

# Improving drug safety predictions by reducing poor analytical practices

Stanley E. Lazic<sup>1,\*</sup>, Dominic P. Williams<sup>2</sup>

<sup>1</sup>*Prioris.ai Inc., 459–207 Bank Street, Ottawa, K2P 2N2, Canada*

<sup>2</sup>*Functional and Mechanistic Safety, Clinical Pharmacology and Safety Sciences, AstraZeneca, R&D, Cambridge, CB4 0WG, UK*

\*Corresponding author: [stan.lazic@cantab.net](mailto:stan.lazic@cantab.net)

## Abstract

Predicting the safety of a drug from preclinical data is a major challenge in drug discovery, and progressing an unsafe compound into the clinic puts patients at risk and wastes resources. Methods and analytic decisions known to provide poor predictions are common in drug safety pharmacology and related fields, which include creating arbitrary thresholds, binning continuous values, giving all assays equal weight, and multiple reuse of information. In addition, the metrics used to evaluate models often omit important criteria and assessing how models perform on new data are often insufficient. Prediction models with these problems are unlikely to perform well, and published models suffer from many of these issues. We describe these problems in detail, often demonstrate their negative consequences, and propose simple solutions that are standard in other disciplines where predictive modelling is used.

## Introduction

In early-stage drug discovery, compounds are tested in assays to determine their likelihood of causing organ toxicity or adverse events in the clinic. The assays often test specific mechanisms such as blocking key ion channels, inhibiting mitochondrial function, causing cell death, or damaging DNA. Based on these assay results and other information about the compounds such as their structural, physical, and chemical properties, project teams decide whether to progress a compound to animal studies, and eventually to human trials. At each stage of drug discovery pipeline, decisions are made based on current data, and quantitative methods are developed to help make these decisions<sup>1</sup>.

Many current practices make inefficient use of the data and can be misleading. Below we describe these practices, illustrate their limitations, and demonstrate the benefits of alternative

approaches. These recommendations complement existing guidelines and best practices on *in silico* toxicology models<sup>2,3</sup>.

## Problems with current practices

### Using safety margins as predictors

Compounds are often tested in concentration-response assays and the results are typically summarised by an  $XC_{50}$  value, which is a measure of the potency of the compound, and which lies between the lowest and highest concentrations tested. Another important quantity for safety assessment is the peak or maximum concentration that a compound reaches in the blood ( $C_{max}$ ) – either actual or predicted.  $C_{max}$  is related to the clinical dose, and compounds given at higher doses tend to have a greater risk for toxic side effects, making  $C_{max}$  one of the best predictors of clinical toxicity<sup>4-6</sup>.  $XC_{50}$  values are often divided by  $C_{max}$  values to give a safety margin or safety factor, where higher values indicate greater safety. Although safety margins are useful and easy to interpret, they have five problems when used to make predictions. The first problem is minor: an assumption is created that all combinations of  $XC_{50}$  and  $C_{max}$  values with the same margin have the same risk. For example, 1/0.1, 10/1, 100/10 all have a margin of 10, and therefore are assumed to have the same risk. This may be a reasonable assumption, but it should be verified.

A more damaging consequence of using margins is that mechanistic assays with zero predictive ability for clinical toxicity can appear useful. For example, mechanistic assays are used to design out specific liabilities such as mitochondria inhibition or BSEP transporter protein inhibition. Even if the primary purpose of these assays is for drug design, they are often used in models to predict clinical toxicity, but the margin values are used instead of the  $XC_{50}$  values. Using all available data is a good idea, but we show in Figure 1 with simulated data how margins can mislead.  $XC_{50}$  values for two assays are independently simulated from a uniform distribution for 100 compounds, and an outcome of clinical interest is also independently simulated from a uniform distribution. Figure 1A and 1B show no association between the assay values and the outcome, indicating that the assays have zero predictive ability. Next we simulate  $C_{max}$  values that strongly predict the outcome (Fig. 1C), and then calculate the safety margins. Both margins are now good predictors of the outcome (Fig. 1D and 1E), but this is completely driven by  $C_{max}$ , the assays have contributed nothing to the prediction.

A third problem is that the safety margin variables are now correlated, owing to their dependence on  $C_{max}$  (Fig. 1F), whereas the original assay values are uncorrelated (not shown, but the values were simulated to be independent). Correlated variables are a problem when creating prediction models because the models become less stable, the precision of the estimates decreases, and the apparent importance of the variables can decrease.

A fourth problem is that the predictive information contained in the  $C_{max}$  values is used multiple times and thus  $C_{max}$  receives greater weight. Even if the assays are predictive,  $C_{max}$  is one piece of information that should be used only once, not multiple times for each margin. Furthermore,  $C_{max}$  is often measured with considerable uncertainty due to variability between

patients, or it may be predicted from *in vitro* and *in silico* data. This measurement noise is then added to all assay values.

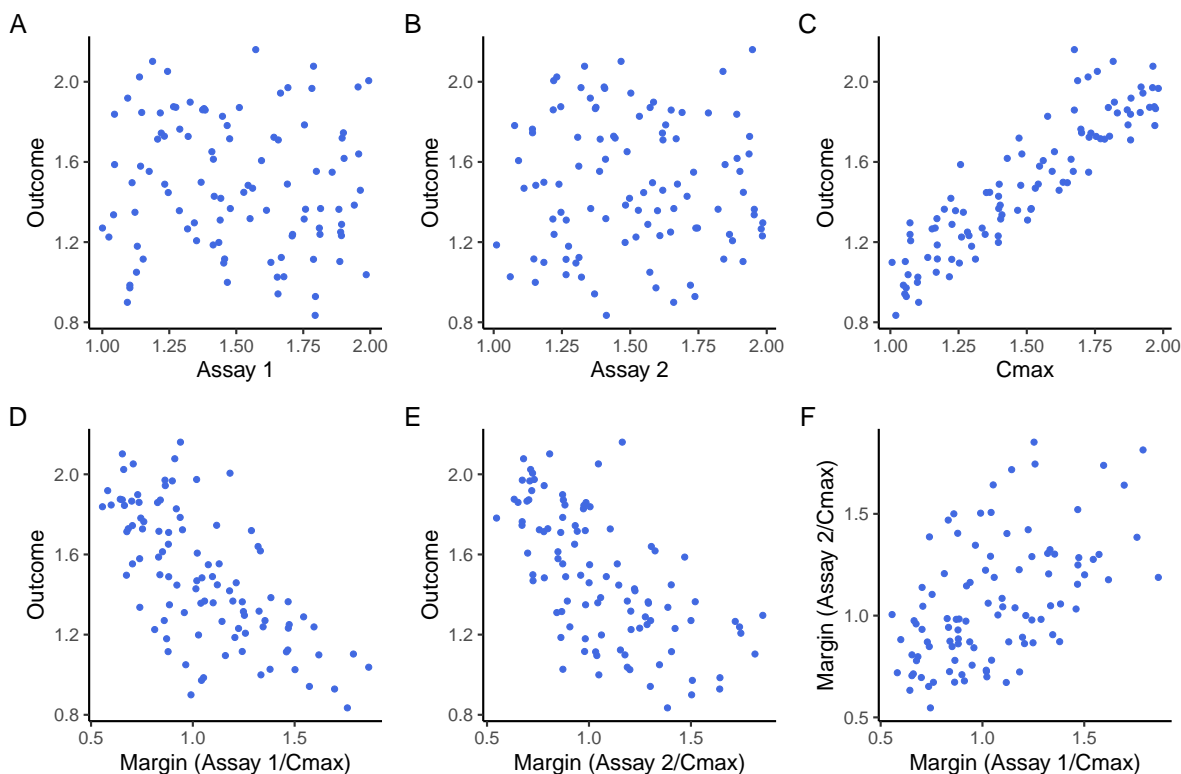


Figure 1: Using safety margins can mislead. Neither Assay 1 (A) nor Assay 2 (B) predict the outcome, but  $C_{\max}$  is strongly predictive (C). Safety margins for Assay 1 (D) and Assay 2 (E) appear predictive, but these relationships are completely driven by  $C_{\max}$ . The assay results are now correlated, making prediction models unstable (F).

Finally, by combining two variables into one, we've reduced our flexibility to make predictions. For example, Figure 1A in Lazic et al. plots margin values for a cardiotoxicity model, and the only option to separate the two classes of compounds is the location of a vertical line<sup>5</sup>. Figure 1B in that paper plots  $C_{\max}$  and the  $IC_{50}$  values from a hERG inhibition assay on two axes, and now it is easier to separate the groups with a line in two dimensions, as both the slope and intercept provide greater flexibility. Furthermore, we have even more flexibility because the separating boundary need not be a straight line – a curved line might be preferable. Reducing the dimensions of the data by calculating a ratio of two values can be beneficial when faced with too many predictor variables, but the trade-off is reduced flexibility. Given all these problems with safety margins, they cannot be recommended for inputs into a predictive model. Instead,  $C_{\max}$  and  $XC_{50}$  values should be used directly.

## Binning and arbitrary thresholds

A common practice is to categorise compounds into “Active” or “Not active” in an assay based on their  $XC_{50}$  values or safety margins. This makes as much scientific sense as measuring people’s height with a noisy ruler and then categorising them into “giants” and “dwarfs” based on an arbitrary cutoff. There is little gained by such a procedure and much is lost, as statisticians have repeatedly argued<sup>7-17</sup>.

$XC_{50}$  values, like human height, are continuous variables (at least to the measured precision) and should be treated as such. Binning compounds into Active/Not active has three problems. First, considerable information is discarded. All compounds within a category are considered equal, but scientific sense tells us that a compound just below a threshold will differ from another compound in the same category that is far from the threshold. Similarly, compounds just on either side of a threshold do not have dramatically different risks. Second, binning creates categories that do not exist in Nature, but are then taken as real. This is known as the reification fallacy, or the fallacy of misplaced concreteness, where a concept is mistaken for a physical property. This leads to surprise when compounds that are “negative” in all assays show clinical toxicity later on. Being negative in an assay is not a physical property of a compound, it’s an arbitrary designation we’ve assigned. A compound with a weak signal in several assays may be deemed safe because the  $XC_{50}$  values were below a threshold, but this does not reflect the true risk. An additive or multiplicative combination of several weak signals could potentially indicate a high-risk compound.

Third, binning ignores the often substantial uncertainty in the assay values – 95% confidence intervals for  $XC_{50}$  values can sometimes span several orders of magnitude. If an assay were run again, compounds close to the threshold could fall on the opposite side. This procedure therefore introduces misclassification errors in the data. Thresholds are only needed when a decision or action is taken, which occurs after a prediction has been made (or based on assay values directly if no prediction is made; for example, selecting compounds to progress from a high-throughput screen). When building predictive models, the numeric assay values should be used.

## Deriving arbitrary “scores”

To rank compounds or make a decision, the information contained in multiple variables that predict clinical safety need to be combined. These variables are often on different scales and are not directly comparable; for example, assay results are often  $XC_{50}$  values, physical and chemical properties such as cLogP may have no units, and compounds satisfying Lipinski’s Rule-of-5 may be binary Yes/No variables. One motivation for binning continuous variables is that all become binary and can be combined into a score by summing them, but ensuring that the value of “1” is interpreted as “higher risk” for all variables:

$$\text{Score} = \text{Variable}_1 + \text{Variable}_2 + \dots + \text{Variable}_n. \quad (1)$$

The score can range from zero to the number of variables  $n$ , and the hope is that the

score predicts the clinical outcome. We refer to this as the “bin-and-sum” method, which was previously used at AstraZeneca and other companies<sup>18,19</sup>. Fortunately, the above score equation is the beginning of a standard statistical model. In Equation 2 below we convert the above score equation into a statistical model by using the measured variable values directly to predict the clinical outcome  $y$ , which could be QT<sub>c</sub> prolongation for a cardiotoxicity model or liver enzyme levels for a liver toxicity model (we assume  $y$  is a continuous variable for this discussion, but the equation below can be easily modified for other types of outcomes). We also introduce parameters ( $\beta$ 's, which are called weights in the machine learning literature) that are estimated from the data. The above bin-and-sum method also implicitly has parameters, but they are fix to  $\beta_0 = 0$  and the other  $\beta$ 's to a value of one and therefore have no influence. Why fix the  $\beta$ 's? It is better to learn their optimal values from the data. Note how in the equation below we are using the variables to directly predict the outcome  $y$ , not to sum to a score which we hope will be associated with the outcome in a later step. Also, the variables do not need to be on the same scale or have the same units.

$$y = \beta_0 + \beta_1 \text{Variable}_1 + \beta_2 \text{Variable}_2 + \dots + \beta_n \text{Variable}_n \quad (2)$$

To complete the statistical model we also need an error term ( $\varepsilon$ ), which captures the difference between the prediction based on the variables and the true value of the outcome:

$$y = \beta_0 + \beta_1 \text{Variable}_1 + \beta_2 \text{Variable}_2 + \dots + \beta_n \text{Variable}_n + \varepsilon. \quad (3)$$

The bin-and-sum method has several problems. First, the predictor variables are binned, which was discussed in the previous section. Second, all variables are given equal weight in calculating the score, implying that all variables are equally effective in predicting the outcome, but this is unrealistic. Third, the outcome – the target of our prediction – is not used to create the prediction rule, which is an inefficient use of the data. Fourth, even simple interactions between variables cannot be captured. Figure 2 illustrates this point with simulated data. 100 compounds are defined as either toxic (red triangles) or safe (blue circles) based on a clinical outcome. These compounds are run on two assays (Fig. 2A and 2B) and the measured values are obtained (e.g. XC<sub>50</sub> values). Both assays are predictive as they (imperfectly) separate the toxic/safe groups, and the vertical dashed lines are optimal thresholds that minimise the misclassification error. If a compound is to the right of the dashed line it is considered “positive” in the assay and given a value of 1, otherwise a value of 0. Combining the assay results by summing these values results in a total score of either 0, 1, or 2. Compounds with a score of 2 fall in the top right quadrant of Figure 2C. These compounds are all toxic, but many toxic compounds are in other quadrants. Compounds with a score of at least 1 fall into the top left, top right, and bottom right quadrants. Although all the toxic compounds are in these three quadrants, so are some safe compounds.

Table 1 reports the sensitivity, specificity, and accuracy of each assay and their combination using two rules: call a compound “toxic” if it is positive in at least one assay, or call a compound “toxic” if it is positive in both assays. This approach is sub-optimal compared with fitting a statistical model to the data, which perfectly classifies the compounds, indicted by the solid

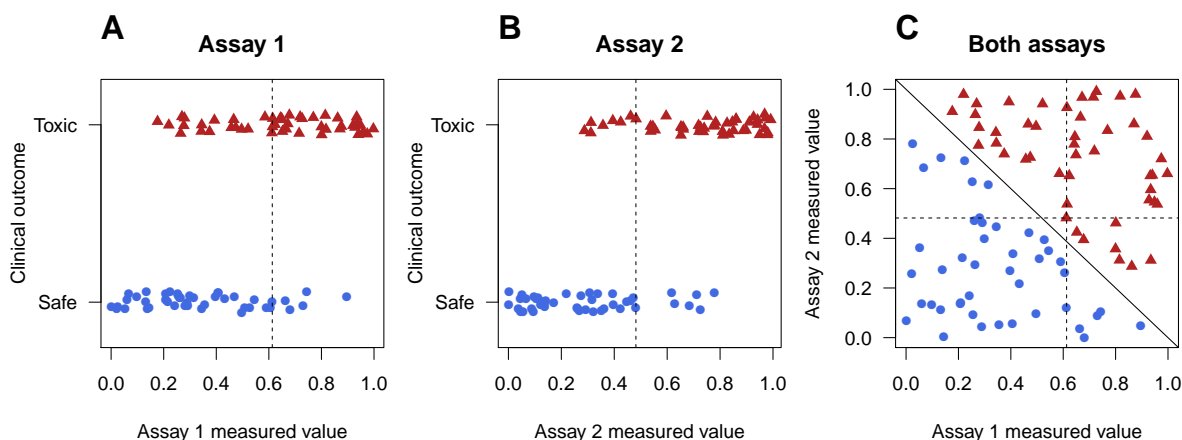


Figure 2: The bin-and-sum method misses simple patterns. Two assays predict a clinical outcome (A, B), but summing the number of positive assay results for each compound provides poor predictions. These compounds can be perfectly separated with a statistical model, represented by the diagonal solid line (C). Dashed lines are the optimal univariate thresholds.

diagonal line (Fig. 2C). The dashed lines are the same optimal univariate thresholds obtained from Figure 2A and 2B.

[Insert Table 1]

The final problem with the bin-and-sum method is that it cannot capture cases when two or more weak signals are highly predictive of toxicity. Alternatively, values for two or more assays may be correlated because they are measuring a similar process or mechanism, such as cell death. Compounds will tend to be positive or negative on all the correlated assays, leading to “double-” or “triple-counting”. Both of these effects can be captured with interaction terms in a statistical model.

## Methods to prevent overfitting are inadequate

Overfitting occurs when a prediction model or prediction rule is too flexible and accounts for idiosyncratic aspects of the available data that will not be observed in future data. An overfit model will predict the current data well but will predict future data poorly, and it is very easy to overfit. Figure 3 illustrates overfitting, where data were simulated for 100 compounds, 50 toxic and 50 safe. Four predictor variables ( $x_1$  to  $x_4$ ) were generated from a uniform distribution and have no predictive value. Even so, we can build a decision tree with 70% accuracy using the following decision rules (Fig. 3A). First, if a compound has a value for  $x_3$  less than 0.17, classify it as toxic (right branch). 14 out of 18 compounds with  $x_3 < 0.17$  are toxic. If a compound has a value for  $x_3 \geq 0.17$  then variable  $x_1$  is used for the next decision. If  $x_1 \geq 0.9$  then a compound is classified as safe, and all 7 compounds in this far-left branch are correctly classified. And so on for the other decision points. This tree had several constraints to minimise the complexity, otherwise, we could keep creating decision

points until we perfectly classify all compounds. The constraints are: (1) a maximum of four splits could be made (potentially one for each variable), a split was not made if there were fewer than 20 compounds in a branch, and a split was not made if it would lead to a final node with less than five compounds. Despite these constraints, we achieved 70% accuracy when the predictor variables are random noise.

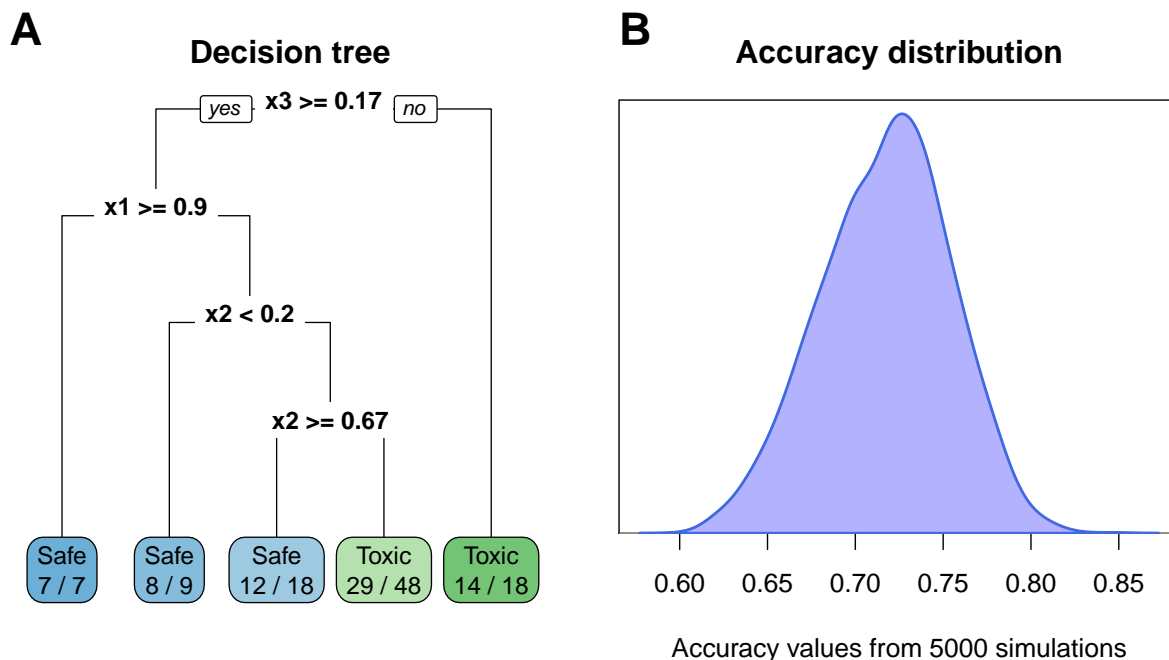


Figure 3: Overfitting. A decision tree predicts a toxic or safe outcome using four variables ( $x_1$  to  $x_4$ ) that are pure noise with 70% accuracy (A). Simulating this process with 5000 datasets gives a distribution of accuracy values with a mean of 72% and a range of 60% to 85% (B).

To better appreciate how easy it is to get a respectable prediction with pure noise, 5000 datasets were generated and analysed as above, and the accuracy values are plotted in Figure 3B. The average accuracy was 72%, with a range of 60% to 85%.

A common way to prevent overfitting is to split the data into a training set that is used to develop a model and a test set that is used to evaluate the final model. When building a model, many analytical decisions must be made and values for parameters that are not learned from the data must be selected. Cross-validation and bootstrapping are common methods to ensure that decisions don't lead to overfitting<sup>20</sup>. A benefit of the recent interest in machine learning and artificial intelligence is a greater awareness of assessing predictive performance on an independent test set. If standard machine learning methods are used for these prediction problems then often standard workflows will include an assessment of overfitting. However, many prediction models in industry are still hand-crafted or created by those with little training in predictive modelling.

## Ignoring uncertainty

Uncertainty exists in almost all aspects of making predictions, and for most toxicology and safety pharmacology applications, uncertainty needs to be incorporated into models and reported with the final predictions. For example, if a climate model predicts a 2°C increase in global temperature in the next 20 years, one's interpretation differs if the uncertainty in this value is 1.8 to 2.2°C versus -8 to 12°C. In the latter case the uncertainty is so great that no sensible actions can be taken (should we prevent global warming or global cooling?).

As mentioned above, the data are a source of uncertainty; assay values are often measured with considerable error and some predictor variables may themselves be predictions, such as  $C_{\max}$  values. Parameters in statistical models ( $\beta$ 's in Equations 2 and 3 above) are also uncertain. Since predictions from such models are based on the  $\beta$  values, greater uncertainty in the parameters will lead to greater uncertainty in the predictions. As the sample size increases, uncertainty in parameters decreases, but more samples enables more complex models to be fit, and so parameter uncertainty rarely reduces to zero in practice, even with a large sample size.

The models used to make predictions are also uncertain; we do not have a “true” model that we take as given and use for predictions. Many models may provide similar predictions, with no single model dominating. For example, if we have ten  $x$  variables (assays, physical/chemical properties,  $C_{\max}$ , etc.) in Equation 3, which combination of these variables should we use in a regression model? Do we include any interactions to allow for the multiplicative effect of weak signals, and if so, between which predictors? Are all the relationships linear between the  $x$ 's and the outcome, or do we expect some “U” or inverted-U relationships? We also need not restrict ourselves to a regression model but can also consider random forests, support vector machines, or neural network models. Although many predictive models can be created, predictions are typically only made from one model and thus model uncertainty is ignored. Model uncertainty can be incorporated into predictions using ensembles, stacking, or model-averaging<sup>21,22</sup>.

Finally, even if all the above sources of uncertainty are reduced to near zero, the predictions themselves will still be uncertain because we can rarely make perfect predictions of clinical outcomes from preclinical assays and the properties of compounds. Unfortunately, many prediction models only report a point prediction or the “best estimate”, and this is used to make decisions. Bayesian methods provide a principled approach to quantify, integrate, and report uncertainty in predictions, and which are now used in production for cardiotoxicity and liver toxicity<sup>5,6</sup>, and are increasingly being developed for other applications<sup>23–28</sup>. Note that we are not referring to naive Bayes classifiers, which do not place priors on all unknowns and update them with data<sup>29</sup>, but rather fully Bayesian versions of standard statistical or machine learning methods, such as logistic regression and neural networks, as well as Gaussian processes and Bayesian Additive Regression Trees (BART).



## Narrowly focusing on a few metrics

After a model is developed, its performance must be assessed, but the metrics commonly examined are limited. The discussion below mainly applies to classification tasks, where the aim is to predict a binary outcome (safe/toxic), but is also relevant for categorical outcomes (safe/apoptosis/necrosis) and ordered categorical outcomes (safe/mild/severe). Four types of metrics should be examined: classification, discrimination, calibration, and overall fit.

*Classification* metrics are the most common and include accuracy, sensitivity, and specificity. Although easy to interpret, they can be misleading when the classes are unbalanced. For example, if 90% of compounds are non-toxic, we can achieve 90% accuracy by ignoring any experimental data and always predicting “safe”. A better option is to report the balanced accuracy, which is the average accuracy of the two classes. Table 2 shows the result of a hypothetical prediction model. The standard accuracy calculation is 92%:

$$\text{Accuracy} : \frac{85 + 7}{85 + 7 + 3 + 5} \times 100 = 92\%.$$

[Insert Table 2]

Such a high accuracy may appear impressive, but we can get 90% accuracy by always predicting a compound is safe. The balanced accuracy is only 82%, which indicates that the model is worse than always predicting the most frequent class, at least based on the accuracy criterion:

$$\text{Balanced accuracy} : \left( \frac{85}{85 + 5} + \frac{7}{7 + 3} \right) / 2 \times 100 = 82\%.$$

A second problem is that to calculate these metrics we have silently converted a continuous prediction probability into a class using a threshold that may not be meaningful. A logical threshold is 0.5, with compounds classified as toxic if the prediction is  $> 0.5$ , and safe if the prediction is  $\leq 0.5$ . However, another threshold might better reflect the costs of misclassifying a safe compound as toxic, and vice versa. Furthermore, this is another example of losing information by binning. Many predictive methods such as logistic regression, naive Bayes, and neural networks provide a continuous value for the probability of toxicity, and models are better assessed on this. Consider a model (or a person) that predicts the probability of toxicity for a compound is 0.55, and another model predicts 0.95. If the compound is indeed toxic, and if we take 0.5 as a threshold for toxicity, both models made correct predictions, but the second made a much stronger prediction and should therefore be considered a better model. Similarly, if the drug is safe, the model that predicted a 0.95 probability of toxicity should be penalised more than the model that predicted 0.55. Both models made wrong predictions, but one was more wrong than the other. Methods that assess the quality of a prediction are known as scoring rules or scoring functions, and good scoring rules distinguish between the above situations<sup>30</sup>. Accuracy fails, but so do sensitivity and specificity, which also depend on the balance of the classes and require converting a continuous probability into a binary category. Furthermore, sensitivity and specificity do not tell us what we want to know when making a

prediction for a new compound, which is the probability that a compound is safe, given that it is predicted to be safe:  $P(\text{safe} \mid \text{predicted safe})$ . Instead they tell us the probability that a compound is predicted to be safe, given that it is actually safe:  $P(\text{predicted safe} \mid \text{actually safe})$ <sup>31</sup>. Although classification metrics tend receive considerable attention, they do not tell the full story.

*Discrimination* is the ability of a model to distinguish between classes. A common measure is the concordance index (c-index), which equals the area under the receiver-operating curve (AUC) when the variable to be predicted has two classes<sup>32</sup>. The c-index is a number between 0 and 1 giving the probability that a randomly selected toxic compound has a predicted score higher than a randomly selected non-toxic compound. Perfect discrimination gives a c-index of 1 and a model with no predictive ability will give a value of 0.5. The c-index requires no arbitrary threshold and is unaffected by class imbalances. We use data from Pollard et al. to illustrate discrimination and the other metrics<sup>33</sup>. This dataset is used to predict QT prolongation, a binary (Yes = 1 / No = 0) outcome based on clinical  $C_{\text{max}}$  values and  $IC_{50}$  values from a hERG inhibition assay. Further details on the data and model can be found in Lazic et al.<sup>5</sup>, and in the supplementary material. The c-index for the model equals 0.94, and for comparison, accuracy = 89.7%, sensitivity = 90.9%, and specificity = 88.9%. Although uncommon, the c-index can be represented as a number between 0-100% if the other metrics are also on a percentage scale.

*Calibration* measures the agreement between the predictions and the true outcomes. A model is calibrated when the predicted probabilities match the observed probabilities, for all values of the predicted probabilities. Even if a model has reasonable accuracy, it can still over- or under-predict the true probabilities, and this is often assessed graphically. Figure 4A plots the predicted probabilities on the x-axis and the observed proportions on the y-axis, with perfectly calibrated predictions falling along the diagonal red line. The curved blue line shows the results from the model and the grey band indicates the 95% confidence interval. The blue line is close to the red, but with few compounds, substantial uncertainty is present. Calibration curves and methods to summarise them are discussed in<sup>20,31</sup>.

The *overall fit* of a model is typically used by machine learning algorithms to optimise the model, but it less common to report these metrics as they are not intuitive and difficult to interpret. One metric that is relatively easy to understand is the Brier score (BS)<sup>34</sup>. It is appropriate for predicting binary outcomes with models that produce a probability (see Semenova et al. for an extension to an ordered categorical outcome<sup>28</sup>). The BS is the squared difference between the prediction – a value between 0 and 1 – and the actual category, which is either 0 or 1:

$$BS = \frac{\sum_{i=1}^N (\text{predicted}_i - \text{actual}_i)^2}{N}.$$

$N$  is the number of compounds and  $i$  indexes the compounds. The greater the difference between the actual and predicted values, the larger the Brier score, and so low scores are better. We get a BS for each compound, and by summing all the scores and dividing by  $N$  we get an average BS for the dataset. The grey squares in Figure 4B are the true values for QT

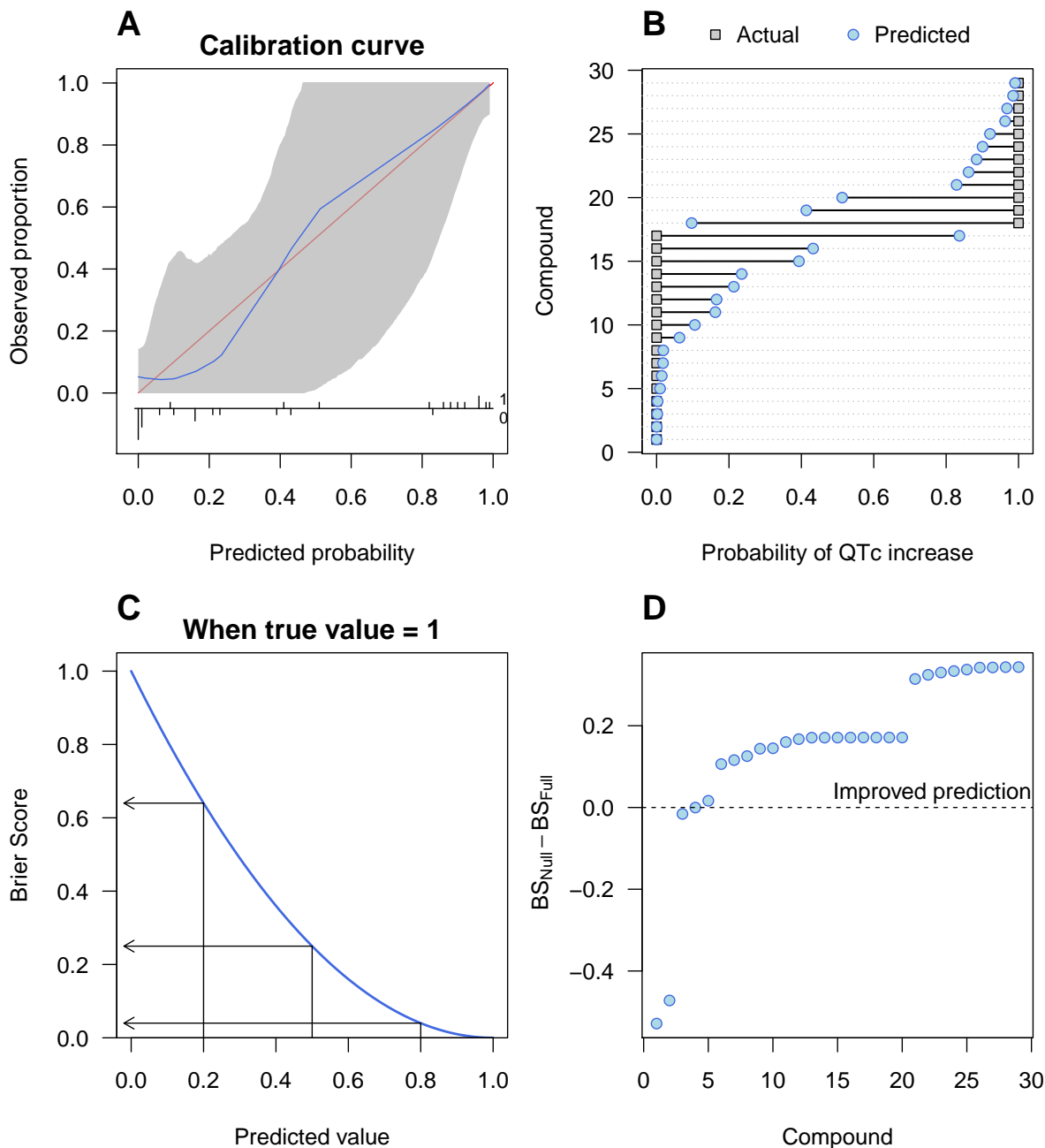


Figure 4: Model evaluation metrics. A calibration curve plots the predicted versus actual outcomes, and a well-calibrated model (blue line) will be close to the red diagonal line (A). The Brier score is calculated as the squared difference between the actual and predicted values (B). The relationship between predictions and Brier score when the true class is 1 (C). The difference in BS between a null and full model shows the improvement in BS for most compounds (D).

prolongation from the Pollard et al. data and the blue circles are the predicted values<sup>33</sup>. The square of the length of the line connecting the two points is the BS for that compound.

Figure 4D shows the difference in BS for two models. The first is a “null” model that does not contain any predictor variables. The model can still achieve a 59% prediction accuracy because 59% of compounds do not increase the QT interval. The second model uses  $C_{\max}$  and hERG  $IC_{50}$  and is called the “full” model. Since the BS is expected to be higher for the null model (worse predictions), positive values for the difference  $BS_{\text{Null}} - BS_{\text{Full}}$  indicate compounds that were better predicted when using  $C_{\max}$  and hERG  $IC_{50}$ , compared with not using this information. The average BS for the null model is 0.24 and for the full model is 0.09, indicating a large improvement. Even though the prediction for three compounds is worse when including the  $C_{\max}$  and hERG data, the average performance is better. The average BS estimates the overall model fit, and the individual Brier scores indicate which compounds are hard to predict. These can be further investigated to understand why the model was unable to predict their outcome. For example, they may inhibit other cardiac ion channels.

Another way to express the Brier score is to scale it between 0 and 1 using the following equation<sup>35</sup>:

$$BS_{\text{scaled}} = 1 - \frac{BS}{BS_{\text{Max}}}$$

$BS_{\text{Max}}$  is the maximum possible BS calculated when we always predict the most frequent class, and averaged over all of the compounds (the null model). A scaled BS of 0 means we predict no better than always choosing the most frequent class, and a value of 1 means we always predict the true class with maximum probability. It can also be interpreted as an  $R^2$  value, or the proportion of variance in the outcome that is explained by the predictors<sup>36</sup>. For this example, the scaled Brier score is 0.62.

The commonly reported metrics of accuracy, sensitivity, and specificity have limitations and do not fully describe the performance of a prediction model. Hence, we suggest also examining discrimination, calibration, and overall model fit metrics when making decisions<sup>37</sup>.

## Not assessing the domain of applicability

Once a suitable prediction model has been built and is put into production, it can still perform poorly if the new compounds differ from those used to build the model. A final problem therefore is that the similarity between new (test) compounds and old (training) compounds is often not assessed; in other words, the model may be applied outside of its relevant domain. New compounds may differ structurally from the training compounds as new chemical space is explored, and this may translate into different physicochemical properties and different assay results. If new compounds differ by having values outside of the training data in multidimensional space, prediction becomes a form of extrapolation, which is a dangerous procedure. Another situation is when a new compound is located within the training data, but in a region with few observations. The model is interpolating (which is less dangerous than extrapolating), but in regions of sparse data, and so predictions may be worse than in

regions of dense data. Some models such as Gaussian Processes can account for both of these situations by allowing the uncertainty in the prediction to increase, but this is not possible for most standard machine learning models.

Assessing the similarity of a new compound to the old compounds should therefore be routine when making predictions. Many methods are available to detect if a data point is unlike others, known as outlier detection or anomaly detection<sup>38</sup>. To illustrate this point and one solution, the hERG IC<sub>50</sub> and C<sub>max</sub> data and the model from the previous section are used<sup>5</sup>. The idea is to characterise the hERG IC<sub>50</sub> and C<sub>max</sub> values with a mixture of Gaussian distributions, which is a form of clustering or unsupervised machine learning<sup>39</sup>. Then for a new compound, we calculate its distance to the nearest cluster. If this distance is large, we can flag the compound as an outlier and be more suspicious of the prediction. We use Mahalanobis distance, which takes the shape and orientation of the clusters into account when calculating the distance to the cluster centre.

Figure 5A plots the hERG IC<sub>50</sub> and C<sub>max</sub> training data (black circles) and the shaded blue regions show the three clusters that describe the data. The optimal number of clusters was determined by fitting 1-4 Gaussian distributions to the data and assessing their fit while penalising overly complex descriptions. We do not interpret the clusters and merely use them to describe the structure of the data. The four red diamonds labelled A–D are the hypothetical new compounds and were not used to form the clusters.

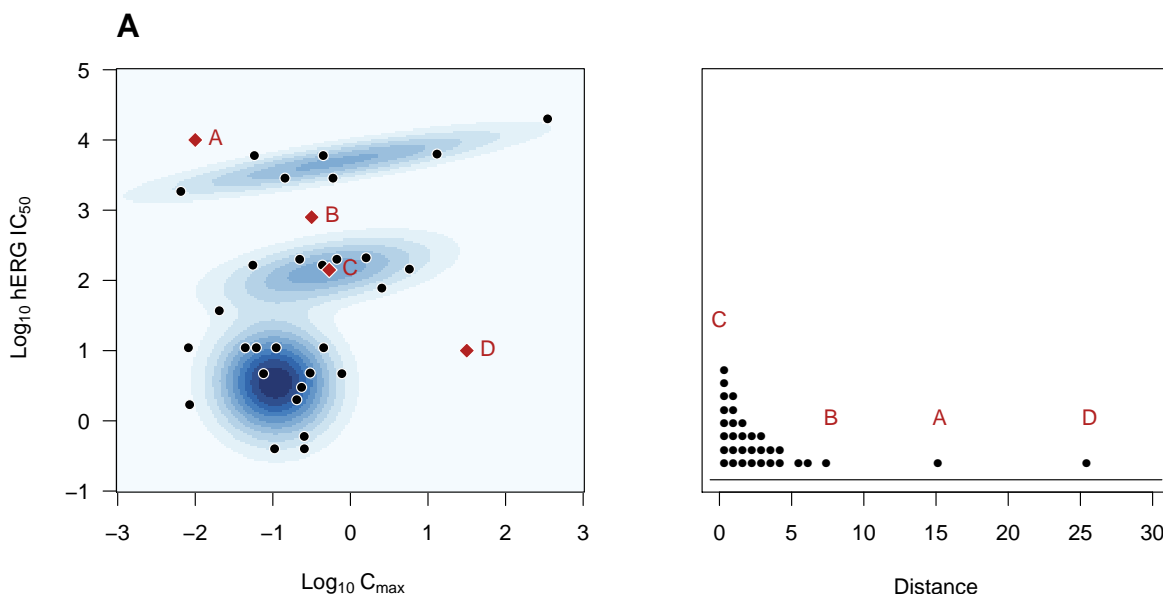


Figure 5: Outlier detection. The training data (black circles) are described by a mixture of three Gaussian distributions (shaded regions; A). The distance of four new data points (red diamonds, A–D) to the nearest cluster centre are calculated (B). Compounds with a large distance are unlikely to be from one of the clusters and their predictions may be unreliable.

Figure 5B plots the distance to the nearest cluster centre for the original and the four new compounds. Compound D is furthest from the original data in Figure 5A and has the largest

distance. Note that compound D does not have the lowest hERG IC<sub>50</sub> value nor the largest C<sub>max</sub> value, and so looking at one variable at a time would not indicate that compound D is unusual. Compound A is closer to the bulk of the data but still near the edge, and thus has a large distance. Compound B lies well within the interior of the data but in a region with few observations. It therefore has a larger distance. Compound C lies at the mean of the middle cluster and therefore has the smallest distance of all the compounds, both old and new (although this cannot be seen in Fig. 5B due to the small amount of binning required to stack the data points).

With only two predictor variables, visualising an unusual data point is simple, but it is much more difficult in higher dimensions. Formal methods to flag unusual compounds are therefore useful to highlight compounds for which predictions may be unreliable<sup>38,40,41</sup>.

## Discussion and recommendations

Why do poor practices persist? First, inertia: methods become established in a field or a company and they become the standard operating procedure, despite their shortcomings, and especially when a method's limitations are discovered years after adoption. Second, many of the above practices are simple to use and easy to understand. Many people who develop and use these predictive models may not have a statistics, data science, or machine learning background, and therefore rely on simpler but less suitable methods. Finally, there is little learning from failure at the institutional level in pharmaceutical companies, where these models are developed and used. A long lag exists between making a preclinical decision and learning the eventual clinical outcome. Often years pass, and with staff turnover, there may be few people around to reflect on the decision or method they used to make the prediction. By documenting the decision and reasons for it – ideally in a machine-readable format – this information can be fed back to improve decisions years later.

We discussed several common practices that can lead to problems when predicting toxicity. These problems happen at all stages of model development, including preprocessing (converting continuous variables into binary variables using arbitrary thresholds), feature engineering (using margins), model building (using scores, overfitting, ignoring uncertainty), assessment (evaluation metrics), and production (similarity of old and new compounds). This paper provided a high-level summary of current practices and the references cited throughout provide more theoretical background and technical details on implementing these methods. In addition, the data and R-code are provided as supplementary material to facilitate the uptake of these methods.

Many recommendations are simple to apply: avoid binning variables, use the measured continuous values instead; avoid using margins in predictive model, use IC<sub>50</sub> and C<sub>max</sub> values directly; and when evaluating the suitability of a predictive model, consider four types of metrics: classification, discrimination, calibration, and overall fit. If classes are unbalanced, prefer balanced accuracy to overall accuracy. Other recommendations will require a collaboration with quantitative researchers, such as avoid hand-crafted arbitrary scores and use statistical or machine learning models to predict outcomes. Similarly, assessing the

domain of applicability is more difficult as many options are available, but any method is better than none. Quantifying and propagating uncertainty is a large topic that was only briefly mentioned, but Bayesian methods are the standard way of handling uncertainty and are becoming more popular for predictive models in safety pharmacology<sup>23–28</sup>. Large and small pharmaceutical and biotechnology companies, contract research organisations, and other industries involved in drug development need to understand the limitations of using these assays out of context. A robust statistical framework is required to maximise predictivity, prevent errors of judgement, and significant cost from incorrect use.

## References

1. Loiodice, S., Nogueira da Costa, A. & Atienzar, F. Current trends in in silico, in vitro toxicology, and safety biomarkers in early drug development. *Drug Chem Toxicol* **42**, 113–121 (2019).
2. Gleeson, M. P. *et al.* The challenges involved in modeling toxicity data in silico: a review. *Curr. Pharm. Des.* **18**, 1266–1291 (2012).
3. Myatt, G. J. *et al.* In silico toxicology protocols. *Regul. Toxicol. Pharmacol.* **96**, 1–17 (2018).
4. Shah, F. *et al.* Setting clinical exposure levels of concern for drug-induced liver injury (DILI) using mechanistic in vitro assays. *Toxicol. Sci.* **147**, 500–514 (2015).
5. Lazic, S. E., Edmunds, N. & Pollard, C. E. Predicting drug safety and communicating risk: benefits of a Bayesian approach. *Toxicol. Sci.* **162**, 89–98 (2018).
6. Williams, D. P., Lazic, S. E., Foster, A. J., Semenova, E. & Morgan, P. Predicting drug-induced liver injury with Bayesian machine learning. *Chem. Res. Toxicol.* **33**, 239–248 (2020).
7. Cohen, J. The cost of dichotomization. *Applied Psychological Measurement* **7**, 249–253 (1983).
8. Maxwell, S. & Delaney, H. Bivariate median splits and spurious statistical significance. *Quantitative Methods in Psychology* **113**, 181–190 (1993).
9. Owen, S. V. & Froman, R. D. Why carve up your continuous data? *Res Nurs Health* **28**, 496–503 (2005).
10. Wainer, H., Gessaroli, M. & Verdi, M. Finding what is not there through the unfortunate binning of results: The Mendel effect. *CHANCE* **19**, 49–52 (2006).
11. Royston, P., Altman, D. G. & Sauerbrei, W. Dichotomizing continuous predictors in multiple regression: A bad idea. *Stat Med* **25**, 127–141 (2006).
12. Walraven, C. van & Hart, R. G. Leave 'em alone - why continuous variables should be analyzed as such. *Neuroepidemiology* **30**, 138–139 (2008).
13. Fedorov, V., Mannino, F. & Zhang, R. Consequences of dichotomization. *Pharm Stat* **8**,

50–61 (2009).

14. Irwin, J. & McClelland, G. Negative consequences of dichotomizing continuous predictor variables. *Journal of Marketing Research* **40**, 366–371 (2003).

15. Kenny, P. W. & Montanari, C. A. Inflation of correlation in the pursuit of drug-likeness. *J Comput Aided Mol Des* **27**, 1–13 (2013).

16. Kuss, O. The danger of dichotomizing continuous variables: A visualization. *Teaching Statistics* **35**, 78–79 (2013).

17. Lazic, S. E. Four simple ways to increase power without increasing the sample size. *Lab. Anim.* **52**, 621–629 (2018).

18. Thompson, R. A. *et al.* In vitro approach to assess the potential for risk of idiosyncratic adverse reactions caused by candidate drugs. *Chem. Res. Toxicol.* **25**, 1616–1632 (2012).

19. Aleo, M. D. *et al.* Moving beyond Binary Predictions of Human Drug-Induced Liver Injury (DILI) toward Contrasting Relative Risk Potential. *Chem. Res. Toxicol.* **33**, 223–238 (2020).

20. Steyerberg, E. W. *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating.* (Springer, 2019).

21. Burnham, K. P. & Anderson, D. R. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach.* (Springer, 2002).

22. Yao, Y., Vehtari, A., Simpson, D. & Gelman, A. Using stacking to average Bayesian predictive distributions (with discussion). *Bayesian Anal.* **13**, 917–1007 (2018).

23. Obrezanova, O. & Segall, M. D. Gaussian processes for classification: QSAR modeling of ADMET and target activity. *J Chem Inf Model* **50**, 1053–1061 (2010).

24. Johnstone, R. H., Bardenet, R., Gavaghan, D. J. & Mirams, G. R. Hierarchical Bayesian inference for ion channel screening dose-response data. *Wellcome Open Res* **1**, 6 (2016).

25. Fronczyk, K. & Kottas, A. Risk assessment for toxicity experiments with discrete and continuous outcomes: A Bayesian nonparametric approach. *Journal of Agricultural, Biological and Environmental Statistics volume* **22**, 585–601 (2017).

26. Feng, D., Svetnik, V., Liaw, A., Pratola, M. & Sheridan, R. P. Building Quantitative Structure-Activity Relationship Models Using Bayesian Additive Regression Trees. *J Chem Inf Model* **59**, 2642–2655 (2019).

27. Hatherell, S. *et al.* Identifying and characterizing stress pathways of concern for consumer safety in next-generation risk assessment. *Toxicol. Sci.* **176**, 11–33 (2020).

28. Semenova, E., Williams, D. P., Afzal, A. M. & Lazic, S. E. A Bayesian neural network for toxicity prediction. *Computational Toxicology (in press)*, (2020).

29. Barber, D. *Bayesian Reasoning and Machine Learning.* (Cambridge University Press, 2012).

30. Gneiting, T. & Raftery, A. E. Strictly proper scoring rules, prediction, and estimation.



*Journal of the American Statistical Association* **102**, 359–378 (2007).

31. Harrell, F. E. *Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression, and Survival Analysis*. (Springer, 2015).

32. Harrell, F. E., Califf, R. M., Pryor, D. B., Lee, K. L. & Rosati, R. A. Evaluating the yield of medical tests. *JAMA* **247**, 2543–2546 (1982).

33. Pollard, C. E. *et al.* An analysis of the relationship between preclinical and clinical QT interval-related data. *Toxicol. Sci.* **159**, 94–101 (2017).

34. Brier, G. W. Verification of forecasts expressed in terms of probability. *Month. Weather Rev* **78**, 1–3 (1950).

35. Steyerberg, E. W. *et al.* Assessing the performance of prediction models: a framework for traditional and novel measures. *Epidemiology* **21**, 128–138 (2010).

36. Hu, B., Palta, M. & Shao, J. Properties of R(2) statistics for logistic regression. *Stat Med* **25**, 1383–1395 (2006).

37. Moons, K. G. *et al.* Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): explanation and elaboration. *Ann. Intern. Med.* **162**, 1–73 (2015).

38. Aggarwal, C. C. *Outlier Analysis*. (Springer, 2017).

39. Bouveyron, C., Celeux, G., Murphy, T. B. & Raftery, A. E. *Model-Based Clustering and Classification for Data Science: With Applications in R*. (Cambridge University Press, 2019).

40. Netzeva, T. I. *et al.* Current status of methods for defining the applicability domain of (quantitative) structure-activity relationships. The report and recommendations of ECVAM Workshop 52. *Altern Lab Anim* **33**, 155–173 (2005).

41. Carrio, P., Pinto, M., Ecker, G., Sanz, F. & Pastor, M. Applicability Domain ANALysis (ADAN): a robust method for assessing the reliability of drug property predictions. *J Chem Inf Model* **54**, 1500–1511 (2014).