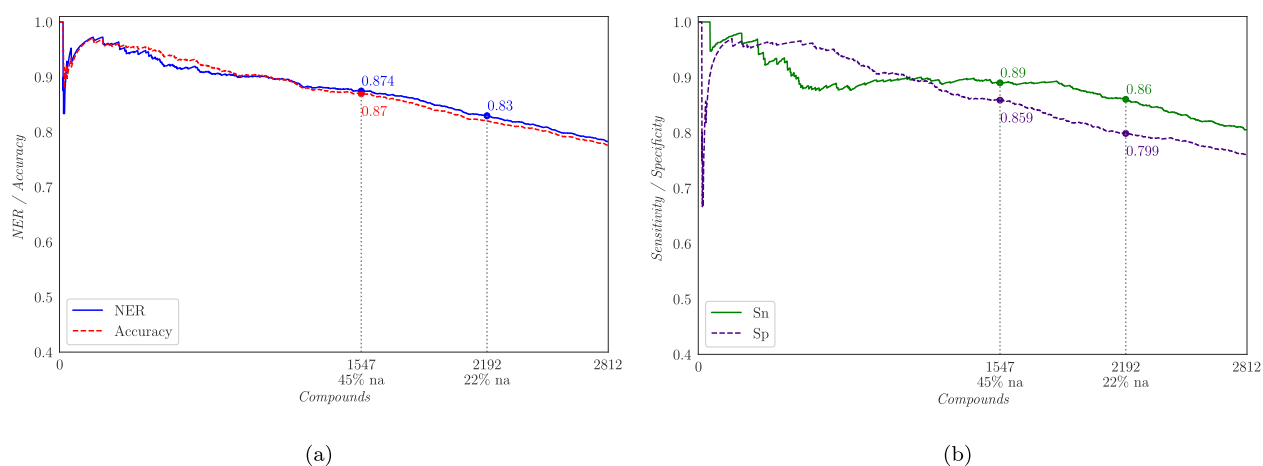
The plots showing the comprehensive performance of the QSAR models with regard to its AD model for CYP3A4 Held-out set are presented in Figs. 4 and 5. Similarly as it was done for CYP2C9, Fig. 4 shows the results for the mean of all of our trials, while Fig. 5 reports the results when our best model is considered.

##### 4.3. Ready biodegradability

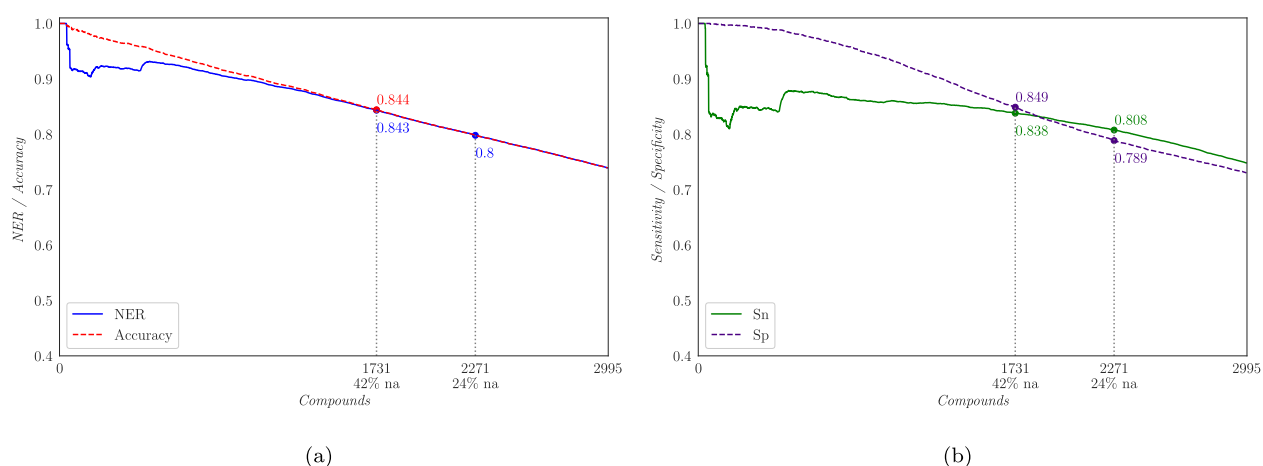
Table 3 shows the results for RB. Our best Enriched model, i.e., *Best\_E*, exhibits higher  $NER$  in both Validation and Held-out sets than *Best\_O*, our best model trained on the Original dataset. Both of these models show higher  $NER$  values than for *Consensus 2*.



**Fig. 2.** Average MAP performance of all trials on the Held-out set of Enriched CYP2C9. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.



**Fig. 3.** MAP performance of the best trial on the Held-out set of Enriched CYP2C9. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.



**Fig. 4.** Average MAP performance of all trials on the Held-out set of Enriched CYP3A4. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.

The plots displaying the performance of the QSAR models with regard to its AD model for the Enriched model on the RB Held-out set are presented for the mean of all of our trials (Fig. 6), and for the best trial, i.e., *Best\_E* (Fig. 7).

#### 4.4. Post hoc feature analysis

We performed a *post hoc* feature analysis in order to gain insight on which molecular descriptors are the most relevant to our

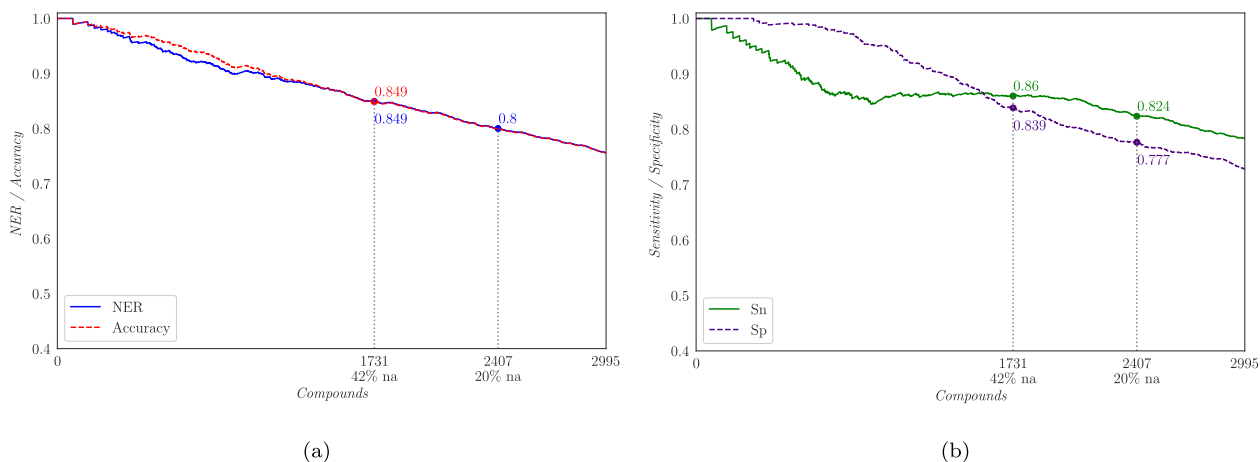


Fig. 5. MAP performance of the best trial on the Held-out set of Enriched CYP3A4. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.

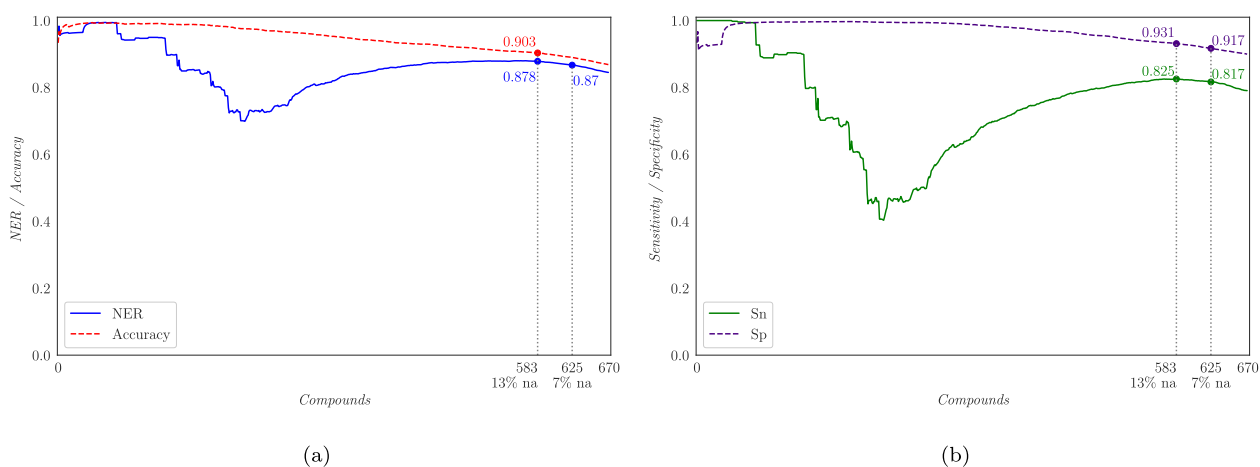


Fig. 6. Average MAP performance of all trials on the Held-out set of Enriched RB. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.

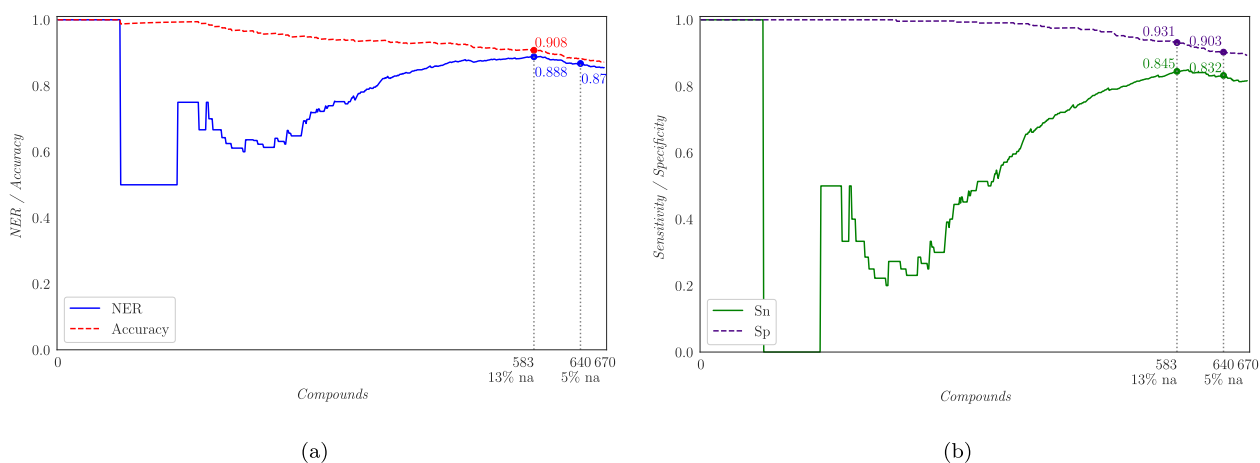


Fig. 7. MAP performance of the best trial on the Held-out set of Enriched RB. (a) NER and Accuracy are shown. (b) Sensitivity and Specificity are shown.

models. We propose a novel visualization for summarizing main patterns of features that were found to be relevant across multiple trials by means of a heat map. The heat maps that encode the results corresponding to the feature analysis for datasets CYP2C9, CYP3A4 and RB are presented in Figs. 8–10, respectively. Each row corresponds to a different trial of our model using its own seed for initialization of weights and random variables. These trials are sorted by performance, where the top row represents

the best trial. Each column on the maps represents molecular descriptors, where only the 20 most relevant descriptors of each model according to our measure were considered. Descriptors marked with an asterisk are also part of the Original version of the dataset. Descriptor names starting with 'ECFP' represent Extended Connectivity Fingerprint fragments, which are followed by a number that represents the location of the bit that was identified by our method as relevant. The rank that a descriptor



**Table 3**

Results on the Validation and Held-out sets of Ready Biodegradability (RB). Consistently with the results reported by Mansouri et al. [41] the percentage of not assigned compounds (%na) was set to 15% for Validation set, and 13% for Held-out set.

Ready Biodegradability (RB)		Validation set				Held-out set			
		NER	Sn	Sp	Acc	NER	Sn	Sp	Acc
Original	Consensus 2	0.91	0.88	0.94	–	0.87	0.81	0.94	–
	Average_O	0.92	0.94	0.90	0.91	0.88	0.85	0.91	0.90
	Best_O	0.91	0.91	0.91	0.91	0.88	0.85	0.92	0.90
Enriched	Average_E	0.94	0.93	0.90	0.94	0.88	0.83	0.93	0.90
	Best_E	<b>0.94</b>	<b>0.95</b>	<b>0.92</b>	<b>0.93</b>	<b>0.89</b>	<b>0.85</b>	<b>0.93</b>	<b>0.91</b>

occupies in the relevance order of a particular trial is encoded with the number inside the corresponding heat map cell. Likewise, darker colors are applied to cells depicting higher relevance descriptors in a specific trial, whereas lighter colors apply to less relevant descriptors.

## 5. Discussion

In this section, we review the results presented previously in Section 4. We discuss the performance of the models, as well as the results of our embedded AD model and *post hoc* feature analysis technique.

### 5.1. Neural-based classifiers versus consensus-based classifiers

As it is shown in Table 1, our models for prediction of CYP2C9 drug interaction outperform the results reported by Nembri et al. [40]. The *NER* values of *Average\_E* and *Average\_O* are consistently superior to the reference results. On the internal Validation set of CYP2C9, the average *Sn* and *Sp* values in both the Original and the Enriched version of the dataset are higher and more balanced than those achieved by *Consensus 1*, which indicate that our model has successfully overcome the class imbalance of the dataset as it was able to correctly predict both active and inactive compounds with similar accuracy. When evaluated on the Held-out set of CYP2C9, our models also improve the performance of the reference results, although the differences between our models and *Consensus 1* are smaller than those obtained on the Validation set. Cohesively, a mild imbalance between *Sn* and *Sp* is observed, which is consistent with the results reported by Nembri et al. [40] on the Held-out set. The best trial on the Enriched version of the CYP2C9 dataset, i.e., *Best\_E*, attained a *NER* value of 0.93 for the internal validation set and a value of 0.87 when tested on the Held-out set, which is an improvement of 0.04 over the same results for *Consensus 1*.

In Table 2 the predictive performance of our models exceed the results reported using *Consensus 1*. Similarly to what it was observed for CYP2C9, the average *NER* of our models is higher than that reported for *Consensus 1* in the Original and Enriched versions of the dataset, while also showing balanced results between *Sn* and *Sp* average values. On the Enriched Held-out set, while the predictive performance slightly decreases compared to the results on the Validation set, balanced results between *Sn* and *Sp* values are obtained. This observation does not hold for the average results in the Original version of the dataset. It is worth noticing that average results were computed by taking into consideration all trials of the model, including those undertrained due to the model apparently getting caught on local minima. In a similar way as it happened for CYP2C9, *Best\_E* obtained the highest *NER* for both the internal validation set and Held-out set – 0.93 and 0.85, respectively – which represents an increase of 0.05 over the same results for *Consensus 1*.

From Table 3 we can see that the predictive performance of our models improves the results using *Consensus 2* [41] for the RB dataset. The average *NER* of our models in the Validation set is higher than that of *Consensus 2* in both the Original and Enriched versions of the dataset, and at the same time showing more balanced *Sn* and *Sp* average values. The performance of the model on the Held-out set is higher in the Original dataset than in its Enriched version. Balance between *Sn* and *Sp* values is not observed for this partition, where *Sn* is consistently lower than *Sp* in all the experiments. It is noteworthy that this imbalance is also present in *Consensus 2*, which suggests an issue with the Held-out set data that makes prediction of inactive compounds to be inaccurate when compared to results for the Validation set. Nonetheless, the best trial on the Enriched version of RB, i.e., *Best\_E*, attained the highest *NER* in both Validation set and Held-out set—0.94 and 0.89, respectively.

For both datasets CYP3A4 and CYP2C9, a high consistency is observed between the results got for both Validation and Held-out sets. These results show that the obtained models have strong generalization capabilities. Besides, for both Held-out sets no large disparity is observed between *Sn* and *Sp* values of *Best\_E*, which in turn implies that the proposed models are able to classify active and inactive compounds unbiasedly. At a high-level analysis of the results, the results for the three datasets show that their Enriched versions lead to models with higher predictive performance and more balanced Specificity and Sensitivity than just using the Original versions with the descriptors chosen by means of a feature selection approach. Attaining balanced performance in a binary classification problem is a desirable quality in a QSAR model [81]. The results of the Enriched models suggest that there is relevant data encoded on molecular descriptors that were not present in the Original versions of the datasets. Consequently, our experimental results imply that neural networks are able to learn in large dimensionality scenarios, and that performing a feature selection step could lead to valuable information loss, and hence to a decrease in the predictive performance. Furthermore, our work proves that neural-based QSAR models are capable of surpassing the benchmarked consensus-based models. Therefore, the use of neural networks constitutes a strong approach for QSAR modeling.

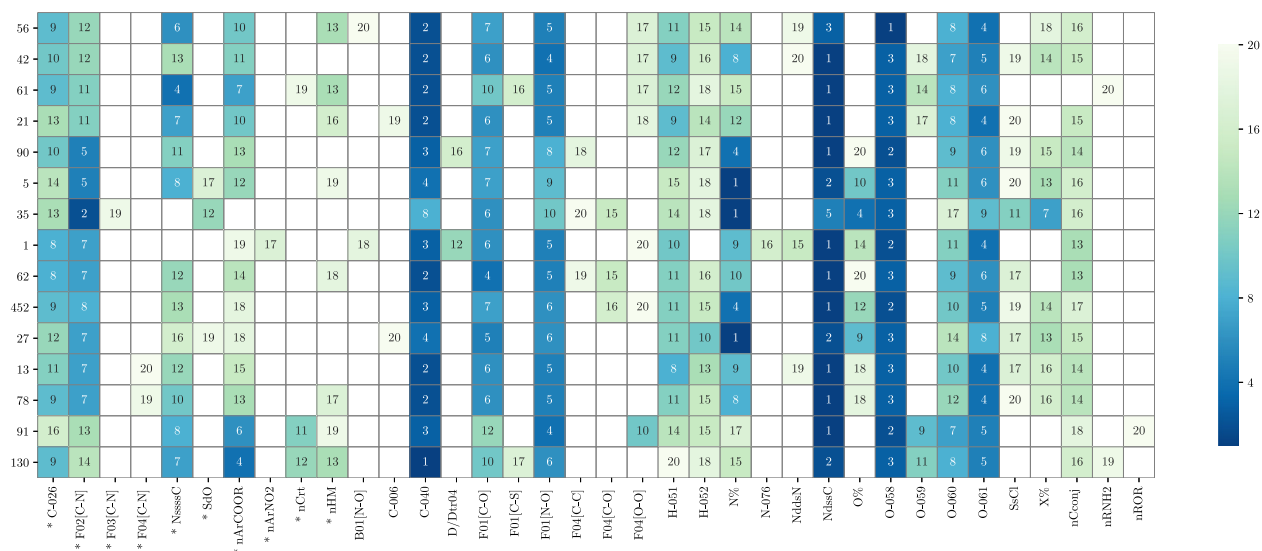
### 5.2. Embedded applicability domain technique

We proposed using an embedded AD model based on the class probabilities calculated by the output layer of our prediction model. This approach was applied on each QSAR model and evaluated on the Validation and Held-out sets by measuring the extent by which misclassification is correlated to the predicted class confidence.

The plots for CYP2C9, which are displayed in Figs. 2 and 3, show that as the number of compounds in the AD increases, i.e., to the right on the horizontal axis, the values for all metrics tend to decay continuously, which implies that the predictive performance of the model is in fact correlated with our definition of AD.

It is fair to say that these curves are not smooth for the left-most compounds. Some peaks are observed in *NER* and *Acc* plots (Figs. 2-a and 3-a), while slope fluctuations in the *Sn* and *Sp* curves are observed (Figs. 2-b and 3-b). For the best trial *Best\_E*, as it can be seen from Fig. 3, a strong downward peak is observed on the *Sp* curve for the left-most compounds. This was caused by a few inactive compounds that were misclassified with a high output probability, which is clearly an unexpected result. It is worth noticing, however, that the curve fluctuation stabilizes shortly after this peak is observed. Fig. 2 shows that in average all trials exhibited similar behavior to the best trial,





**Fig. 10.** Post hoc feature analysis on Ready biodegradability (RB). The rows represent different trials of our model sorted by performance, and the columns represent the 20 most relevant molecular descriptors. Both the color and the value in the cells represent the rank of a descriptor in terms of its relevance for a particular trial.

and CYP3A4 datasets. The difference between these sets of compounds would explain the generalization problems of *Best\_E* on RB Held-out set.

This generalization problem due to data distribution differences gets exacerbated as the dimensionality of the data increases, so the models built upon the Enriched version are disfavored in contrast to those models for the Original version.

From the figures discussed above it is clear that our models were able to reach the same *NER* values as reported by Nembri et al. [40] and Mansouri et al. [41] for the three datasets, yet dismissing many fewer compounds as not assigned (%*na*) than the models proposed therein. For CYP2C9, a *NER* value of 0.83 with 45%*na* is reported for *Consensus 1*, while *Best\_E* attains the same *NER* value dismissing only 22% of compounds on the Held-out set. In the case of CYP3A4, a *NER* value of 0.8 with 42%*na* is reported for *Consensus 1*, while *Best\_E* achieves the same *NER* value with only 20%*na* on the Held-out set. Finally, on the Held-out set of RB dataset, a *NER* value of 0.87 with 13%*na* is reported for *Consensus 2*, while *Best\_E* attains the same *NER* value with only 5%*na*. A similar analysis could be performed by taking into consideration %*na* reported by Nembri et al. [40] and Mansouri et al. [41] for the three datasets, since our models systematically reached higher *NER* values for the same amount of discarded compounds than *Consensus 1* and *Consensus 2* in both Validation and Held-out sets.

### 5.3. Post hoc feature analysis

From the heat maps in Figs. 8–10, one interesting aspect in all three models is that we can pinpoint molecular descriptors that were highly influential to all trials of the same model. For instance, for CYP2C9, the ECFP fragment *ECFP\_393* was the most relevant feature for eleven out of the fifteen trials, being nine of those trials among the best performing ones. Molecular descriptors *H-046* and *nRNR2* were also frequently selected in the different trials; the latter descriptor is also present in the Original CYP2C9 dataset. In the case of CYP3A4, the fragment *ECFP\_885* was a highly influential feature during the training phase of the model, as it was considered among the top three most relevant descriptors in all trials. Descriptors *SdssC* and *H-049* were also signaled as relevant to the majority of trials, according to our measure. Interestingly, the molecular descriptor *D/Dtr04* was identified as an important feature occupying the first place in five trials, although these trials were among the worst trials. For

dataset RB, the molecular descriptor *NdsC* is chosen as the most relevant for most models, as it was considered the most relevant descriptor for eleven out of the fifteen trials. Descriptors *O-058* and *C-040* were also signaled as important features, occupying the top three positions for the majority of trials.

Another interesting aspect that can be observed from these heat maps is that similarly performing trials tend to choose the same descriptors and in similar order of relevance. For instance, in Fig. 8 descriptors are: *N-071*, *NssssCH*, *C-034* and *H-051* were spotted as relevant only by the best performing trials and occupied similar positions on the relevance rankings of every trial. Similarly, *SsssN*, *T(N..N)*, *H-047*, *MLOGP2*, *F01[C-N]* and the *P\_VSA* family of descriptors were deemed as influential in low-performing trials. The ECFP fragments were mostly included by the best-performing trials. The same phenomenon is observed in Fig. 9: *ECFP\_509*, *ECFP\_599* and *ECFP\_862* were mostly signaled by our measure in the best trials, while descriptors *D/Dtr05*, *D/Dtr11*, *SAdon*, *SsOH*, *SsssN*, *nArOR* and the *P\_VSA* family of descriptors were found to be somewhat relevant in the low-performing trials.

Among the molecular descriptors identified as relevant for each model by our technique, there are some descriptors that are also in the Original versions of the datasets. The largest number of descriptors shared between these two sets is observed for dataset RB, where 10 out of 36 of the features signaled as the most important were also present in the Original RB dataset. Out of these 10, only 4 descriptors were highly relevant to the majority of trials, yet occupying medium-to-low importance positions in all models. Fig. 8 shows that only 3 out of the 36 descriptors are present in both the Original and Enriched CYP2C9 dataset; however, as mentioned before, the descriptor *nRNR2* was identified as one of the most relevant descriptors for all of the trials by our technique. Lastly, in Fig. 9 it is shown that no molecular descriptors present in the Original CYP3A4 dataset were marked as relevant for the Enriched model. It is worth noticing that both models *Consensus 1* developed for CYP2C9 and CYP3A4 datasets by Nembri et al. [40] take into account ECFP as inputs to one of their base models, hence all of the ECFP fragments are considered to be present in the Original version of these two datasets.

Taking into consideration that all of our models outperformed the reference models reported by Nembri et al. [40] and Mansouri et al. [41], while at the same time identifying relevant molecular descriptors not included in the Original datasets, it is possible to conclude that meaningful information for the model might be

encoded in such molecular descriptors, and hence that relevant data could have been lost in the Original feature selection process. Furthermore, by means of this technique it was possible to identify molecular descriptors that were relevant to our models, and to find interesting relationships among them. Therefore, the proposed technique for *post hoc* feature analysis represents a way of providing interpretability to our neural-based models. One observation of the heat map visualizations is that they are practically limited by the maximum number of compounds that can be visually analyzed at the same time. Yet we note that an analysis on the top-20 or 30 features can be carried out with no problems as it was described previously.

## 6. Conclusions

QSAR modeling has become a key stage in the complex drug discovery process throughout the years. Upon the recent increase in the volume and quality of accessible datasets as well as computational power, more complex machine learning algorithms have established as current state of the art in QSAR modeling. Consensus approaches have consistently proven their efficacy for bioactivity prediction, but tend to lack interpretability and suffer from the limitations of their base classifiers. DNNs have not been widely adopted as a standard for QSAR modeling yet, although their effectiveness in solving high-dimensional problems make them a suitable technique for this area.

While DNN-based models attain higher predictive performance than other established techniques, they have their own challenges, such as low interpretability and proneness to overfitting. In this work we developed three neural-based QSAR models, which outperformed the state-of-the-art results for the three properties under study. At the same time we address the interpretability drawback without the need for performing feature selection. In addition, in this work we posed a strategy for analyzing the applicability domain of a neural-based QSAR model based on network output probabilities, which was shown to be correlated to the likelihood of correct classification. We also provided a technique based on an aggregation of the network weights for identifying the most relevant molecular descriptors and fingerprint fragments in a *post hoc* manner, which provides a sense of interpretability to our models. As future work we plan to investigate the impact of multi-task training, as a way of improving the performance of neural-based QSAR models.

## Declaration of competing interest

No author associated with this paper has disclosed any potential or pertinent conflicts which may be perceived to have impending conflict with this work. For full disclosure statements refer to <https://doi.org/10.1016/j.asoc.2019.105777>.

## Acknowledgments

This work is kindly supported by CONICET, Argentina, grant PIP 112-2012-0100471 and UNS, Argentina, grant PGI 24/N042. Authors thank MinCyT for its grant "PIDRI/PRH-2017-0007". Authors also thank Dr. Gustavo Vazquez for his help in the calculation of molecular descriptors.

## Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.asoc.2019.105777>.

## References

- [1] R. Vasundhara Devi, S. Siva Sathya, Mohane Selvaraj Coumar, Evolutionary algorithms for de novo drug design – A survey, *Appl. Soft Comput.* 27 (2015) 543–552, <http://dx.doi.org/10.1016/j.asoc.2014.09.042>.
- [2] Alfonso E. Márquez-Chamorro, Gualberto Asencio-Cortés, Cosme E. Santiesteban-Toca, Jesús S. Aguilar-Ruiz, Soft computing methods for the prediction of protein tertiary structures: A survey, *Appl. Soft Comput.* 35 (2015) 398–410, <http://dx.doi.org/10.1016/j.asoc.2015.06.024>.
- [3] Hongming Chen, Ola Engkvist, Yin Hai Wang, Marcus Olivecrona, Thomas Blaschke, The rise of deep learning in drug discovery, *Drug Discovery Today* 23 (6) (2018) 1241–1250, <http://dx.doi.org/10.1016/j.drudis.2018.01.039>.
- [4] Stephani Joy Y. Macalino, Vijayakumar Gosu, Sunhye Hong, Sun Choi, Role of computer-aided drug design in modern drug discovery, *Arch. Pharm. Res.* 38 (9) (2015) 1686–1701, <http://dx.doi.org/10.1007/s12272-015-0640-5>.
- [5] Yulan Wang, Jing Xing, Yuan Xu, Nannan Zhou, Jianlong Peng, Zhaoping Xiong, Xian Liu, Xiaomin Luo, Cheng Luo, Kaixian Chen, Mingyue Zheng, Hualiang Jiang, In silico ADME/T modelling for rational drug design, *Q. Rev. Biophys.* 48 (04) (2015) 488–515, <http://dx.doi.org/10.1017/S0033583515000190>.
- [6] Martin Eklund, Ulf Norinder, Scott Boyer, Lars Carlsson, Choosing feature selection and learning algorithms in QSAR, *J. Chem. Inf. Model.* 54 (3) (2014) 837–843, <http://dx.doi.org/10.1021/ci400573c>.
- [7] Ignacio Ponzoni, Víctor Sebastián-Pérez, Carlos Requena-Triguero, Carlos Roca, María J. Martínez, Fiorella Craverio, Mónica F. Díaz, Juan A. Páez, Ramón Gómez Arrayás, Javier Adrio, Nuria E. Campillo, Hybridizing feature selection and feature learning approaches in QSAR modeling for drug discovery, *Sci. Rep.* 7 (1) (2017) 2403, <http://dx.doi.org/10.1038/s41598-017-02114-3>.
- [8] Mohammad Goodarzi, Bieke Dejaegher, Yvan Vander Heyden, Feature selection methods in QSAR studies, *J. AOAC Int.* 95 (3) (2012) 636–651.
- [9] Kunal Roy, Pravin Ambure, Rahul B. Aher, How important is to detect systematic error in predictions and understand statistical applicability domain of QSAR models? *Chemometr. Intell. Lab. Syst.* 162 (2017) 44–54, <http://dx.doi.org/10.1016/j.chemolab.2017.01.010>.
- [10] Waldemar Klingspohn, Miriam Mathea, Antonius ter Laak, Nikolaus Heinrich, Knut Baumann, Efficiency of different measures for defining the applicability domain of classification models, *J. Cheminform.* 9 (1) (2017) 44, <http://dx.doi.org/10.1186/s13321-017-0230-2>.
- [11] Supratik Kar, Kunal Roy, Jerzy Leszczynski, Applicability domain: A step toward confident predictions and decidability for QSAR modeling, in: *Methods in Molecular Biology* (Clifton, N.J.), Vol. 1800, Springer, 2018, pp. 141–169, [http://dx.doi.org/10.1007/978-1-4939-7899-1\\_6](http://dx.doi.org/10.1007/978-1-4939-7899-1_6).
- [12] Watshara Shoombuatong, Philip Prathipati, Wiwat Owasirikul, Apilak Worachartcheewan, Saw Simeon, Nuttapat Anuwongcharoen, Jarl E.S. Wikberg, Chanin Nantasenamat, Towards the revival of interpretable QSAR models, in: *Advances in QSAR Modeling*, Springer, 2017, pp. 3–55, [http://dx.doi.org/10.1007/978-3-319-56850-8\\_1](http://dx.doi.org/10.1007/978-3-319-56850-8_1).
- [13] Kabiruddin Khan, Emilio Benfenati, Kunal Roy, Consensus QSAR modeling of toxicity of pharmaceuticals to different aquatic organisms: Ranking and prioritization of the DrugBank database compounds, *Ecotoxicol. Environ. Saf.* 168 (2019) 287–297, <http://dx.doi.org/10.1016/j.ecoenv.2018.10.060>.
- [14] Davide Ballabio, Francesca Grisoni, Viviana Consonni, Roberto Todeschini, Integrated QSAR models to predict acute oral systemic toxicity, *Mol. Inf.* (2018) <http://dx.doi.org/10.1002/minf.201800124>, [minf.201800124](https://doi.org/10.1002/minf.201800124).
- [15] Dimitar Dobchev, Mati Karelson, Have artificial neural networks met expectations in drug discovery as implemented in QSAR framework? *Expert Opin. Drug Discovery* 11 (7) (2016) 627–639, <http://dx.doi.org/10.1080/17460441.2016.1186876>.
- [16] Igor I. Baskin, Vladimir A. Palyulin, Nikolai S. Zefirov, Neural networks in building QSAR models, in: *Artificial Neural Networks*, Humana Press, 2006, pp. 133–154, [http://dx.doi.org/10.1007/978-1-60327-101-1\\_8](http://dx.doi.org/10.1007/978-1-60327-101-1_8).
- [17] Igor I. Baskin, David Winkler, Igor V. Tetko, A renaissance of neural networks in drug discovery, *Expert Opin. Drug Discovery* 11 (8) (2016) 785–795, <http://dx.doi.org/10.1080/17460441.2016.1201262>.
- [18] Audrey Cayot, Davy Laroche, Anne Disson-Dautriche, Anais Arbault, Jean-François Mailliefert, Paul Ornetti, Cytochrome P450 interactions and clinical implication in rheumatology, *Clin. Rheumatol.* 33 (9) (2014) 1231–1238, <http://dx.doi.org/10.1007/s10067-014-2710-3>.
- [19] David Rosner, Gerald Markowitz, Persistent pollutants: A brief history of the discovery of the widespread toxicity of chlorinated hydrocarbons, *Environ. Res.* 120 (2013) 126–133, <http://dx.doi.org/10.1016/j.envres.2012.08.011>.
- [20] Carlos J.S. Passos, Donna Mergler, Human mercury exposure and adverse health effects in the Amazon: a review, *Cad. Saud. Publ.* 24 Suppl 4 (2008) s503–20.
- [21] M.E. Pina, P. Coimbra, P. Ferreira, P. Alves, A.I. Figueiredo, M.H. Gil, Polymeric materials in ocular drug delivery systems, in: *Handbook of Polymers for Pharmaceutical Technologies*, John Wiley & Sons, Inc., Hoboken, NJ, USA, 2015, pp. 439–458, <http://dx.doi.org/10.1002/9781119041412.ch16>.



- [22] J.P. Hughes, S. Rees, S.B. Kalindjian, K.L. Philpott, Principles of early drug discovery, *Br. J. Pharmacol.* 162 (6) (2011) 1239–1249, <http://dx.doi.org/10.1111/j.1476-5381.2010.01127.x>.
- [23] Sheng Tian, Junmei Wang, Youyong Li, Dan Li, Lei Xu, The application of in silico drug-likeness predictions in pharmaceutical research, *Adv. Drug Deliv. Rev.* 86 (2015) 2–10, <http://dx.doi.org/10.1016/j.addr.2015.01.009>.
- [24] Gerhard Hessler, Karl-Heinz Baringhaus, Gerhard Hessler, Karl-Heinz Baringhaus, Artificial intelligence in drug design, *Molecules* 23 (10) (2018) 2520, <http://dx.doi.org/10.3390/molecules23102520>.
- [25] Roberto Todeschini, Davide Ballabio, Matteo Cassotti, Viviana Consonni, N3 and BNN: Two new similarity based classification methods in comparison with other classifiers, *J. Chem. Inf. Model.* 55 (11) (2015) 2365–2374, <http://dx.doi.org/10.1021/acs.jcim.5b00326>.
- [26] Sorin Avram, Alina Bora, Liliana Halip, Ramona Curpăn, Modeling kinase inhibition using highly confident data sets, *J. Chem. Inf. Model.* 58 (5) (2018) 957–967, <http://dx.doi.org/10.1021/acs.jcim.7b00729>.
- [27] Richard L. Marchese Robinson, Anna Palczewska, Jan Palczewski, Nathan Kidley, Comparison of the predictive performance and interpretability of random forest and linear models on benchmark data sets, *J. Chem. Inf. Model.* 57 (8) (2017) 1773–1792, <http://dx.doi.org/10.1021/acs.jcim.6b00753>.
- [28] Antonio Lavecchia, Machine-learning approaches in drug discovery: methods and applications, *Drug Discovery Today* 20 (3) (2015) 318–331, <http://dx.doi.org/10.1016/j.drudis.2014.10.012>.
- [29] Jiansong Fang, Ranyao Yang, Li Gao, Shengqian Yang, Xiaocong Pang, Chao Li, Yangyang He, Ai-Lin Liu, Guan-Hua Du, Consensus models for CDK5 inhibitors in silico and their application to inhibitor discovery, *Mol. Divers.* 19 (1) (2015) 149–162, <http://dx.doi.org/10.1007/s11030-014-9561-3>.
- [30] Tailong Lei, Youyong Li, Yunlong Song, Dan Li, Huiyong Sun, Tingjun Hou, ADMET Evaluation in drug discovery: 15. Accurate prediction of rat oral acute toxicity using relevance vector machine and consensus modeling, *J. Cheminform.* 8 (1) (2016) 6, <http://dx.doi.org/10.1186/s13321-016-0117-7>.
- [31] Vinicius M. Alves, Alexander Golbraikh, Stephen J. Capuzzo, Kammy Liu, Wai In Lam, Daniel Robert Korn, Diane Pozefsky, Carolina Horta Andrade, Eugene N. Muratov, Alexander Tropsha, Multi-descriptor read across (MuDRA): A simple and transparent approach for developing accurate quantitative structure–Activity relationship models, *J. Chem. Inf. Model.* 58 (6) (2018) 1214–1223, <http://dx.doi.org/10.1021/acs.jcim.8b00124>.
- [32] Daniel P. Russo, Kimberley M. Zorn, Alex M. Clark, Hao Zhu, Sean Ekins, Comparing multiple machine learning algorithms and metrics for estrogen receptor binding prediction, *Mol. Pharm.* 15 (10) (2018) 4361–4370, <http://dx.doi.org/10.1021/acs.molpharmaceut.8b00546>.
- [33] M.A. Hall, G. Holmes, Benchmarking attribute selection techniques for discrete class data mining, *IEEE Trans. Knowl. Data Eng.* 15 (6) (2003) 1437–1447, <http://dx.doi.org/10.1109/TKDE.2003.1245283>.
- [34] Hao Lin, Hui Ding, Feng-Biao Guo, Jian Huang, Prediction of subcellular location of mycobacterial protein using feature selection techniques, *Mol. Divers.* 14 (4) (2010) 667–671, <http://dx.doi.org/10.1007/s11030-009-9205-1>.
- [35] Kunal Roy, Asim Sattwa Mandal, Development of linear and nonlinear predictive QSAR models and their external validation using molecular similarity principle for anti-HIV indolyl aryl sulfones, *J. Enzyme Inhib. Med. Chem.* 23 (6) (2008) 980–995, <http://dx.doi.org/10.1080/14756360701811379>.
- [36] Berith F. Jensen, Christian Vind, Søren B. Padkjær, Per B. Brockhoff, Hanne H.F. Refsgaard, In silico prediction of cytochrome P450 2D6 and 3A4 inhibition using Gaussian kernel weighted k-nearest neighbor and extended connectivity fingerprints, including structural fragment analysis of inhibitors versus noninhibitors, *J. Med. Chem.* (2007) <http://dx.doi.org/10.1021/JM060333S>.
- [37] Feixiong Cheng, Yue Yu, Jie Shen, Lei Yang, Weihua Li, Guixia Liu, Philip W. Lee, Yun Tang, Classification of cytochrome P450 inhibitors and noninhibitors using combined classifiers, *J. Chem. Inf. Model.* 51 (5) (2011) 996–1011, <http://dx.doi.org/10.1021/ci200028n>.
- [38] Pranav Shah, Alexey Zakharov, R. Scott Obach, Anton Simeonov, Cornelis Hop, Dac-Trung Guyen, Eric Gonzalez, Hongmao Sun, Xin Xu, Development of a multitask deep learning QSAR model using data from individual cytochrome P450 isozymes, *Drug Metab. Pharmacokinet.* 33 (1) (2018) S35–S36, <http://dx.doi.org/10.1016/j.dmpk.2017.11.131>.
- [39] David Rogers, Mathew Hahn, Extended-connectivity fingerprints, *J. Chem. Inf. Model.* 50 (5) (2010) 742–754, <http://dx.doi.org/10.1021/ci100050t>.
- [40] Serena Nembri, Francesca Grisoni, Viviana Consonni, Roberto Todeschini, In silico prediction of cytochrome P450–drug interaction: QSARs for CYP3A4 and CYP2C9, *Int. J. Mol. Sci.* 17 (6) (2016) 914, <http://dx.doi.org/10.3390/ijms17060914>.
- [41] Kamel Mansouri, Tine Ringsted, Davide Ballabio, Roberto Todeschini, Viviana Consonni, Quantitative structure–activity relationship models for ready biodegradability of chemicals, *J. Chem. Inf. Model.* 53 (4) (2013) 867–878, <http://dx.doi.org/10.1021/ci4000213>.
- [42] Alberto Fernández, Robert Rallo, Francesc Giral, Prioritization of in silico models and molecular descriptors for the assessment of ready biodegradability, *Environ. Res.* 142 (2015) 161–168, <http://dx.doi.org/10.1016/j.envres.2015.06.031>.
- [43] Lidia Ceriani, Ester Papa, Simona Kovarich, Robert Boethling, Paola Gramatica, Modeling ready biodegradability of fragrance materials, *Environ. Toxicol. Chem.* 34 (6) (2015) 1224–1231, <http://dx.doi.org/10.1002/etc.2926>.
- [44] Davide Ballabio, Fabrizio Biganzoli, Roberto Todeschini, Viviana Consonni, Qualitative consensus of QSAR ready biodegradability predictions, *Toxicol. Environ. Chem.* (2016) 1–24, <http://dx.doi.org/10.1080/02772248.2016.1260133>.
- [45] Garrett B. Goh, Khusheem Sakloth, Charles Siegel, Abhinav Vishnu, Jim Pfafndtner, Multimodal Deep Neural Networks using Both Engineered and Learned Representations for Biodegradability Prediction. arXiv preprint [arXiv:1808.04456](https://arxiv.org/abs/1808.04456), aug 2018. doi: [arXiv:1808.04456v2](https://arxiv.org/abs/1808.04456v2).
- [46] María Jimena Martínez, Julieta Sol Dussaut, Ignacio Ponzoni, Biclustering as strategy for improving feature selection in consensus QSAR modeling, *Electron. Notes Discrete Math.* 69 (2018) 117–124, <http://dx.doi.org/10.1016/j.endm.2018.07.016>.
- [47] Erik Gawehn, Jan A. Hiss, Gisbert Schneider, Deep learning in drug discovery, *Mol. Inf.* 35 (1) (2016) 3–14, <http://dx.doi.org/10.1002/minf.201501008>.
- [48] David A. Winkler, Neural networks as robust tools in drug lead discovery and development, *Mol. Biotechnol.* 27 (2) (2004) 139–168, <http://dx.doi.org/10.1385/MB:27:2:139>.
- [49] Junshui Ma, Robert P. Sheridan, Andy Liaw, George E. Dahl, Vladimir Svetnik, Deep neural nets as a method for quantitative structure–activity relationships, *J. Chem. Inf. Model.* 55 (2) (2015) 263–274, <http://dx.doi.org/10.1021/ci500747n>.
- [50] Eelke B. Lenselink, Niels ten Dijke, Brandon Bongers, George Papadatos, Herman W.T. van Vlijmen, Wojtek Kowalczyk, Adriaan P. Ijzerman, Gerard J.P. van Westen, Beyond the hype: deep neural networks outperform established methods using a ChEMBL bioactivity benchmark set, *J. Cheminform.* 9 (1) (2017) 45, <http://dx.doi.org/10.1186/s13321-017-0232-0>.
- [51] Alexios Koutsoukas, Keith J. Monaghan, Xiaoli Li, Jun Huan, Deep learning: investigating deep neural networks hyper-parameters and comparison of performance to shallow methods for modeling bioactivity data, *J. Cheminform.* 9 (1) (2017) 42, <http://dx.doi.org/10.1186/s13321-017-0226-y>.
- [52] Ester Papa, Simona Kovarich, Paola Gramatica, Development, validation and inspection of the applicability domain of QSPR models for physicochemical properties of polybrominated diphenyl ethers, *QSAR Comb. Sci.* 28 (8) (2009) 790–796, <http://dx.doi.org/10.1002/qsar.200806183>.
- [53] Kunal Roy, Supratik Kar, Pravin Ambure, On a simple approach for determining applicability domain of qsar models, *Chemometr. Intell. Lab. Syst.* 145 (2015) 22–29, <http://dx.doi.org/10.1016/j.chemolab.2015.04.013>.
- [54] Ruifeng Liu, Hao Wang, Kyle P. Glover, Michael G. Feasel, Anders Walqvist, Dissecting machine-learning prediction of molecular activity: Is an applicability domain needed for quantitative structure–activity relationship models based on deep neural networks? *J. Chem. Inf. Model.* (2018) <http://dx.doi.org/10.1021/acs.jcim.8b00348>, [acs.jcim.8b00348](https://doi.org/10.1021/acs.jcim.8b00348).
- [55] Francois Berenger, Yoshihiro Yamanishi, A distance-based boolean applicability domain for classification of high throughput screening data, *J. Chem. Inf. Model.* (2019) <http://dx.doi.org/10.1021/acs.jcim.8b00499>, [acs.jcim.8b00499](https://doi.org/10.1021/acs.jcim.8b00499).
- [56] Pavel Polishchuk, Interpretation of quantitative structure–activity relationship models: Past, present, and future, *J. Chem. Inf. Model.* 57 (11) (2017) 2618–2639, <http://dx.doi.org/10.1021/acs.jcim.7b00274>.
- [57] Ravid Schwartz-Ziv, Naftali Tishby, Opening the Black Box of Deep Neural Networks via Information. arXiv preprint [arXiv:1703.00810](https://arxiv.org/abs/1703.00810), mar2017.
- [58] Alfredo Vellido, José David Martín-Guerrero, Paulo J. G. Lisboa, Making machine learning models interpretable, *ESANN*, 2012.
- [59] Seyran Khademi, Xiangwei Shi, Tino Mager, Ronald Siebes, Carola Hein, Victor de Boer, Jan van Gemert, Sight-seeing in the eyes of deep neural networks, in: 2018 IEEE 14th International Conference on E-Science (E-Science), IEEE, 2018, pp. 407–408, <http://dx.doi.org/10.1109/eScience.2018.00125>.
- [60] Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Methods for interpreting and understanding deep neural networks, *Digit. Signal Process.* 73 (2018) 1–15, <http://dx.doi.org/10.1016/j.dsp.2017.10.011>.
- [61] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, Wojciech Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, in: Oscar Deniz Suarez (Ed.), *PLOS ONE* 10 (7) (2015) e0130140, <http://dx.doi.org/10.1371/journal.pone.0130140>.
- [62] Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin, "why should i trust you?", in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16, ACM Press, New York, New York, USA, ISBN: 9781450342322, 2016, pp. 1135–1144, <http://dx.doi.org/10.1145/2939672.2939778>.

- [63] Michael Tsang, Dehua Cheng, Yan Liu, Detecting Statistical Interactions from Neural Network Weights. arXiv preprint [arXiv:1705.04977](https://arxiv.org/abs/1705.04977), may 2017.
- [64] Yadi Zhou, Sunita Cahya, Steven A. Combs, Christos A. Nicolaou, Jibo Wang, Prashant V. Desai, Jie Shen, Exploring tunable hyperparameters for deep neural networks with industrial ADME data sets, *J. Chem. Inf. Model.* (2019) [http://dx.doi.org/10.1021/acs.jcim.8b00671](https://doi.org/10.1021/acs.jcim.8b00671), acs.jcim.8b00671.
- [65] Kode srl, Dragon (software for molecular descriptor calculation). Pisa, Italy, 2016. URL [https://chm.kode-solutions.net/products\\_dragon.php](https://chm.kode-solutions.net/products_dragon.php). Version 7.0.
- [66] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, Xiaoqiang Zheng, TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, URL <https://www.tensorflow.org/>, Software available from tensorflow.org, 2015.
- [67] Vinod Nair, Geoffrey E. Hinton, Rectified Linear Units Improve Restricted Boltzmann Machines, undefined, 2010.
- [68] David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams, Learning representations by back-propagating errors, *Nature* 323 (6088) (1986) 533–536, [http://dx.doi.org/10.1038/323533a0](https://doi.org/10.1038/323533a0).
- [69] Diederik P. Kingma, Jimmy Ba, Adam: A Method for Stochastic Optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980), dec 2014.
- [70] Xavier Glorot, Yoshua Bengio, Understanding the difficulty of training deep feedforward neural networks, mar 2010. URL <http://proceedings.mlr.press/v9/glorot10a.html>.
- [71] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, Delving deep into rectifiers: Surpassing human-level performance on imagenet classification, in: Proceedings of the IEEE International Conference on Computer Vision, 2015, pp. 1026–1034.
- [72] Dmytro Mishkin, Jiri Matas, All you need is a good init. arXiv preprint [arXiv:1511.06422](https://arxiv.org/abs/1511.06422), nov 2015.
- [73] Dan Hendrycks, Kevin Gimpel, Generalizing and Improving Weight Initialization, undefined, 2016.
- [74] Sergey Ioffe, Christian Szegedy, Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, undefined, 2015.
- [75] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov, Dropout: A simple way to prevent neural networks from overfitting, *J. Mach. Learn. Res.* 15 (2014) 1929–1958.
- [76] Arthur E. Hoerl, Robert W. Kennard, Ridge regression: Biased estimation for nonorthogonal problems, *Technometrics* 12 (1) (1970) 55–67, [http://dx.doi.org/10.1080/00401706.1970.10488634](https://doi.org/10.1080/00401706.1970.10488634).
- [77] Joanna Jaworska, Nina Nikolova-Jeliazkova, Tom Aldenberg, QSAR Applicability domain estimation by projection of the training set descriptor space: a review., *Altern. Lab. Animals : ATLA* 33 (5) (2005) 445–459.
- [78] Shane Weaver, M. Paul Gleeson, The importance of the domain of applicability in QSAR modeling, *J. Mol. Graph. Modelling* 26 (8) (2008) 1315–1326, [http://dx.doi.org/10.1016/j.jmglm.2008.01.002](https://doi.org/10.1016/j.jmglm.2008.01.002).
- [79] Christopher D. Manning, Prabhakar. Raghavan, Hinrich. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, 2008, p. 482.
- [80] María Jimena Martínez, Ignacio Ponzoni, Mónica F Díaz, Gustavo E Vazquez, Axel J Soto, Visual analytics in cheminformatics: user-supervised descriptor selection for QSAR methods, *J. Cheminform.* 7 (1) (2015) 39, [http://dx.doi.org/10.1186/s13321-015-0092-4](https://doi.org/10.1186/s13321-015-0092-4).
- [81] Alexey V. Zakharov, Megan L. Peach, Markus Sitzmann, Marc C. Nicklaus, QSAR Modeling of imbalanced high-throughput screening data in pubchem., *J. Chem. Inf. Model.* 54 (3) (2014) 705–712, [http://dx.doi.org/10.1021/ci400737s](https://doi.org/10.1021/ci400737s).