

# Quantifying sources of uncertainty in drug discovery predictions with probabilistic models

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)

## 1 Abstract

2 Knowing the uncertainty in a prediction is critical when making expensive investment de-  
3 cisions and when patient safety is paramount, but machine learning (ML) models in drug  
4 discovery typically provide only a single best estimate and ignore all sources of uncer-  
5 tainty. Predictions from these models may therefore be over-confident, which can put  
6 patients at risk and waste resources when compounds that are destined to fail are fur-  
7 ther developed. Probabilistic predictive models (PPMs) can incorporate uncertainty in  
8 both the data and model, and return a distribution of predicted values that represents  
9 the uncertainty in the prediction. PPMs not only let users know when predictions are  
10 uncertain, but the intuitive output from these models makes communicating risk eas-  
11 ier and decision making better. Many popular machine learning methods have a PPM  
12 or Bayesian analogue, making PPMs easy to fit into current workflows. We use toxicity  
13 prediction as a running example, but the same principles apply for all prediction mod-  
14 els used in drug discovery. The consequences of ignoring uncertainty and how PPMs ac-  
15 count for uncertainty are also described. We aim to make the discussion accessible to a  
16 broad non-mathematical audience, but also provide code for computational researchers  
17 ([https://github.com/stanlazic/ML\\_uncertainty\\_quantification](https://github.com/stanlazic/ML_uncertainty_quantification)).

## 18 Introduction

19 At each stage of the drug discovery pipeline, researchers decide to progress or halt com-  
20 pounds using both qualitative judgements and quantitative methods. This is formally a  
21 prediction problem, where, given some information, a prediction is made about a future  
22 observable outcome. Standard predictive or machine learning models such as random  
23 forests, support vector machines, or neural networks only report point-estimates, or a sin-  
24 gle “best” value for a prediction; they provide no information on the uncertainty of the  
25 prediction. Prediction uncertainty is important when (1) the range of plausible values is  
26 as important as the best estimate, (2) you need to know that the model cannot confidently  
27 make a prediction, (3) you need to reliability distinguish between ranked items, or (4)  
28 the cost of an incorrect decision is large; for example, when making expensive investment  
29 decisions or when assessing patient safety.

30 Figure 1 shows predicted clinical blood aspartate transaminase (AST) levels – an in-  
31 dicator of liver toxicity – for two hypothetical compounds. Assume that levels below 8  
32 (arbitrary units) are considered safe, and that only one compound can be taken forward  
33 for clinical trials. Based only on the best estimate, compound A (blue) is preferable (Fig.  
34 1A). Knowing the prediction uncertainty changes the picture (Fig. 1B). 14% of compound  
35 A’s distribution is in the unsafe shaded region, while only 2% of compound B’s distribution  
36 is in the shaded region. Based on these distributions, compound B maybe the better can-  
37 didate to progress to clinical trials; or a project team may decide to run more experiments  
38 to reduce compound A’s uncertainty.

39 We define a probabilistic predictive model (PPM) as any machine learning model that  
40 returns a distribution for the prediction instead of a single value. PPMs differ in how they  
41 represent uncertainty. At one end, fully probabilistic Bayesian models specify a distribu-  
42 tion for the outcome and for all unknown parameters in the model. These models are the  
43 gold-standard for quantifying uncertainty but they can be computationally expensive. At  
44 the other end are method that return multiple predicted values without specifying a prob-  
45 ability distribution. Examples include fitting the same model on multiple bootstrapped  
46 datasets, or for models with a stochastic component, fitting the same model using different  
47 random number generator seeds [31]. Although these models are often computationally  
48 tractable, the connection between the distribution of predicted values and uncertainty is  
49 unclear. For example, does varying something trivial such as the random number genera-  
50 tor seed adequately capture our uncertainty in a prediction? Between these extremes are  
51 approaches that try to obtain the benefits of fully probabilistic models using approxima-  
52 tions, reformulations, or computational shortcuts. For example, instead of using a fully  
53 Bayesian deep neural network, Kristiadi et al. were able to obtain many of the benefits by  
54 making only the final layer of the network Bayesian [29]. Also in the neural network liter-

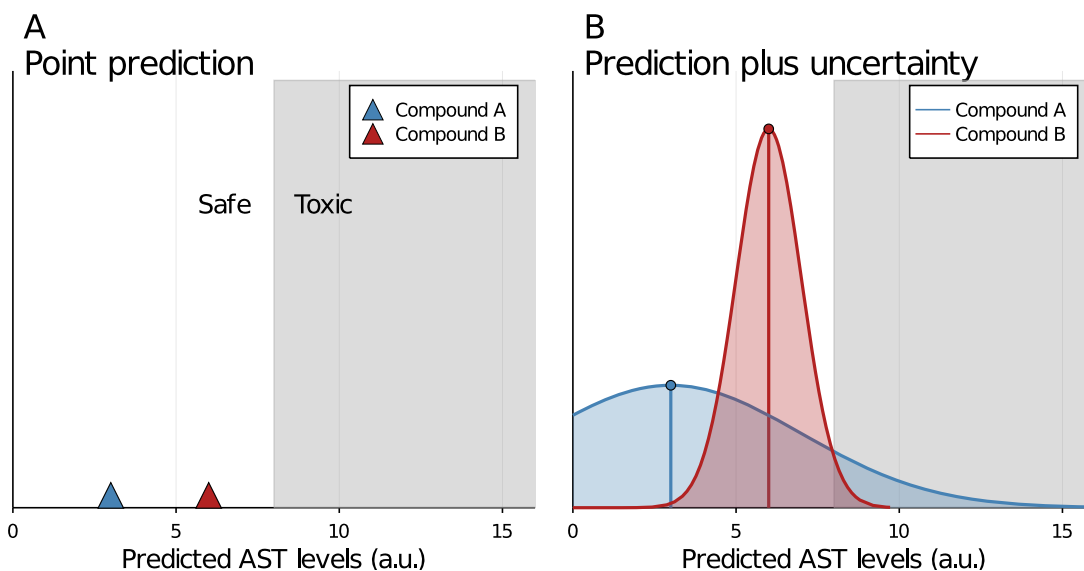


Figure 1: Prediction uncertainty. Predicted blood AST levels for two compounds, with Compound A appearing better (A). Values above the safety threshold of 8 (grey shaded region) will not be progressed. Reporting prediction uncertainty shows that 18% of Compound A is above the threshold but only 2% of Compound B, indicating that Compound B is better (B). AST = Aspartate transaminase; a.u. = arbitrary units.

ature, Gal and Ghahramani approximated parameter uncertainty by randomly inactivating nodes at prediction time – a procedure known as Monte Carlo dropout [15], and Teye and colleagues used the mean and variances calculated from batch normalisation steps for the same purpose [63]. These approximate methods are an active area of research [3, 27, 46, 70] (see Mervin et al., for a recent review [40]), and will be critical for making PPMs more widely adopted, but here we focus on fully probabilistic models as they better highlight the key areas of uncertainty.

All prediction models can be written as

$$y|x \quad (1)$$

and read as “ $y$  given  $x$ ” – where  $y$  is the outcome to be predicted and  $x$  are one or more variables used to predict  $y$ . Both  $x$  and  $y$  are available when training a model, and when a model is deployed, new  $x$  values are observed and used to predict the unknown  $y$ ’s. PPMs additionally provide a probability distribution for  $y$ , denoted as

$$P(y | x) \tag{2}$$

and read as “the probability of  $y$  given  $x$ ”. A distribution for  $y$  enables us to calculate any metric of interest, such as the best guess for  $y$  (e.g. mean, median, or mode), thereby providing the same information as standard methods. But in addition, prediction intervals (PI) can be calculated around the best estimate, or the probability that  $y$  is greater or less than a predefined threshold can be calculated, as was done in Figure 1B.

Below we describe the sources of uncertainty and the advantages of PPMs. We aim to make the discussion accessible to a broad non-mathematical audience. Equations are provided to make ideas concrete for mathematical readers, but can be skipped without loss of understanding of the remaining text. In addition, code is provided for computational researchers ([https://github.com/stanlazic/ML\\_uncertainty\\_quantification](https://github.com/stanlazic/ML_uncertainty_quantification)) and implemented in Julia using Turing [1, 16].

## Sources of uncertainty

The definition of a probabilistic model in Equation 2 lacked a crucial component, which is the model itself:

$$P(y | x, \text{Model}). \tag{3}$$

The sources of uncertainty are now clear: they can reside in the data ( $x, y$ ) or in the model. A prediction for  $y$  not only depends on  $x$ , but on the model used to connect  $x$  and  $y$ . Hence, predictions are conditional on a model, and uncertainty in the model should lead to greater uncertainty in the prediction. Models are composed of the following components, which we discuss in greater detail below:

1. **Data.** The outcome  $y$  and the predictor or input variables  $x$ .
2. **Distribution function.** The distribution that represents our uncertainty in  $y$ , also called the likelihood or data generating distribution:  $G(\cdot)$ .
3. **Mean function.** The functional or structural form of the model describing how  $y$  changes as  $x$  changes:  $f_\mu(\cdot)$ .

- 91 4. **Variance function.** Describes how the uncertainty in  $y$  varies with  $x$ :  $f_\sigma(\cdot)$ .
- 92 5. **Parameters.** The unknown coefficients or weights for the mean ( $\theta_\mu$ ) and variance
- 93 ( $\theta_\sigma$ ) functions that are estimated from the data.
- 94 6. **Hyperparameters.** Parameters or other options used to define a model that are not
- 95 estimated from the data but fixed or selected by the analyst:  $\phi$ .
- 96 7. **Link functions.** Nonlinear transformations of the mean ( $l_\mu(\cdot)$ ) and/or variance func-
- 97 tion ( $l_\sigma(\cdot)$ ) used to keep values within an allowable range.

98 Combining these seven components gives the generic formulation of a PPM (Eq. 4).  
 99 Although this equation is abstract, it captures where uncertainty can reside. In this section  
 100 we carefully describe the terms in the equation and then provide a concrete example.

$$\begin{aligned}
 y &\sim G(\mu, \sigma) \\
 \mu &= l_\mu(f_\mu(x; \theta_\mu)) \\
 \sigma &= l_\sigma(f_\sigma(x; \theta_\sigma))
 \end{aligned}
 \tag{4}$$

101 Starting with the first line of Equation 4,  $y$  is a future value that we want to predict  
 102 and it could represent a clinical outcome, an IC<sub>50</sub> value, or a physicochemical property  
 103 of a compound such as solubility. PPMs require that we specify a distribution for our  
 104 uncertainty in  $y$ , which we can informally think of as the distribution from which  $y$  was  
 105 generated. We denote this distribution as  $G(\cdot)$  and the “ $(\cdot)$ ” notation indicates that  $G$  is a  
 106 function with inputs ( $\mu$  and  $\sigma$  in this case), but places the focus on  $G$  and not the inputs,  
 107 thereby reducing clutter. Common distributions include the normal/Gaussian (for continu-  
 108 ous symmetric data), Student-t (continuous symmetric data with outliers), Bernoulli (0/1  
 109 data), and Poisson (count data). The mean of the chosen distribution is given by  $\mu$  and  
 110 many distributions have a second parameter,  $\sigma$ , that controls the spread or width of the  
 111 distribution. This parameter is critical for PPMs because it describes the uncertainty in  
 112  $y$ . The  $\sim$  symbol can be read as “is distributed as” or “is generated from”. To make this  
 113 concrete, we might represent our uncertainty in  $y$  for a given compound with a Gaussian  
 114 distribution that has a mean  $\mu = 2.45$  – which represents our best estimate of  $y$  – and a  
 115 standard deviation of  $\sigma = 2.1$ . This would be written as  $y \sim \text{Normal}(2.45, 2.1)$ .

116 But where did our best estimate  $\mu$  come from? This is defined on the second line of  
 117 Equation 4.  $x$  is the input data used to predict  $y$  and it could represent the compound  
 118 structure (encoded as a binary fingerprint for example), assay results, or physicochemical

properties. The prediction task reduces to using  $x$  to predict  $y$ , but they need to be connected through a statistical or machine learning model. The structural or functional form of this model is denoted by  $f_\mu(\cdot)$  and it could represent the structure of a simple linear regression model, the architecture of a neural network, an ensemble of trees, or a differential equation representing a pharmacokinetic model. We refer to  $f_\mu(\cdot)$  as the *mean function* because it tells us the predicted mean value of  $y$  for a given value of  $x$ . The mean function contains parameters ( $\theta_\mu$ ) that are estimated from the data, and the parameters could represent the coefficients of a linear model or the weights of a neural network. The “learning” in machine learning refers to estimating values for the parameters that maximises predictive performance, given the data and functional form of the model.  $\theta_\mu$  usually represents multiple parameters in the mean function; for example, it would represent both the intercept and slope in a simple regression model. The subscript  $\mu$  on  $f$  and  $\theta$  indicates that the function and parameters refer to the mean, since we also have a function and parameters for the variance,  $\sigma$ , which we described further below.

One problem is that  $\mu$  has no constraints, and can be any value calculated from  $f_\mu(\cdot)$ , which may lead to impossible predictions. For example, if we’re predicting the probability that a compound is toxic,  $f_\mu(\cdot)$  needs to be between zero and one – values outside this range do not make sense. Hence, we need to transform  $f_\mu(\cdot)$  with a *link function*  $l_\mu(\cdot)$  to put it within a permissible range. One option is to use the equation  $1/(1 + \exp(-f_\mu(\cdot)))$  as a link function, which compresses  $f_\mu(\cdot)$  into the 0-1 range, but other functions are also possible. A link function is unnecessary when no restriction on  $f_\mu(\cdot)$  is required, and  $l_\mu(\cdot)$  can be dropped from the formula. (Useful mnemonics:  $f$  = Function;  $G$  = data Generating distribution;  $l$  = Link function; with subscripts  $\mu$  and  $\sigma$  referring to the mean and variance, respectively).

The third line of Equation 4 shows that the variance or uncertainty in a predicted value of  $y$  can be specified in the same way as we specify the mean or best estimate of  $y$ . Many models assume a constant value for  $\sigma$ , which is the “homogeneity of variance” assumption in traditional statistical models. However, a model’s uncertainty in  $y$  may vary for different values of  $x$ , which can be captured with a *variance function*  $f_\sigma(\cdot)$ . Since variances are positive, a link function for the variance ( $l_\sigma(\cdot)$ ) is needed to constrain  $\sigma$  to be greater than zero.

Fully Bayesian PPMs also require us to specify the uncertainty in the parameters  $\theta_\mu$  and  $\theta_\sigma$  before analysing the data, and the hyperparameters ( $\phi$ ) refer to this specification. Many of the approximate methods avoid this step and  $\phi$  may not correspond to any part of the model but defines other options for the learning algorithm and therefore it is not included in Equation 4.

Usually the  $x$  and  $y$  data are all we have, and we need to choose  $G$ ,  $l_\mu$ ,  $f_\mu$ ,  $l_\sigma$ ,  $f_\sigma$ , and

156  $\phi$  based on background knowledge, preliminary plots of the data, or trying several options  
 157 and empirically assessing which is best. For example, an initial model might assume that  
 158  $G$  is Gaussian and that  $f_\mu$  follows a hypothesised mechanistic relationship based on a  
 159 pharmacokinetic model. Then, we estimate or learn the values of  $\theta_\mu$  based on training  
 160 data, and assess the prediction on a separate test data set. To make these ideas concrete,  
 161 Figure 2 shows simulated data for 100 compounds, where  $y$  is a clinical outcome and  $x$  is  
 162 an assay result. Assume that higher values of  $y$  indicate greater toxicity.

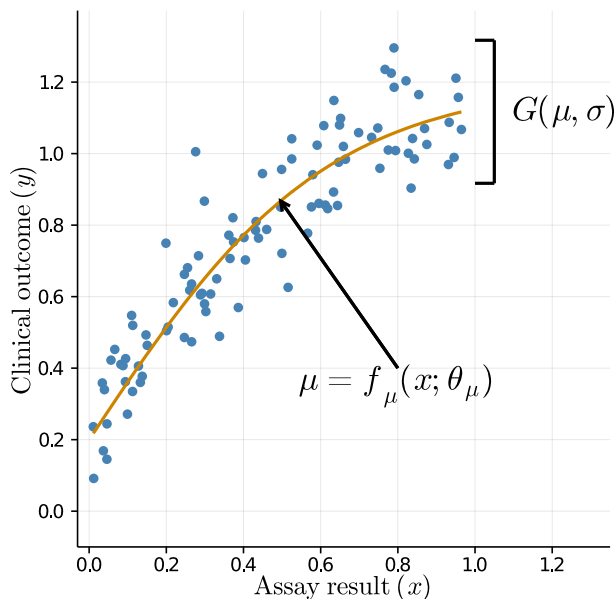


Figure 2: Simulated data with the true relationship between  $x$  and  $y$  given by  $\mu = f_\mu(x; \theta_\mu)$  (orange curve) and based on Equation 5.  $G$  is the Gaussian data generating distribution with a mean  $\mu$  and a constant variance  $\sigma$ , which models the spread of points around the line. Link functions for  $\mu$  and  $\sigma$  are not used and hence are not shown.

163 A nice feature of PPMs is that they are generative, meaning that they can generate  
 164 or simulate data. Indeed, simulation and learning are opposite sides of the same coin:  
 165 learning takes the fixed data and infers likely values of the parameters that could have  
 166 generated the data, whereas simulation fixes the parameters and generates the data. The  
 167 model in Equation 5 generated the data in Figure 2 and we will use it as a running example  
 168 throughout

$$y \sim \text{Normal}(\mu, \sigma) \tag{5}$$

$$\mu = \theta_2 + \frac{1 - e^{-\theta_1 x}}{1 + e^{-\theta_1 x}}.$$

169 The first line of the equation is read as: “the outcome  $y$  is generated ( $\sim$ ) from a Gaus-  
 170 sian or Normal distribution with a mean of  $\mu$  and a standard deviation of  $\sigma$ ”. But how  
 171 does  $y$  depend on  $x$ ? The second line shows how  $x$  enters and how it depends on two  
 172 parameters:  $\theta_1$  and  $\theta_2$  ( $e$  is a constant, not a parameter). We set this equation equal to  $\mu$   
 173 and can substitute it for  $\mu$  in the first line of Equation 5 giving

$$y \sim \text{Normal}\left(\theta_2 + \frac{1 - e^{-\theta_1 x}}{1 + e^{-\theta_1 x}}, \sigma\right).$$

174 Writing the equation in one line makes the relationship between  $x$  and  $y$  clearer, but  
 175 multi-lined equations are easier to read with more complex models. For most prediction  
 176 models the parameters are uninteresting, but to help interpret the model,  $\theta_2$  is the  $y$ -  
 177 intercept (value of  $y$  when  $x = 0$ ), and  $\theta_1$  controls how quickly the line in Figure 2 reaches  
 178 the upper limit of  $y = 1.2$  as  $x$  gets large.

179 To simulate a value for  $y$ , we need to (1) select parameter values, and we use the  
 180 following:  $\theta_1 = 3.25$ ,  $\theta_2 = 0.2$ , and  $\sigma = 0.1$ ; (2) select a value of  $x$ , which enables us to  
 181 calculate  $\mu$ ; then (3) draw a random number from a Gaussian distribution with a mean of  
 182  $\mu$  and standard deviation of  $\sigma$ . This can be repeated any number of times to obtain the  
 183  $P(y|x)$  distribution, and for different values of  $x$ . The data in Figure 2 were generated  
 184 for 100  $x$  values uniformly distributed between 0 and 1. Note that  $\theta_\mu$  in Equation 4 is a  
 185 place-holder for several variables, which correspond to  $\theta_1$  and  $\theta_2$  in Equation 5. Given this  
 186 data, we now illustrate where the seven sources of uncertainty enter.

## 187 Mean function uncertainty

188 Uncertainty in the mean function  $f_\mu(\cdot)$  arises because we rarely know the true form of  
 189 the relationship between  $x$  (assay) and  $y$  (outcome) – that is, we don’t know the form of  
 190 Equation 5. Uncertainty in the mean function is also called model uncertainty, but this term  
 191 is ambiguous because models have multiple components. Choices for the mean function  
 192 include which predictors, interaction terms, transformations, basis expansions, hierarchies,



193 and time-varying components to include in the model. Assume we only observe the data  
 194 in Figure 2, several models we might consider for the relationship between  $x$  and  $y$  are:

$$\begin{aligned}
 \text{Linear :} & \quad y = \theta_0 + \theta_1 x \\
 \text{Quadratic :} & \quad y = \theta_0 + \theta_1 x + \theta_2 x^2 \\
 \text{2 Parameter Exponential :} & \quad y = \theta_2 (1 - e^{-\theta_1 x}) \\
 \text{3 Parameter Exponential :} & \quad y = \theta_3 + \theta_2 (1 - e^{-\theta_1 x}) \\
 \text{Michaelis - Menten :} & \quad y = \theta_1 x / (\theta_2 + x).
 \end{aligned}$$

195 None of these are the true model (Eq. 5), but except for the linear model, they all  
 196 saturate at high values of  $x$  or are concave and therefore capture the main trend in the  
 197 data. Which model should we use? Typically, only a single model is selected and predic-  
 198 tions are made from that. If one model is clearly better than the others, there may be  
 199 little lost by using one model for predictions. However, if two or more models fit the data  
 200 equally well, making predictions from only one will underestimate the prediction uncer-  
 201 tainty. Fortunately, we are not forced to choose one model but can fit several and combine  
 202 their predictions. To illustrate, we will use the quadratic, 2-parameter exponential, and  
 203 3-parameter exponential models. The three models are fit to the data (Fig. 3A–C) and  
 204 predictions are extrapolated to show both how similar the fits are where there is data, and  
 205 how different the fits are when extrapolating. The shaded regions show the 95% prediction  
 206 intervals (PI), and if a model is suitable, we expect 95% of the data to fall in the shaded  
 207 region.

208 To better compare the predictions, the mean functions  $\mu$  for three models are plot-  
 209 ted together in Figure 3D. The models make similar predictions within the range of the  
 210 data, except at very low values of  $x$ , where the 2-parameter exponential model predicts  
 211 smaller values. Figure 3E shows the model-averaged prediction, and note how uncertainty  
 212 is greater at high values of  $x$ , where the models make different predictions. To better ap-  
 213 preciate how model averaging incorporates uncertainty, Figure 3F shows the width of the  
 214 95% PI for the three models and the averaged model. Note how the averaged model has a  
 215 wider PI at low values of  $x$  than any of the original models. This occurs because the three  
 216 models make different predictions at low values and hence this region is more uncertain.  
 217 For intermediate values of  $x$ , all three models make similar predictions and the averaged  
 218 PI is wider than some models and lower than others. When extrapolating to larger values  
 219  $x$ , the model averaged PI width quickly becomes the widest, reflecting both the diverging  
 220 predictions of the individual models and the greater uncertainty in the quadratic model.

221 Predictions are always conditional on a model, and if we don't know the true model,

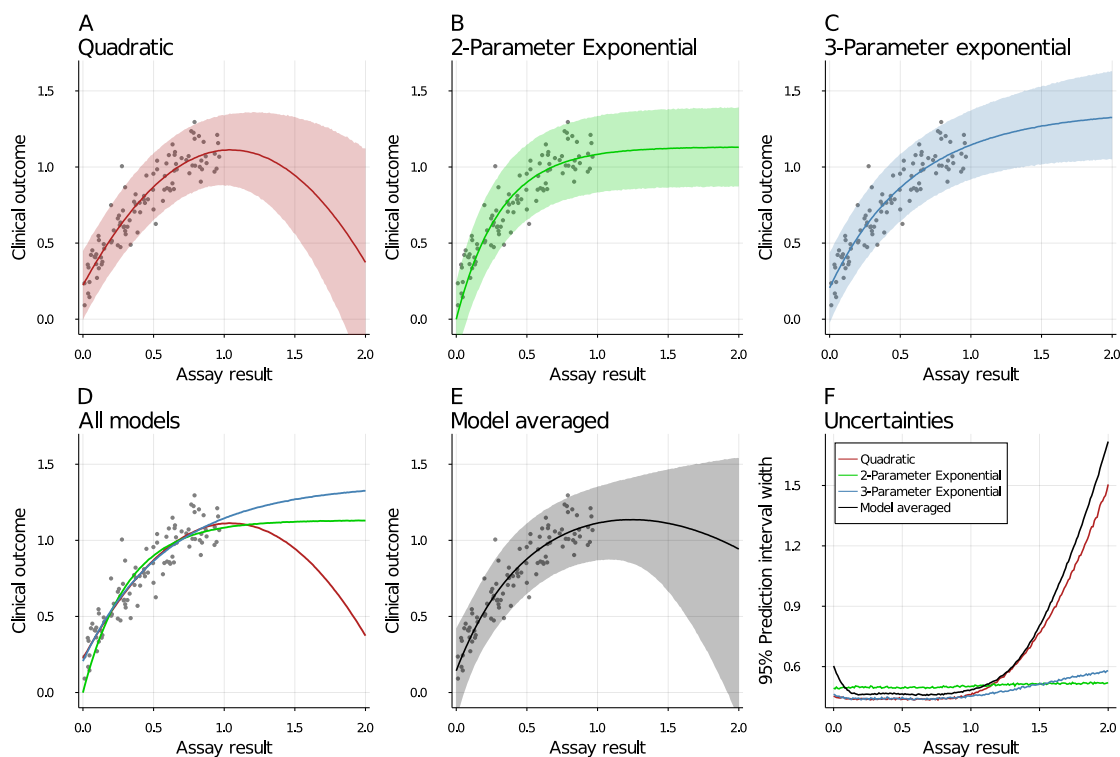


Figure 3: Model averaging. Three models fit the data well (A-C), even though none are the true model. They make different predictions at low assay values and when extrapolating to higher values. The mean predictions are superimposed for easier comparison (D). Model averaged prediction (E). Comparison of prediction interval widths (F). Shaded regions are the 95% prediction intervals.

predictions from the wrong model are likely to be overconfident. We extrapolated well beyond the data to illustrate how model averaging accounts for model uncertainty. Such extrapolation may seem unrealistic, but the message is that whenever models make different predictions for the same values of  $x$ , the uncertainty in the model-specific predictions will be overconfident, as we see to a lesser degree for low values of the assay.

## Parameter uncertainty

Most predictive models have parameters or weights that are learned from the data and control how the predictor variables ( $x$ ) are related to the outcome variable ( $y$ ) – these parameters are the  $\theta$  symbols in the five models considered above. Most machine learning

231 methods only use the single best value of each parameter when making a prediction. But  
 232 since the parameters are learned from the data, they are uncertain, and this uncertainty  
 233 should be propagated into the prediction. Parameter uncertainty decreases as the sample  
 234 sizes increases, so to better illustrate the effect of parameter uncertainty on predictions, a  
 235 smaller dataset was made by taking every eighth data point from the previous example.

236 Figure 4 uses the quadratic model, which has four parameters  $(\theta_0, \theta_1, \theta_2, \sigma)$ . The grey  
 237 shaded region shows the 95% prediction interval from a Bayesian model that accounts for  
 238 parameter uncertainty. The dashed black lines show the 95% PI from a classic quadratic  
 239 regression model which ignores parameter uncertainty, and note how they are slightly  
 240 narrower. The mean function is identical for both the Bayesian and classic model.

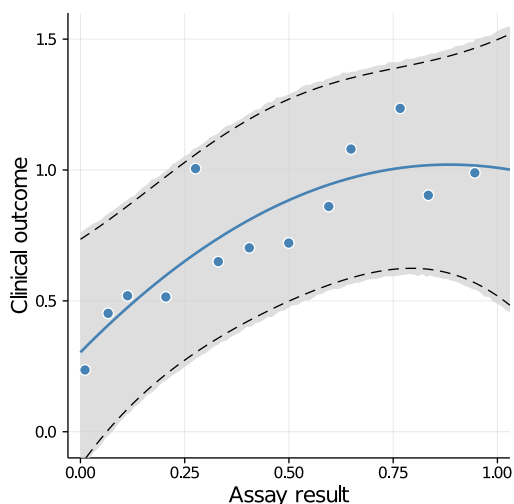


Figure 4: Parameter uncertainty. The grey shaded region is the 95% prediction interval for a Bayesian model and the narrower black dashed lines for a classical model. The classic interval is narrower because it doesn't account for uncertainty in the parameters.

241 The difference in PI width may seem negligible when focusing on the mean prediction,  
 242 but ignoring parameter uncertainty gives approximately 7% narrower PIs. This may be  
 243 important with “point of departure” calculations when the tails of the distributions are  
 244 more important than the means [50]. Assume clinical outcomes above 1.2 are considered  
 245 problematic and the assay result for one compound is  $x = 0.5$ . When fully accounting for  
 246 parameter uncertainty, there is a 5.9% chance that the true value of the clinical outcome is  
 247 above 1.2, versus a 4.3% chance when ignoring uncertainty. The ratio of these numbers is  
 248 1.37, indicating that the tail area is nearly 1.4 times larger when accounting for uncertainty.

249 Models typically have many more parameters than this example (the state-of-the-art  
250 Generative Pre-trained Transformer 3 (GPT-3) deep learning language model has 175 bil-  
251 lion parameters [6]) and simply collecting more data is often not an option to reduce  
252 parameter uncertainty because more data enables more complex models to be fit (e.g. in-  
253 cluding nonlinear terms and interactions), which then increases the number of parameters.

## 254 Hyperparameter uncertainty

255 Hyperparameters are a diverse set of tunable options that affect the training and predic-  
256 tions. Unlike parameters, they are not estimated from the data but selected by the analyst;  
257 examples include the amount of regularisation in a lasso model, the number of trees in a  
258 random forest model, or the cost function in a support vector machine. Suitable values are  
259 typically found by trying several options and selecting the best using crossvalidation. In  
260 the hyperparameter category we can also include options that are rarely part of a formal  
261 selection process such as the choice of optimisation algorithm or random number seed for  
262 models with a stochastic component. For fully Bayesian models we can also include param-  
263 eters for prior distributions, which are not updated by the data. These hyperparameters  
264 are selected pragmatically to provide good predictions, but other sets of hyperparameter  
265 values might give equally good predictions, on average, but slightly different predictions  
266 for each test compound. Hence, uncertainty in hyperparameter values is rarely taken into  
267 account. A further complication is that most hyperparameters are not related to any bi-  
268 ological or chemical quantity of interest and hence it is unclear what the uncertainty is  
269 actually about. Nevertheless, Lakshminarayanan and colleagues showed that by running  
270 many models with a different random seed, the ensemble of predictions performed bet-  
271 ter than a single model, and the distribution of predicted values provided a measure of  
272 uncertainty [31].

## 273 Data uncertainty

274 One of the main sources of uncertainty – and which is almost universally ignored – is the  
275 uncertainty in the data, both in the predictors ( $x$ ) and in the outcome to be predicted  
276 ( $y$ ). The uncertainty can be in the training data used to build the model, in the new data  
277 to be predicted, or both. The main sources of data uncertainty are measurement error,  
278 misclassification error, binning, censoring and truncation, and missing values data. Each  
279 of these are discussed below.

280 Predictors are often experimental measurements and are therefore subject to measure-

281 ment error, or they are samples from a larger population and are therefore subject to sam-  
282 pling error (e.g. only cells in the field of view are measured, not all cells in a well, and if a  
283 different subset of cells were selected, a different measured value would be obtained). Pre-  
284 dictors may also be calculated quantities such as  $IC_{50}$  values estimated from dose-response  
285 or concentration-response curves, and hence are uncertain. Furthermore, some predictors  
286 such as cLogP are the output of other (imperfect) prediction models and therefore are also  
287 uncertain. Finally, some predictors are not measured directly but are estimated from a  
288 standard curve, which introduces additional uncertainty because the curves may not be  
289 not perfectly calibrated.

290 All these are examples of classic measurement error, defined as  $x_{\text{measured}} = x_{\text{true}} + \text{error}$ ,  
291 where the measured  $x$  value is the true value corrupted by some error or noise. Errors can  
292 also be multiplicative, where  $x_{\text{measured}} = x_{\text{true}} \times \text{error}$ . Another type of error is Berk-  
293 son error, where samples or experimental units are assumed to have the same exposure  
294 but actually differ. For example, several wells in a microtitre plate are all given the same  
295 concentration of a compound, but in practice the true concentration may differ due to vari-  
296 ations in the amount of compound dispensed or if wells closer to the edge of a plate have  
297 evaporated and thus have a higher effective concentration. Berkson error often results  
298 from converting continuous values into bins or groups. For example, compounds are cate-  
299 gorised as active versus inactive, despite having a range of activity values. Or, compounds  
300 are classified as having no, mild, or severe toxicity, even though compounds will have a  
301 range of toxicity levels *within* each category. Binning can also lead to misclassification  
302 error, where a compound is placed into the incorrect category. This can occur if the mea-  
303 sured assay value differed from the true value and fell on the opposite side of a threshold.  
304 Hence, binning is strongly discouraged [33, 35]. Misclassification can also occur due to  
305 incorrect diagnoses, labelling errors, or data-entry errors. The standard response to these  
306 known and often large sources of data uncertainty is to ignore it, effectively assuming that  
307  $x_{\text{measured}} = x_{\text{true}}$ .

308 It is well known that ignoring error in  $x$  can bias parameter estimates, but for predic-  
309 tion models the parameters are usually not of interest. Even though noisy data can lead to  
310 biased parameter estimates, the model is still consistent for the prediction, meaning that  
311 as the sample size increases, the prediction will be correct, on average [9, 20]. This likely  
312 explains why error in  $x$  has received little attention in the predictive modelling and ma-  
313 chine learning literature. However, if we're interested in the uncertainty in the prediction,  
314 then making good predictions on average is not good enough, we need to ensure that the  
315 prediction uncertainty is calibrated.

316 Data are *censored* when they are known only up to a boundary value, but not beyond,  
317 and therefore only partial information is available. For example, assays typically have  
318 upper and lower limits of detection (LoD) and uncertainty arises because the exact value

319 is unknown, but it is typically treated as the “true” measured value.

320 Data are *truncated* when values outside of a range are omitted, and the number of  
321 omitted values is unknown. For example, for objects to be segmented as a cell in a standard  
322 image analyses, they must have a minimum cell size. Smaller cells will not be included in  
323 the analysis, and both the estimated cell size and properties that are correlated with cell  
324 size can differ from their true value, thus introducing both uncertainty and bias.

325 Missing data is the final source of data uncertainty and can arise for many reasons.  
326 Imputation is a common approach to deal with missing data, where a plausible value is  
327 generated and substituted for the missing value. The imputed value is then taken as the  
328 true value, ignoring that it was generated and not measured. A simple way to account for  
329 uncertainty in an imputed value is to impute many values – known as multiple imputation  
330 – and the variation in the imputed values captures the uncertainty [37, 65]. Predictions  
331 are the made for each imputed value and the predictions combined.

332 To illustrate data uncertainty, Figure 5A shows the same simulated data but with un-  
333 certainty in both the predictor and outcome variable. Here we assume that the error in  $x$   
334 differs for each sample because it depends on the precision of the measurement, whereas  
335 the uncertainty in  $y$  is constant; for example, we know that the measurements are accurate  
336 to  $\pm$  some fixed amount. Variable and fixed uncertainty are accounted for in the same way,  
337 and we include both for illustration purposes.

338 The uncertainty in  $x$  and  $y$  can be handled in two ways. The first is to directly model the  
339 errors in a Bayesian analysis [20, 39, 41, 51]. For example, if  $x$  is measured with error, the  
340 true value can be inferred with  $x_{\text{measured}} \sim N(x_{\text{true}}, \sigma_{\text{error}})$  where  $\sigma_{\text{error}}$  is the uncertainty  
341 in the measured value of  $x$  – usually the standard error. However, fitting such models  
342 may be difficult as each sample (row) adds as many parameters as there are variables with  
343 measurement error (columns).

344 A second simpler option is to generate multiple data sets, where the  $x$  and  $y$  values are  
345 drawn from a distribution [2]. For example, suppose an  $x$  variable for one compound has  
346 a measured value of 0.5 with a standard error of  $\pm 0.06$ . We can sample, say, 10 values  
347 from a normal distribution with a mean of 0.5 and a standard deviation of 0.06, thereby  
348 generating 10 new datasets where the observed value of  $x$  is replaced with the generated  
349 value. This is a form of multiple imputation and Blackwell, Honaker, and King describe  
350 a more sophisticated method of generating new data by taking the correlations between  
351 variables into account [2]. Each dataset is then analysed separately, and the predictions  
352 from each analysis are combined. Variations between the different datasets will lead to  
353 different parameter estimates, which in turn will lead to different predictions, and the  
354 ensemble of predictions captures the uncertainty in  $x$ . The models can be fit to the separate

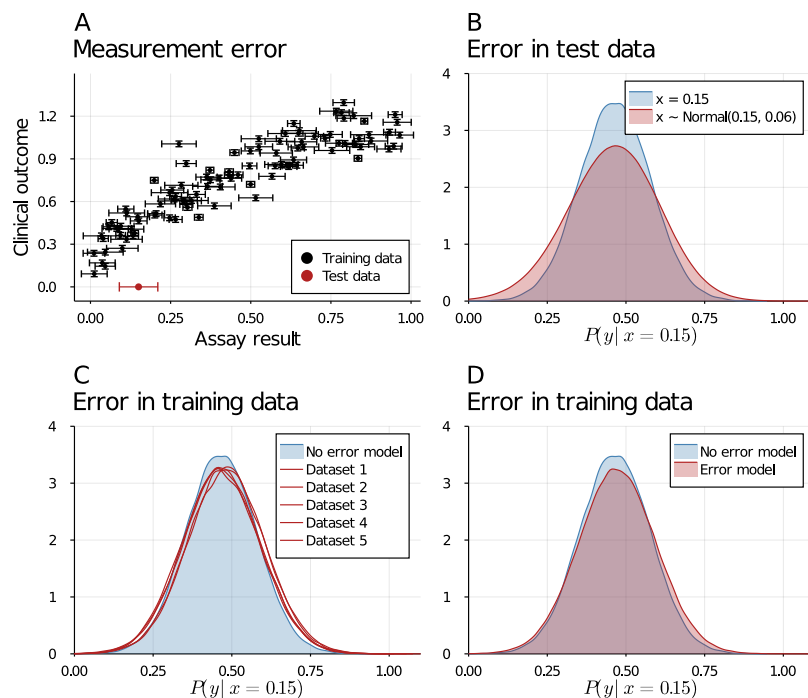


Figure 5: Data uncertainty. Error bars indicate 1 standard error of the estimated assay value and clinical outcome (A). The red point at  $x = 0.15$  is the new value to be predicted. Predictions are more uncertain when measurement error in the test data is included (assuming the training data is measured without error; B). Predictions for the test data when accounting for measurement error in the training data using multiple generated data (assuming no error in the test data; C). Averaged predictions from the generated data shows greater prediction uncertainty compared with ignoring measurement error in the training data (D).

355 datasets in parallel, and so the computation time is the same as fitting the model to one  
 356 dataset.

357 We use the second method to illustrate the effect of ignoring measurement error in  
 358 two ways. First, we assume the training data is measured without error, but the test data  
 359 is measured with error. Then we assume that the training data is measured with error  
 360 but the test data is not. In both cases we compare the result to the standard approach of  
 361 ignoring measurement error in both the training and test data. The test data is a single  
 362 new compound with an assay value of  $x = 0.15 \pm 0.06$ , shown as the red point in Figure  
 363 5A (plotted at  $y = 0$ , but the objective is to predict  $y$ ). The 3-parameter exponential model  
 364 is used in this example. The narrow blue distribution in Figure 5B corresponds to the

standard approach of ignoring uncertainty in both the training and test data, and the red distribution shows the greater uncertainty in the prediction for  $y$  when the measurement error in  $x$  is included. To obtain this distribution, 1000 samples were drawn from a normal distribution with a mean of 0.15 and standard deviation of 0.06. A prediction was made for each of these 1000 samples which reflects the uncertainty in  $x$ .

The blue distribution in Figure 5C is again the standard analysis, and the five red lines show the slightly different predictions from each of the five datasets that account for measurement error in the training data (the test data was assumed to be error-free). Predictions from the five datasets are averaged and shown as the red distribution in Figure 5D, which is slightly wider than the standard analysis from the blue distribution. The additional uncertainty appears negligible in this example, especially compared with uncertainty in the test data (Fig. 5B). However, the variation between datasets is expected to increase with (1) greater uncertainty in the variables, (2) more variables with measurement error included in the model, and (3) more parameters in the model with the total sample size remaining fixed. The effect of measurement error in the training data can be assessed during model development and validation, and if the additional prediction uncertainty is negligible, then the final production model might ignore it.

## Distribution function uncertainty

Another source of uncertainty is the distribution function  $G(\cdot)$ , which represents our uncertainty in a predicted value of  $y$ . Dozens of distributions are available but the list can be narrowed down based on background knowledge of the outcome. For example, if the outcome is binary such as absent/present, safe/toxic, or alive/dead, then a binomial distribution is appropriate; if the outcome is a count such as the number of seizures, then a Poisson or negative binomial distribution are two common options; if the data are positive values and skewed such as liver enzyme levels, then a log-normal or gamma distribution may be suitable; if the outcome is an ordered category such as none/mild/severe, then an ordered categorical distribution would be appropriate; if the outcome is continuous and unbounded with no outliers, then a Gaussian distribution may be suitable; and if there are outliers, a Student-t distribution might be appropriate. Many Bayesian textbooks have appendices that list the common distributions and their properties [18, 36, 38].

Choosing between distributions is made easier because many distributions are special cases of other distributions. For example, both the Gaussian and Cauchy distributions are special cases of the Student-t distribution, the Poisson distribution is a special case of the negative binomial distribution, and the exponential distribution is a special case of a gamma distribution. Hence, we often don't need to choose between a set mutually



exclusive options, but can select the more general distribution and allow the model to determine if one of the special cases is more appropriate. The more general distributions usually have only one additional parameter and therefore do not make the model much more complex. However, not all potentially suitable distributions are related (e.g. gamma and lognormal) and hence two or more models may need to be compared. The data for our running example was generated from a Gaussian distribution and which we have been using for all the models throughout. Hence using the more general Student-t distribution will inform us that the Gaussian is suitable, and so the results are not shown.

A key consideration when selecting a distribution function is the bounds of the data. In our running example, the clinical outcome has a minimum value of zero, but is being modelled with a Gaussian distribution. Since a Gaussian distribution is defined for both positive and negative numbers, there is nothing to prevent negative predictions. A model is clearly inappropriate if it predicts impossible values. Fortunately, we can easily define truncated versions of standard distributions, and so we could specify a Gaussian distribution with a lower bound of zero. Figure 6 shows an example using the 2-parameter exponential model to predict the clinical outcome for an assay value of  $x = 0.05$  both without (A) and with (B) truncation at  $y = 0$ . Without truncation the model gives 9% chance that  $y$  will be less than zero (Fig. 6C). With truncation, this probability gets redistributed to positive values (Fig. 6D, the small proportion of the distribution below zero is a plotting artefact). Hence, if training or test data are near boundaries, using truncated versions of standard distributions is sensible. However, if the data are bounded but the values are far from the boundaries, then accounting for such boundaries may be unnecessary.

## Link function uncertainty

Link functions are required when the predicted mean  $\mu = f_{\mu}(\cdot)$  is bounded by an upper and/or lower limit. For example, when predicting the probability of an event, the predicted value must lie between 0 and 1. Values returned from the mean function are unconstrained and can lie well outside this range. Hence, a link function is used to transform the values to respect the bounds. The logit, probit, cauchit, and complementary log-log functions all take unconstrained numbers and compress them into the 0 – 1 range (Fig. 7). There is no “correct” link function and each provides a different mapping from the unconstrained input to the constrained output and hence gives a different prediction, especially for large values of the input. Link functions are also required for the variances, since variances cannot be negative values, and exponential or power links are often used.

Link functions are analogous to activation functions in neural networks, although they are used as nonlinear transformations between neurons and not necessarily to constrain

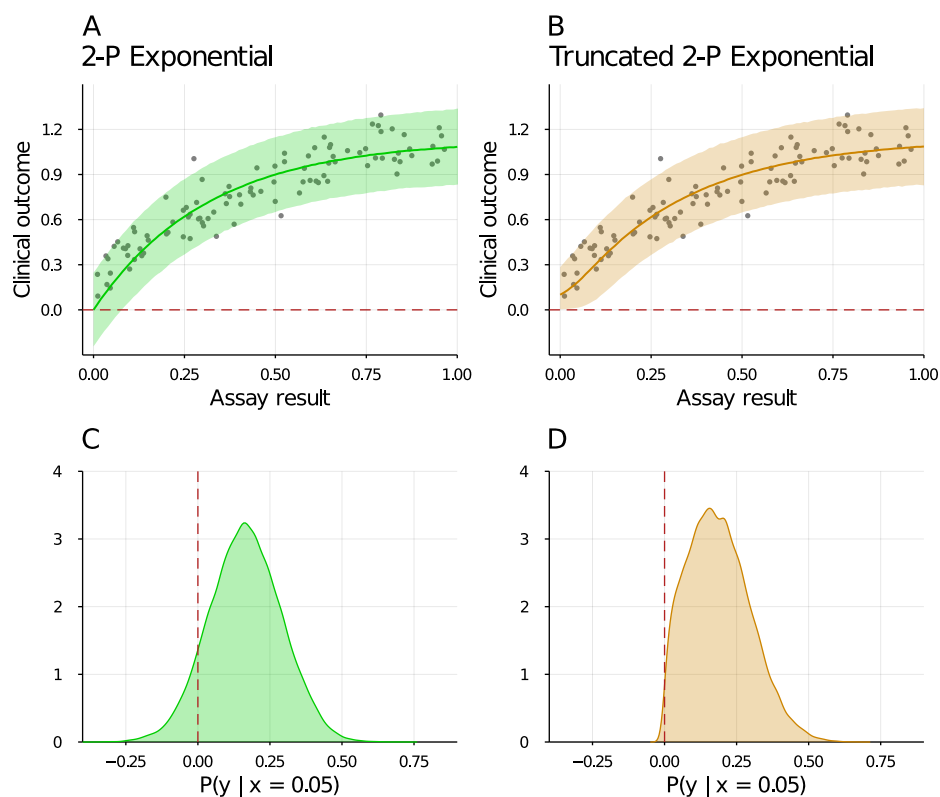


Figure 6: Truncated predictive distributions. Fits and 95% PI for models without (A and C) and with (B and D) a constraint that the predicted values must be positive. Prediction for a new compound given  $x = 0.5$  without the constraint shows that 9.1% of the predicted distribution is negative (C), while the truncated model redistributes the prediction to positive values (D). Red dashed lines are the data boundary.

435 values to allowable values. But the same issue arises in that many activation functions exist  
 436 and different functions will lead to different predictions.

### 437 Variance function uncertainty

438 The variance function models the uncertainty in  $y$  for given values of  $x$ . Another way  
 439 to think of a variance function is that it models the spread of points around the mean  
 440 prediction (Fig. 8). The standard approach assumes that uncertainty in  $y$  is constant (Fig.  
 441 8A). However, when the variance is not constant, a model for  $\sigma$  is required (Fig. 8B).

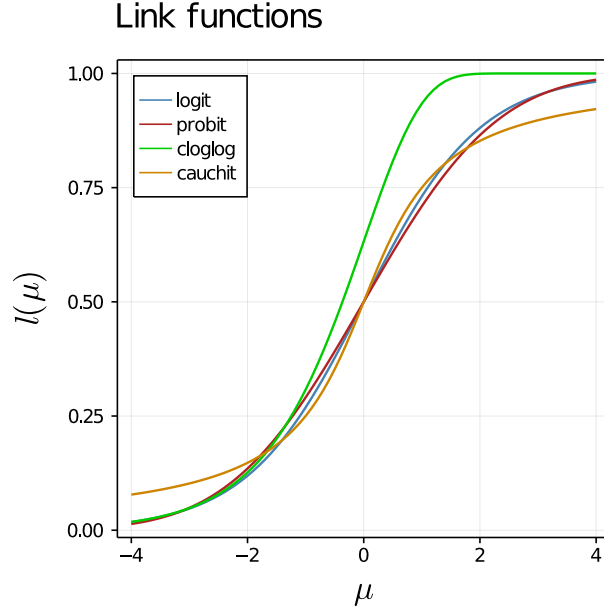


Figure 7: Link function uncertainty. Four functions that map the mean function  $\mu = f(\cdot)$  from  $-\infty$  to  $\infty$  to values between 0 and 1. These link functions are required when the outcome is a probability and must be between 0 and 1. The predicted probabilities therefore differ depending on the link function.

442 Just like modelling  $\mu$  as a function of  $x$ , we now need to model  $\sigma$  as a function of  $x$ .  
 443 This function could be a simple function of one  $x$  variable or a full neural network for  
 444 all  $x$  variables [44]. The latter option involves creating a second neural network for the  
 445 variance, but this doubles the complexity of the model and the training time.

446 In Figure 8B we do not use  $x$  directly, but model  $\sigma$  as a function of  $\mu$  – in other words,  
 447 the uncertainty in  $y$  is proportional to the predicted value of  $y$ . This allows for a simple  
 448 mean function such as  $f_\sigma = \sigma_0 + \sigma_1\mu$ , where  $\sigma_0$  and  $\sigma_1$  are parameters that control the  
 449 relationship between  $\sigma$  and  $\mu$ . Since variances must be positive values, a link function is  
 450 needed to constrain  $f_\sigma$  to be positive, and the softplus function  $l_\sigma = \log(1 + \exp(f_\sigma(\cdot)))$   
 451 is used here. The result is shown in Figure 8B, where the 95% shaded prediction region  
 452 better matches the spread of the data compared with assuming a constant variance (Fig.  
 453 8A).

454 Instead of using a variance function, another option is to transform the outcome vari-  
 455 able (e.g. log, square-root, or inverse), so that the uncertainty in  $y$  is constant. Alter-  
 456 natively, some distributions such as the Poisson and Bernoulli have a defined relationship

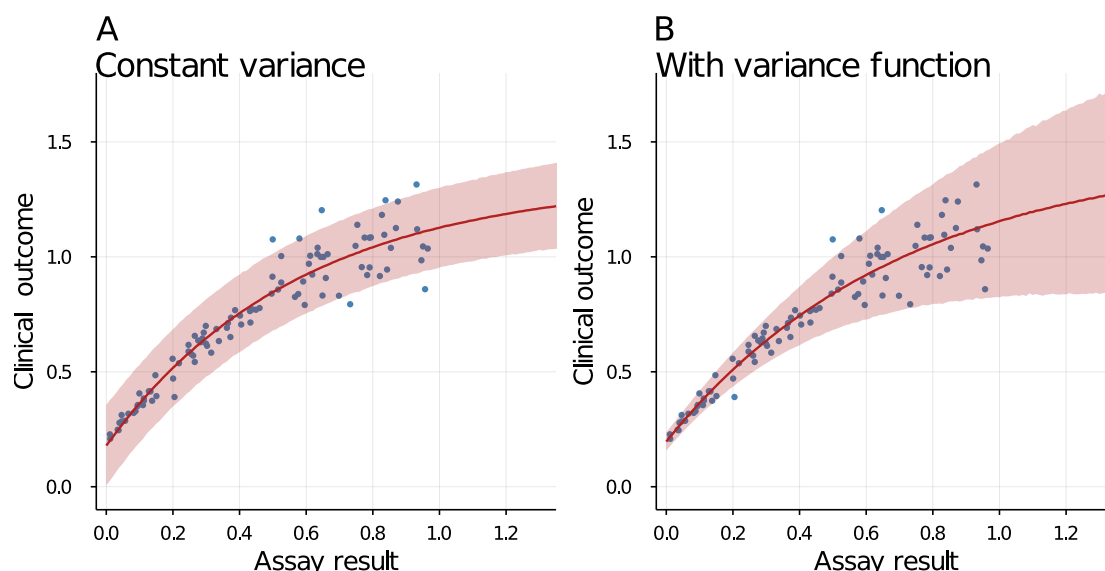


Figure 8: Variance function uncertainty. A constant variance implies that the uncertainty in a prediction is the same for all values  $x$  (A). However, just like the mean prediction can change as function of  $x$ , so can the uncertainty in the prediction (B).

457 between the variance the mean, which allows for non-constant variances, but they are only  
 458 appropriate for certain type of data.

## 459 Sources of uncertainty combined

460 Breaking down the sources of uncertainty into seven items enables us to think about them  
 461 separately and assess their importance when developing a prediction model. A final model  
 462 may include several sources and they can be easily combined. For example, suppose vari-  
 463 ance and link functions were not required but two mean functions and two distribution  
 464 functions performed similarly and therefore four models with each combination of distri-  
 465 bution and mean function are fit to the training data and the predictions averaged. If fully  
 466 Bayesian models are used, parameter uncertainty is already account for. And if the test  
 467 data are measured with error, we can use the approach in 5B to draw multiple samples for  
 468 each test sample and feed them all through the prediction models. The more sources of  
 469 uncertainty accounted for the more complex the prediction model. Hence, sources of un-  
 470 certainty that make little contribution to the overall prediction uncertainty can be ignored.

## 471 Uncertainty for classification tasks

472 The previous examples had a continuous outcome variable, but often outcomes are cate-  
473 gorical such as toxic versus safe. Much of the previous discussion applies, but an important  
474 distinction is between uncertainty in a parameter ( $\mu$ ) and uncertainty in a prediction for  
475 a new observable ( $y$ ) [5, 17]. Greater parameter uncertainty leads to greater prediction  
476 uncertainty, but not for classification tasks, where the objective is to predict which of  $K$   
477 classes a sample belongs to. This point is illustrated in Figure 9.

478 Figure 9A plots data for a 2-group classification task with two predictors ( $x_1, x_2$ ), and  
479 assume the grey triangles are the “toxic” class and the blue circles are the “safe” class.  
480 The black line is the optimal separating boundary. A logistic regression model is used to  
481 separate the classes and the prediction from the model will be a number between 0 and  
482 1, where 1 corresponds toxic and zero corresponds to safe. This prediction is derived  
483 from the mean function and is passed through a link function to constrain the predictions  
484 to lie between 0 and 1. We’ll call these predicted values  $\mu$  and the uncertainty in the  
485 prediction  $\sigma$ . Figure 9B plots  $\mu$  versus  $\sigma$  for each point in Figure 9A, and the inverted-U  
486 relationship is a known feature of such models. But some samples are especially uncertain  
487 and are highlighted in red ( $\sigma \geq 0.8$ ). Samples with intermediate uncertainty ( $\sigma$  between  
488 0.6 and 0.8) are highlighted in yellow. Figure 9C shows that the uncertain samples are all  
489 close to the decision boundary, and that the most uncertain red points lie near the edge  
490 of the data where the location of decision boundary itself is uncertain. The shaded grey  
491 region in Figure 9C represents the uncertainty in the decision boundary, and note how the  
492 uncertainty is wider at the ends compared with the middle.

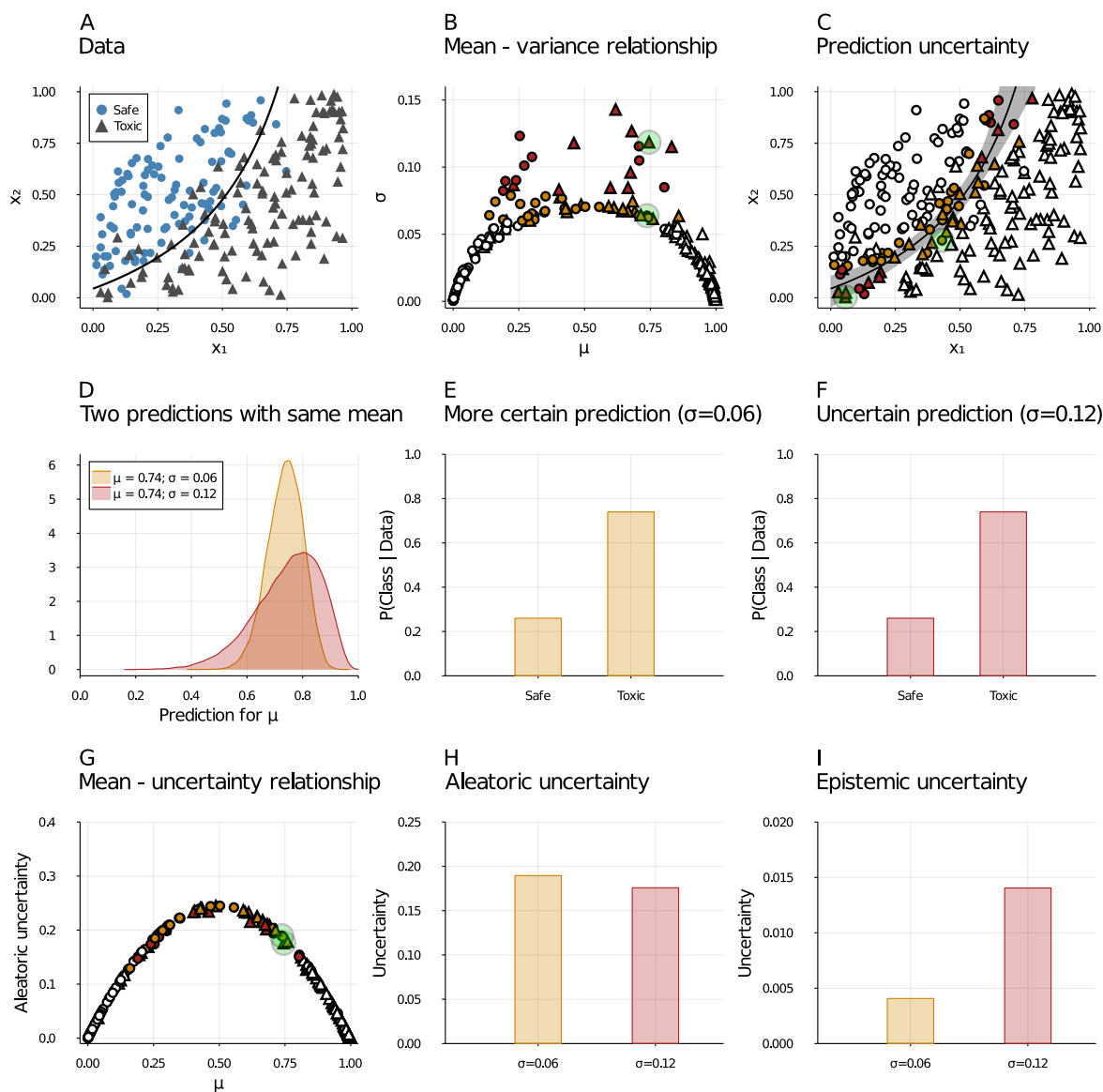


Figure 9: Uncertainty in classification. Simulated data with two features (A). A plot of the mean ( $\mu$ ) and uncertainty ( $\sigma$ ) shows the expected inverted-U relationship (B). The highest variance predictions (red points) are near the decision boundary and at the edge of the data, and high variance predictions (orange points) follow the boundary (C). Two compounds with the same mean but different uncertainties (D), have the same uncertainty in the final predicted value for  $y$  (E,F). Aleatoric uncertainty is fully explained by  $\mu$  (G) and is similar for the two compounds (H). Epistemic uncertainty is nearly 3.5 time greater for the compound with the larger  $\sigma$  (I).

Figure 9D plots the full distribution for two samples that have the same  $\mu$ , but one has twice the uncertainty. Recall that these distributions are for the parameter  $\mu$  and not for the observable outcome  $y$ , which is either safe or toxic. To obtain a prediction for the observable we need a data generating distribution such as the Bernoulli distribution, giving us  $y \sim \text{Bernoulli}(\mu)$ . Figure 9E and F shows the predicted values of  $y$  for these two compounds and note that they are identical, despite different values of  $\sigma$ . At first this may seem strange, but it is the desired behaviour. Consider tossing a coin four times and getting 3 heads, our prediction for the probability of heads is  $3/4 = 0.75$  (and  $1/4 = 0.25$  for tails). Similarly, if we toss a coin 1000 times and get 750 heads, our prediction is still 0.75, even though we are much more certain that the proportion of heads is 0.75 (i.e.  $\sigma$  is much smaller). The width of the distributions in Figure 9D (captured by  $\sigma$ ) provides the *weight of evidence* [26] which quantifies how much the prediction will change as more data are gathered or as the inputs change. For example, a single new observation with four coin tosses alters the probability to either  $3/5 = 0.6$  or  $4/5 = 0.8$ , depending on whether a heads or tails was observed, whereas with 1000 tosses the probability is still 0.75, rounding to two decimal places.

Uncertainty in the predicted classes is often divided into aleatoric and epistemic uncertainty [25, 28]. Aleatoric uncertainty is supposedly due to “inherent randomness” whereas epistemic uncertainty is due to lack of knowledge. Without starting a philosophical debate, we take the position that all uncertainty is due to a lack of knowledge [5, 26]. Nevertheless, we can decompose our uncertainty into two components. The first is how close our point prediction is to zero or one – a more confident prediction would be close to these bounds, and the most uncertain prediction would be 0.5. This corresponds to aleatoric uncertainty. The second component is how confident we are in our point prediction. When a weatherperson states there is a 70% chance of rain tomorrow, they do not mean exactly 70.00000% but  $70\% \pm$  some amount. The uncertainty in the stated value corresponds to epistemic uncertainty and is represented by the width of the distributions in Figure 9D and the parameter  $\sigma$ . Using the approach of Kwon et al. we decompose the two sources of uncertainty for the compounds and plot aleatoric uncertainty versus the mean prediction ( $\mu$ ) in Figure 9G [30]. Note how  $\mu$  completely explains aleatoric uncertainty. The two compounds have similar aleatoric uncertainty since they have a similar value of  $\mu$  (Fig. 9H). However, the compound with the larger value of  $\sigma$  has nearly 3.5 times greater epistemic uncertainty (Fig. 9H). Epistemic uncertainty or the weight of evidence ( $\sigma$ ), provides important information about the uncertainty of a prediction, which is critical for high-stakes decisions.

## 528 Discussion

### 529 Generalisations and extensions

530 The above examples used simple models but this framework can be generalised to more  
531 complex cases. For example, we had functions for the mean and variance, but any parame-  
532 ter in the distribution function can be modelled. For example, a Student-t distribution has  
533 a parameter called the degrees of freedom (df) which controls the heaviness of the tails.  
534 The df could be modelled as a function of  $x$ , just like  $\mu$  or  $\sigma$  [52].

535 The above examples used a single distribution function, but flexibility can be increased  
536 by using mixtures of distributions. For example, outliers can be modelled with a mixture  
537 of Gaussian distributions: one to account for the regular observations and the second  
538 to account for the outliers. Metabolite, gene, and protein levels are non-negative and  
539 often positively skewed, and hence gamma or lognormal distributions may be appropriate.  
540 But these distributions are only defined for values *greater* than zero, and there may be  
541 zeros in the data, which are often dealt with by adding a small value to all data points.  
542 A better option can be to model the data with a two-part model, one which accounts  
543 for the zeros and the other (e.g. gamma or lognormal) which accounts for the non-zero  
544 values. Such “hurdle models” provide this flexibility and also return a parameter that  
545 estimates the proportion of zeros, which may be scientifically interesting [13]. Taking this  
546 idea a step further, Dirichlet Process models allow us to specify as many distributions as  
547 needed to model the data. Instead of specifying a single distribution, we specify a prior  
548 over distributions, and learn them from the data (yes, we can specify a distribution over  
549 distributions! [42]).

550 The above examples also used a single mean function for each model, but it’s possi-  
551 ble to have a distribution of mean functions, which are called Gaussian Process models  
552 [19, 49, 54]. These flexible models can fit complex relationships between  $x$  and  $y$ . Surpris-  
553 ingly, they are not implemented via the mean function, but by generalising the variance  
554 function to make it a covariance function. Covariance functions are not discussed here  
555 but they are also useful for modelling hierarchical or nested data [24, 48], and for mod-  
556 elling dependencies in time or space. Neural networks [21, 56, 68] and Bayesian additive  
557 regression trees (BART) [12, 60] are other options for flexible mean functions.



## 558 Further advantages of PPMs

559 In addition to providing prediction uncertainty, PPMs have several other benefits. Hyperpa-  
560 rameter values are typically selected by trying many options and choosing the combination  
561 that performs best. To avoid overfitting, crossvalidation or a similar approach divides the  
562 training data into smaller subsets, some of which are used for training and others to assess  
563 performance. But with small datasets, crossvalidation can give unstable models and a poor  
564 assessment of performance. Many Bayesian approaches can learn values of some hyperpa-  
565 rameters using all the training data and have a built-in prevention of overfitting [64, 71].  
566 They also incorporate the uncertainty in the hyperparameters in the predictions. Models  
567 can still be compared using only the training data by estimating leave-one-out (LOO) cross-  
568 validation performance, without the computational cost of actually retraining the model  
569 for each sample [66, 67]. Vehtari and colleagues have also developed methods to assess  
570 when a LOO estimate is unreliable, and the model can be retrained only for these samples  
571 [69].

572 Another advantage is that background information such as adverse outcome pathways  
573 [7], constraints on parameters [34], or monotonic relationships [14] can often be incorpo-  
574 rated into the model, which can guide the model to better solutions.

575 Often several structurally similar compounds are available that have different binding  
576 affinities or potencies, but also with different results in the toxicity assays, and a decision  
577 must be taken to designate one compound in the series as the lead. PPMs can not only rank  
578 compounds but also obtain an uncertainty in the ranking, thus enabling decisions makes  
579 to conclude that one compound is reliably better than another [34, 55].

580 Finally, many popular machine learning methods have a PPM or Bayesian analogue,  
581 including regularised linear and generalised linear models (lasso, ridge regression) [10,  
582 45, 47, 53], tree models (random forests, xgboost) [11, 12, 60], support vector machines  
583 [59, 64], and neural networks [43, 56, 68]. Hence, it is often possible to convert your  
584 favourite model into one that provides prediction uncertainty.

## 585 Drawbacks and challenges

586 The main drawback of Bayesian or other PMMs is that they require more work, possibly  
587 twice as much, since getting appropriately calibrated uncertainty is just as hard as getting  
588 accurate predictions. For example, 95% prediction intervals should contain 95% of the  
589 out-of-sample or test data values. [72].

590 For fully Bayesian methods, the computational overhead may be high, making it diffi-  
591 cult to iteratively fit, check, and update models during development (although computa-  
592 tions are often much quicker when making predictions). Storage for parameter values may  
593 be a problem for large models since this equals the number of parameters times number  
594 of Markov chain Monte Carlo draws. These approaches may therefore be harder to scale  
595 to large datasets, but faster and scalable algorithms is an active area of research. Another  
596 solution to large data is to cleverly select a weighted subset of samples that is much smaller  
597 than the original but captures the essential features. This “coreset” approach enables stan-  
598 dard PPM methods to be used on the smaller dataset with little loss of information [8, 22].

599 Finally, not all sources of uncertainty can be captured. Many sources of uncertainty  
600 discussed above arise because many modelling options are available, and different choices  
601 lead to different predictions. All of the choices relate to the prediction model, but many  
602 decisions need to be made outside of the model. We refer to these extra-model choices  
603 as the workflow and they include experimental decisions such as the technology, cell-line,  
604 assay, antibodies, protocol, and so on. Also included are data processing pipelines where  
605 raw data are cleaned, transformed, categorised, coded, and normalised before they are  
606 entered into a prediction model. A single workflow is commonly used, with the untested  
607 assumption that variations in the workflow will lead to the same predictions and results.  
608 However, variations in workflows and analytic decisions do lead to variations results [4,  
609 23, 32, 57, 58, 61, 62].

## 610 **Reporting uncertainty to help risk communication and decision making**

611 The ultimate aim of prediction models in drug discovery is to enable better decision mak-  
612 ing. Thus, not only should predictions be accurate with prediction uncertainty adequately  
613 represented, but the results should be easy to understand by decision makers. Fortunately,  
614 PPMs provide intuitive results for continuous (Fig. 6D), binary (Fig. 9D), categorical, and  
615 ordered categorical outcomes [56, 71], as well as for compound rankings [34, 55]. We  
616 have found that safety pharmacologists and other project members can easily interpret the  
617 predictive distributions provided by PPMs and value the confidence in the predictions that  
618 these distributions provide [34, 71].

619 With recent advances in algorithms, hardware, and software, Bayesian or other PPMs  
620 are now feasible for most – if not all – machine learning problems encountered in drug  
621 discovery. Making PPMs the standard approach for critical ML problems will enable more  
622 informed and better decisions.

## References

- [1] Bezanson J, Edelman A, Karpinski S, Shah VB (2017). Julia: A fresh approach to numerical computing. *SIAM review* 59(1): 65–98.
- [2] Blackwell M, Honaker J, King G (2015). A Unified Approach to Measurement Error and Missing Data: Overview and Applications. *Sociological Methods & Research* 46(3): 303–341.
- [3] Blundell C, Cornebise J, Kavukcuoglu K, Wierstra D (2015). Weight Uncertainty in Neural Networks. In: *Proceedings of the 32nd International Conference on International Conference on Machine Learning*, vol. 37 of *ICML’15*, 1613–1622.
- [4] Botvinik-Nezer R, Holzmeister F, Camerer CF, Dreber A, Huber J, Johannesson M, Kirchler M, Iwanir R, Mumford JA, Adcock RA, Avesani P, Baczkowski BM, Bajracharya A, Bakst L, Ball S, Barilari M, Bault N, Beaton D, Beitner J, Benoit RG, Berkers RMWJ, Bhanji JP, Biswal BB, Bobadilla-Suarez S, Bortolini T, Bottenhorn KL, Bowring A, Braem S, Brooks HR, Brudner EG, Calderon CB, Camilleri JA, Castrellon JJ, Cecchetti L, Cieslik EC, Cole ZJ, Collignon O, Cox RW, Cunningham WA, Czoschke S, Dadi K, Davis CP, Luca AD, Delgado MR, Demetriou L, Dennison JB, Di X, Dickie EW, Dobryakova E, Donnat CL, Dukart J, Duncan NW, Durnez J, Eed A, Eickhoff SB, Erhart A, Fontanesi L, Fricke GM, Fu S, Galván A, Gau R, Genon S, Glatard T, Glerean E, Goeman JJ, Golowin SAE, González-García C, Gorgolewski KJ, Grady CL, Green MA, Guassi Moreira JF, Guest O, Hakimi S, Hamilton JP, Hancock R, Handjaras G, Harry BB, Hawco C, Herholz P, Herman G, Heunis S, Hoffstaedter F, Hogeveen J, Holmes S, Hu CP, Huettel SA, Hughes ME, Iacovella V, Iordan AD, Isager PM, Isik AI, Jahn A, Johnson MR, Johnstone T, Joseph MJE, Juliano AC, Kable JW, Kassinosopoulos M, Koba C, Kong XZ, Kosciuk TR, Kucukboyaci NE, Kuhl BA, Kupek S, Laird AR, Lamm C, Langner R, Lauharatanahirun N, Lee H, Lee S, Leemans A, Leo A, Lesage E, Li F, Li MYC, Lim PC, Lintz EN, Liphardt SW, Losecaat Vermeer AB, Love BC, Mack ML, Malpica N, Marins T, Maumet C, McDonald K, McGuire JT, Melero H, Méndez Leal AS, Meyer B, Meyer KN, Mihai G, Mitsis GD, Moll J, Nielson DM, Nilsson G, Notter MP, Olivetti E, Onicas AI, Papale P, Patil KR, Peelle JE, Pérez A, Pischetta D, Poline JB, Prystauka Y, Ray S, Reuter-Lorenz PA, Reynolds RC, Ricciardi E, Rieck JR, Rodriguez-Thompson AM, Romyn A, Salo T, Samanez-Larkin GR, Sanz-Morales E, Schlichting ML, Schultz DH, Shen Q, Sheridan MA, Silvers JA, Skagerlund K, Smith A, Smith DV, Sokol-Hessner P, Steinkamp SR, Tashjian SM, Thirion B, Thorp JN, Tinghög G, Tisdall L, Thompson SH, Toro-Serey C, Torre Tresols JJ, Tozzi L, Truong V, Turella L, van ’t Veer AE, Verguts T, Vettel JM, Vijayarajah S, Vo K, Wall MB, Weeda WD, Weis S, White DJ, Wisniewski D, Xifra-Porxas A, Yearling EA, Yoon S, Yuan R, Yuen KSL, Zhang L, Zhang X, Zosky JE, Nichols TE, Poldrack RA, Schonberg T (2020). Variability in the analysis of a single neuroimaging dataset by many teams. *Nature* 582: 84–88.

- 661 [5] Briggs W (2016). *Uncertainty: The Soul of Modeling, Probability and Statistics*. New  
662 York, NY: Springer.
- 663 [6] Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, Neelakantan A, Shyam  
664 P, Sastry G, Askell A, Agarwal S, Herbert-Voss A, Krueger G, Henighan T, Child R,  
665 Ramesh A, Ziegler DM, Wu J, Winter C, Hesse C, Chen M, Sigler E, Litwin M, Gray S,  
666 Chess B, Clark J, Berner C, McCandlish S, Radford A, Sutskever I, Amodei D (2020).  
667 Language Models are Few-Shot Learners. *arXiv* .
- 668 [7] Burgoon LD, Angrish M, Garcia-Reyero N, Pollesch N, Zupanic A, Perkins E (2019).  
669 Predicting the Probability that a Chemical Causes Steatosis Using Adverse Outcome  
670 Pathway Bayesian Networks (AOPBNs). *Risk Analysis* 40(3): 512–523.
- 671 [8] Campbell T, Broderick T (2018). Bayesian Coreset Construction via Greedy Iterative  
672 Geodesic Ascent. In: J Dy, A Krause (Eds.) *Proceedings of the 35th International Con-*  
673 *ference on Machine Learning*, vol. 80 of *Proceedings of Machine Learning Research*,  
674 698–706, PMLR.
- 675 [9] Carroll RJ, Ruppert D, Stefanski LA, Crainiceanu CM (2006). *Measurement Error in*  
676 *Nonlinear Models: A Modern Perspective*. Boca Raton, FL: Chapman & Hall/CRC, 2nd  
677 edn.
- 678 [10] Carvalho CM, Polson NG, Scott JG (2009). Handling Sparsity via the Horseshoe. *Pro-*  
679 *ceedings of Machine Learning Research* 5: 73–80.
- 680 [11] Chipman HA, George EI, McCulloch RE (1998). Bayesian CART Model Search. *Jour-*  
681 *nal of the American Statistical Association* 93(443): 935–948.
- 682 [12] Chipman HA, George EI, McCulloch RE (2010). BART: Bayesian additive regression  
683 trees. *The Annals of Applied Statistics* 4(1): 266–298.
- 684 [13] Cragg JG (1971). Some Statistical Models for Limited Dependent Variables with Ap-  
685 plication to the Demand for Durable Goods. *Econometrica* 39(5): 829.
- 686 [14] DePalma G, Craig BA (2017). Bayesian monotonic errors-in-variables models with  
687 applications to pathogen susceptibility testing. *Statistics in Medicine* 37(3): 487–502.
- 688 [15] Gal Y, Ghahramani Z (2016). Dropout as a Bayesian Approximation: Representing  
689 Model Uncertainty in Deep Learning. In: MF Balcan, KQ Weinberger (Eds.) *Proceed-*  
690 *ings of The 33rd International Conference on Machine Learning*, vol. 48 of *Proceedings*  
691 *of Machine Learning Research*, 1050–1059, New York, New York, USA: PMLR.
- 692 [16] Ge H, Xu K, Ghahramani Z (2018). Turing: a language for flexible probabilistic in-  
693 ference. In: *International Conference on Artificial Intelligence and Statistics, AISTATS*  
694 *2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*, 1682–1690.

- 695 [17] Geisser S (1993). *Predictive Inference: An Introduction*. New York, NY: Chapman &  
696 Hall.
- 697 [18] Gelman A, Carlin JB, Stern HS, Rubin DB (2004). *Bayesian Data Analysis*. Boca Raton:  
698 Chapman & Hall/CRC, 2nd edn.
- 699 [19] Gramacy RB (2020). *Surrogates: Gaussian Process Modeling, Design, and Optimization*  
700 *for the Applied Sciences*. CRC Press.
- 701 [20] Gustafson P (2004). *Measurement Error and Misclassification in Statistics and Epidemi-*  
702 *ology: Impacts and Bayesian Adjustments*. Boca Raton: Chapman & Hall/CRC.
- 703 [21] Hirschfeld L, Swanson K, Yang K, Barzilay R, Coley CW (2020). Uncertainty Quantifi-  
704 cation Using Neural Networks for Molecular Property Prediction. *Journal of Chemical*  
705 *Information and Modeling* 60(8): 3770–3780.
- 706 [22] Huggins JH, Campbell T, Broderick T (2016). Coresets for Scalable Bayesian Logistic  
707 Regression. *arXiv* .
- 708 [23] Huntington-Klein N, Arenas A, Beam E, Bertoni M, Bloem JR, Burli P, Chen N, Grieco  
709 P, Ekpe G, Pugatch T, Saavedra M, Stopnitzky Y (2021). The influence of hidden  
710 researcher decisions in applied microeconomics. *Economic Inquiry* .
- 711 [24] Johnstone RH, Bardenet R, Gavaghan DJ, Mirams GR (2016). Hierarchical Bayesian  
712 inference for ion channel screening dose-response data. *Wellcome Open Res* 1: 6.
- 713 [25] Kendall A, Gal Y (2017). What Uncertainties Do We Need in Bayesian Deep Learning  
714 for Computer Vision? In: I Guyon, UV Luxburg, S Bengio, H Wallach, R Fergus,  
715 S Vishwanathan, R Garnett (Eds.) *Advances in Neural Information Processing Systems*,  
716 vol. 30, Curran Associates, Inc.
- 717 [26] Keynes JM (1921). *A Treatise on Probability*. London: Macmillan & Co.
- 718 [27] Khosravi A, Nahavandi S, Creighton D, Atiya AF (2011). Comprehensive Review of  
719 Neural Network-Based Prediction Intervals and New Advances. *IEEE Transactions on*  
720 *Neural Networks* 22(9): 1341–1356.
- 721 [28] Kiureghian AD, Ditlevsen O (2009). Aleatory or epistemic? Does it matter? *Structural*  
722 *Safety* 31(2): 105–112.
- 723 [29] Kristiadi A, Hein M, Hennig P (2020). Being Bayesian, Even Just a Bit, Fixes Over-  
724 confidence in ReLU Networks. In: HD III, A Singh (Eds.) *Proceedings of the 37th Inter-*  
725 *national Conference on Machine Learning*, vol. 119 of *Proceedings of Machine Learning*  
726 *Research*, 5436–5446, PMLR.

- [30] Kwon Y, Won JH, Kim BJ, Paik MC (2020). Uncertainty quantification using Bayesian neural networks in classification: Application to biomedical image segmentation. *Computational Statistics & Data Analysis* 142: 106816.
- [31] Lakshminarayanan B, Pritzel A, Blundell C (2017). Simple and scalable predictive uncertainty estimation using deep ensembles. In: *Advances in Neural Information Processing Systems*, vol. 30, 6402–6413.
- [32] Landy JF, Jia ML, Ding IL, Viganola D, Tierney W, Dreber A, Johannesson M, Pfeiffer T, Ebersole CR, Gronau QF, Ly A, van den Bergh D, Marsman M, Derks K, Wagenmakers EJ, Proctor A, Bartels DM, Bauman CW, Brady WJ, Cheung F, Cimpian A, Dohle S, Donnellan MB, Hahn A, Hall MP, Jiménez-Leal W, Johnson DJ, Lucas RE, Monin B, Montealegre A, Mullen E, Pang J, Ray J, Reiner DA, Reynolds J, Sowden W, Storage D, Su R, Tworek CM, Bavel JJV, Walco D, Wills J, Xu X, Yam KC, Yang X, Cunningham WA, Schweinsberg M, Urwitz M, Collaboration TCHT, Uhlmann EL (2020). Crowdsourcing hypothesis tests: Making transparent how design choices shape research results. *Psychological Bulletin* 146(5): 451–479.
- [33] Lazic SE (2018). Four simple ways to increase power without increasing the sample size. *Lab. Anim.* 52(6): 621–629.
- [34] Lazic SE, Edmunds N, Pollard CE (2018). Predicting drug safety and communicating risk: benefits of a Bayesian approach. *Toxicol. Sci.* 162(1): 89–98.
- [35] Lazic SE, Williams DP (2020). Improving drug safety predictions by reducing poor analytical practices. *Toxicology Research and Application* 4: 239784732097863.
- [36] Lesaffre E, Lawson AB (2012). *Bayesian Biostatistics*. Chichester, UK: Wiley.
- [37] Little RJA, Rubin DB (2020). *Statistical Analysis with Missing Data*. Hoboken, NJ: Wiley, 3rd edn.
- [38] Lunn D, Jackson C, Best N, Thomas A, Spiegelhalter D (2013). *The BUGS Book: A Practical Introduction to Bayesian Analysis*. Boca Raton, FL: CRC Press.
- [39] McElreath R (2016). *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. Boca Raton, FL: CRC Press.
- [40] Mervin LH, Johansson S, Semenova E, Giblin KA, Engkvist O (2021). Uncertainty quantification in drug design. *Drug Discovery Today* 26(2): 474–489.
- [41] Muff S, Riebler A, Held L, Rue H, Saner P (2014). Bayesian analysis of measurement error models using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64(2): 231–252.

- [42] Muller P, Quintana FA, Jara A, Hanson T (2015). *Bayesian Nonparametric Data Analysis*. Springer.
- [43] Neal RM (1996). *Bayesian Learning for Neural Networks*. New York, NY: Springer.
- [44] Nix D, Weigend A (1994). Estimating the mean and variance of the target probability distribution. In: *Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94)*, IEEE.
- [45] Park T, Casella G (2008). The Bayesian Lasso. *Journal of the American Statistical Association* 103(482): 681–686.
- [46] Pearce T, Leibfried F, Brintrup A (2020). Uncertainty in Neural Networks: Approximately Bayesian Ensembling. In: S Chiappa, R Calandra (Eds.) *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, vol. 108 of *Proceedings of Machine Learning Research*, 234–244, PMLR.
- [47] Piironen J, Vehtari A (2017). Sparsity information and regularization in the horseshoe and other shrinkage priors. *Electronic Journal of Statistics* 11(2): 5018–5051.
- [48] Pinheiro JC, Bates DM (2000). *Mixed-Effects Models in S and S-Plus*. London: Springer.
- [49] Rasmussen CE, Williams CKI (2006). *Gaussian Processes for Machine Learning*. Cambridge, MA: MIT Press.
- [50] Reynolds J, Malcomber S, White A (2020). A Bayesian approach for inferring global points of departure from transcriptomics data. *Computational Toxicology* 100138.
- [51] Richardson S, Gilks WR (1993). A Bayesian approach to measurement error problems in epidemiology using conditional independence models. *American Journal of Epidemiology* 138(6): 430–442.
- [52] Rigby R (2020). *Distributions for modelling location, scale, and shape : using GAMLSS in R*. Boca Raton, FL: CRC Press.
- [53] Ročková V, George EI (2018). The Spike-and-Slab LASSO. *Journal of the American Statistical Association* 113(521): 431–444.
- [54] Schulz E, Speekenbrink M, Krause A (2018). A tutorial on Gaussian process regression: Modelling, exploring, and exploiting functions. *Journal of Mathematical Psychology* 85: 1–16.
- [55] Semenova E, Guerriero ML, Zhang B, Hock A, Hopcroft P, Kadamur G, Afzal AM, Lazic SE (2020). Flexible fitting of PROTAC concentration-response curves with Gaussian Processes. *bioRxiv* .

- 792 [56] Semenova E, Williams DP, Afzal AM, Lazic SE (2020). A Bayesian neural network for  
793 toxicity prediction. *Computational Toxicology* 16: 100133.
- 794 [57] Shi L, Campbell G, Jones WD, Campagne F, Wen Z, Walker SJ, Su Z, Chu TM, Good-  
795 said FM, Pusztai L, Shaughnessy JD, Oberthuer A, Thomas RS, Paules RS, Fielden M,  
796 Barlogie B, Chen W, Du P, Fischer M, Furlanello C, Gallas BD, Ge X, Megherbi DB,  
797 Symmans WF, Wang MD, Zhang J, Bitter H, Brors B, Bushel PR, Bylesjo M, Chen M,  
798 Cheng J, Cheng J, Chou J, Davison TS, Delorenzi M, Deng Y, Devanarayan V, Dix DJ,  
799 Dopazo J, Dorff KC, Elloumi F, Fan J, Fan S, Fan X, Fang H, Gonzaludo N, Hess KR,  
800 Hong H, Huan J, Irizarry RA, Judson R, Juraeva D, Lababidi S, Lambert CG, Li L, Li  
801 Y, Li Z, Lin SM, Liu G, Lobenhofer EK, Luo J, Luo W, McCall MN, Nikolsky Y, Pennello  
802 GA, Perkins RG, Philip R, Popovici V, Price ND, Qian F, Scherer A, Shi T, Shi W, Sung  
803 J, Thierry-Mieg D, Thierry-Mieg J, Thodima V, Trygg J, Vishnuvajjala L, Wang SJ, Wu  
804 J, Wu Y, Xie Q, Yousef WA, Zhang L, Zhang X, Zhong S, Zhou Y, Zhu S, Arasappan D,  
805 Bao W, Lucas AB, Berthold F, Brennan RJ, Bunes A, Catalano JG, Chang C, Chen R,  
806 Cheng Y, Cui J, Czika W, Demichelis F, Deng X, Dosymbekov D, Eils R, Feng Y, Fostel  
807 J, Fulmer-Smentek S, Fuscoe JC, Gatto L, Ge W, Goldstein DR, Guo L, Halbert DN,  
808 Han J, Harris SC, Hatzis C, Herman D, Huang J, Jensen RV, Jiang R, Johnson CD,  
809 Jurman G, Kahlert Y, Khuder SA, Kohl M, Li J, Li L, Li M, Li QZ, Li S, Li Z, Liu J, Liu Y,  
810 Liu Z, Meng L, Madera M, Martinez-Murillo F, Medina I, Meehan J, Miclaus K, Moffitt  
811 RA, Montaner D, Mukherjee P, Mulligan GJ, Neville P, Nikolskaya T, Ning B, Page GP,  
812 Parker J, Parry RM, Peng X, Peterson RL, Phan JH, Quanz B, Ren Y, Riccadonna S,  
813 Roter AH, Samuelson FW, Schumacher MM, Shambaugh JD, Shi Q, Shippy R, Si S,  
814 Smalter A, Sotiriou C, Soukup M, Staedtler F, Steiner G, Stokes TH, Sun Q, Tan PY,  
815 Tang R, Tezak Z, Thorn B, Tsyganova M, Turpaz Y, Vega SC, Visintainer R, von Frese  
816 J, Wang C, Wang E, Wang J, Wang W, Westermann F, Willey JC, Woods M, Wu S,  
817 Xiao N, Xu J, Xu L, Yang L, Zeng X, Zhang J, Zhang L, Zhang M, Zhao C, Puri RK,  
818 Scherf U, Tong W, Wolfinger RD, Consortium MAQC (2010). The MicroArray Quality  
819 Control (MAQC)-II study of common practices for the development and validation of  
820 microarray-based predictive models. *Nature biotechnology* 28: 827–838.
- 821 [58] Silberzahn R, Uhlmann EL, Martin DP, Anselmi P, Aust F, Awtrey E, Bahník Š, Bai F,  
822 Bannard C, Bonnier E, Carlsson R, Cheung F, Christensen G, Clay R, Craig MA, Rosa  
823 AD, Dam L, Evans MH, Cervantes IF, Fong N, Gamez-Djokic M, Glenz A, Gordon-  
824 McKeon S, Heaton TJ, Hederos K, Heene M, Mohr AJH, Högden F, Hui K, Johan-  
825 nesson M, Kalodimos J, Kaszubowski E, Kennedy DM, Lei R, Lindsay TA, Liverani  
826 S, Madan CR, Molden D, Molleman E, Morey RD, Mulder LB, Nijstad BR, Pope NG,  
827 Pope B, Prenoveau JM, Rink F, Robusto E, Roderique H, Sandberg A, Schlüter E,  
828 Schönbrodt FD, Sherman MF, Sommer SA, Sotak K, Spain S, Spörlein C, Stafford T,  
829 Stefanutti L, Tauber S, Ullrich J, Vianello M, Wagenmakers EJ, Witkowiak M, Yoon S,  
830 Nosek BA (2018). Many Analysts, One Data Set: Making Transparent How Variations



- 831 in Analytic Choices Affect Results. *Advances in Methods and Practices in Psychological*  
832 *Science* 1(3): 337–356.
- 833 [59] Sollich P (2002). Bayesian methods for support vector machines: evidence and pre-  
834 dictive class probabilities. *Machine Learning* 46(1/3): 21–52.
- 835 [60] Sparapani RA, Logan BR, McCulloch RE, Laud PW (2016). Nonparametric survival  
836 analysis using Bayesian Additive Regression Trees (BART). *Statistics in Medicine*  
837 35(16): 2741–2753.
- 838 [61] Stanton-Geddes J, de Freitas CG, Dambros CdS (2014). In defense of P values: com-  
839 ment on the statistical methods actually used by ecologists. *Ecology* 95: 637–642.
- 840 [62] Steegen S, Tuerlinckx F, Gelman A, Vanpaemel W (2016). Increasing transparency  
841 through a multiverse analysis. *Perspectives on Psychological Science* 11(5): 702–712.
- 842 [63] Teye M, Azizpour H, Smith K (2018). Bayesian Uncertainty Estimation for Batch Nor-  
843 malized Deep Networks. In: J Dy, A Krause (Eds.) *Proceedings of the 35th International*  
844 *Conference on Machine Learning*, vol. 80 of *Proceedings of Machine Learning Research*,  
845 4907–4916, PMLR.
- 846 [64] Tipping ME (2001). Sparse Bayesian learning and the relevance vector machine.  
847 *JMLR* 1(1): 211–244.
- 848 [65] van Buuren S (2012). *Flexible Imputation of Missing Data*. Boca Raton, FL: CRC Press.
- 849 [66] Vehtari A, Gelman A, Gabry J (2016). Erratum to: Practical Bayesian model evalua-  
850 tion using leave-one-out cross-validation and WAIC. *Statistics and Computing* 27(5):  
851 1433–1433.
- 852 [67] Vehtari A, Gelman A, Gabry J (2016). Practical Bayesian model evaluation using  
853 leave-one-out cross-validation and WAIC. *Statistics and Computing* 27(5): 1413–  
854 1432.
- 855 [68] Vehtari A, Lampinen J (2000). Bayesian Neural Networks: Case Studies in Indus-  
856 trial Applications. In: *Soft Computing in Industrial Applications*, 415–424, Springer  
857 London.
- 858 [69] Vehtari A, Simpson D, Gelman A, Yao Y, Gabry J (2015). Pareto Smoothed Importance  
859 Sampling. *arXiv* 1507.02646.
- 860 [70] Welling M, Teh YW (2011). Bayesian learning via stochastic gradient Langevin dy-  
861 namics. In: *Proceedings of the 28th International Conference on International Confer-*  
862 *ence on Machine Learning*, vol. 33 of *ICML’11*, 681–688.

- 863 [71] Williams DP, Lazic SE, Foster AJ, Semenova E, Morgan P (2020). Predicting drug-  
864 induced liver injury with Bayesian machine learning. *Chem. Res. Toxicol.* 33(1): 239–  
865 248.
- 866 [72] Zhang Y, Lee AA (2019). Bayesian semi-supervised learning for uncertainty-  
867 calibrated prediction of molecular properties and active learning. *Chemical Science*  
868 10(35): 8154–8163.