

A probability transducer for machine-learning classifiers

K. Dyrland 
<kjetil.dyrland@gmail.com>

A. S. Lundervold [†]
<alexander.selvikvag.lundervold@hvl.no>


P.G.L. Porta Mana 
<pgl@portamana.org>

(listed alphabetically)

Dept of Computer science, Electrical Engineering and Mathematical Sciences
Western Norway University of Applied Sciences, Bergen, Norway

[†]& Mohn Medical Imaging and Visualization Centre, Bergen, Norway

Draft. 14 April 2022; updated 24 May 2022

 The present work sells probabilities for machine-learning algorithms at a bargain price. It shows that these probabilities are just what those algorithms need in order to increase their sales.

1 The inadequacy of common classification approaches

As the potential of using machine-learning algorithms in important fields such as medicine or drug discovery increases, the machine-learning community ought to keep in mind what the actual needs and inference contexts in such fields are. We must avoid trying (consciously or unconsciously) to convince such fields to change their needs or ignore their contexts just to fit machine-learning solutions that are available and fashionable at the moment. Rather, we must make sure the that solutions fit needs & context, and amend them if they do not.


The machine-learning mindset and approach to problems such as classification in such new important fields is often still inadequate in many respects. It reflects simpler needs and contexts of many early inference problems successfully tackled by machine learning.

A stereotypical ‘cat vs dog’ image classification, for instance, has four very important differences from a ‘disease *I* vs disease *II*’ medical classification, or from an ‘active vs inactive’ drug classification:

(i) Nobody presumably dies or loses large amounts of money if a cat image is misclassified as dog or vice versa. But a person can die if a disease is misdiagnosed; huge capitals can be lost if an ultimately

ineffective drug candidate is pursued. The *gains and losses* – or generally speaking the *utilities* – of correct and incorrect classifications in the former problem and in the two latter problems are vastly different.

(ii) To what purpose do we try to guess whether an image's subject is a cat or a dog? For example because we must decide whether to put it in the folder 'cats' or in the folder 'dogs'. To what purpose do we try to guess a patient's disease or a compound's chemical activity? A clinician does not simply tell a patient "You probably have such-and-such disease. Goodbye!", but has to decide among many different kinds of treatments. The candidate drug compound may be discarded, pursued as it is, modified, and so on. *The ultimate goal of a classification is always some kind of decision*, not just a class guess. In the cat-vs-dog problem there is a natural one-one correspondence between classes and decisions. But in the medical or drug-discovery problems *the set of classes and the set of decisions are very different*, and have even different numbers of elements.

(iii) If there is a 70% probability that an image's subject is a cat, then it is natural to put it in the folder 'cats' rather than 'dogs' (if the decision is only between these two folders). If there is a 70% probability that a patient has a particular health condition, it may nonetheless be better to dismiss the patient – that is, to behave as if there was no condition. This is the optimal decision, for example, when the only available treatment for the condition would severely harm the patient if the condition were not present. Such treatment would be recommended only if the probability for the condition were much higher than 70%. Similarly, even if there is a 70% probability that a candidate drug is active it may nonetheless be best to discard it. This is the economically most advantageous choice if pursuing a false-positive leads to large economic losses. *The target of a classification is not what's most probable, but what's optimal.*  *add something on proper probabilities, connecting to (iv)*

(iv) The relation from image pixels to house-pet subject may be almost deterministic; so we are effectively extrapolating or looking for a function $\text{pet} = f(\text{pixels})$ contaminated by little noise. But the relation between medical-test scores or biochemical features on one side, and disease or drug activity on the other, is typically *probabilistic*; so a function $\text{disease} = f(\text{scores})$ or $\text{activity} = f(\text{features})$ does not even exist. *We are assessing statistical relationships* $P(\text{disease}, \text{scores})$ or $P(\text{activity}, \text{features})$ instead, which include deterministic ones as special cases.

In summary, there is place to improve classifiers so as to (i) quantitatively take into account actual utilities, (ii) separate classes from decisions, (iii) target optimality rather than ‘truth’, (iv) output and use proper probabilities.

Although we know how to address all these issues in principle – the theoretical framework is for example beautifully presented in the first 18 chapters or so of Russell & Norvig’s (2022) text¹ add also refs to [Cheeseman](#) – some of them are extremely difficult to solve. The output of sensible probabilities, issue (iv), is especially difficult. Some machine-learning algorithms for classification, such as support-vector machines, output only a class label. Others, such as deep networks, output a set of real numbers that can bear some qualitative relation to the plausibilities of the classes. But these numbers cannot be reliably interpreted as proper probabilities, that is, as the degrees of belief assigned to the classes by a rational agent¹; or, in terms of ‘populations’², as the expected frequencies of the classes in the hypothetical population of units (degrees of belief and frequencies being related by de Finetti’s theorem³). Algorithms that internally do perform probabilistic calculations, for instance naive-Bayes or logistic-regression classifiers⁴, unfortunately rest on probabilistic assumptions, such as independence and particular shapes of distributions, that are often unrealistic (and their consistency with the specific application is rarely checked). Only particular classifiers such as Bayesian neural networks⁵ output sensible probabilities, but they are computationally very expensive. The stumbling block is the extremely high dimensionality of the feature space, which makes the calculation of the probabilities

$$P(\text{class, feature} \mid \text{training data}) ,$$

(a problem opaquely called ‘density regression’ or ‘density estimation’⁶), computationally unfeasible.

If we solved the issue of outputting proper probabilities then the remaining three issues would be very easy to solve. We would simply need to enumerate the possible decisions, assess their utilities relative to the possible true classes, and finally combine utilities with probabilities

¹ MacKay 1992a; Gal & Ghahramani 2016; Russell & Norvig 2022 chs 2, 12, 13. ² Lindley & Novick 1981; Fisher 1967 § II.4. ³ Bernardo & Smith 2000 ch. 4; Dawid 2013. ⁴ Murphy 2012 § 3.5, ch. 8; Bishop 2006 §§ 8.2, 4.3; Barber 2020 ch. 10, § 17.4. ⁵ Neal & Zhang 2006; Bishop 2006 § 5.7. ⁶ Ferguson 1983; Thorburn 1986; Hjort 1996; Dunson et al. 2007.

to determine the optimal decision. The last operation amounts to a simple matrix multiplication and maximization.

✂ find a connection with imbalanced data

In the present work we propose an alternative solution to calculate proper class probabilities. It consists in a sort of ‘transducer’ or ‘calibration curve’ that transforms the algorithm’s raw output into a probability. It has a low computational cost, can be applied to all commonly used classifiers and even to simple regression algorithms, does not need any changes in algorithm architecture or in training procedures, and is grounded on first principles.

Moreover, this transducer has two great benefits. First, it allows us to calculate both the probability of class conditional on features, and *the probability of features conditional on class*. In other words it allows us to use the classification algorithm in both ‘discriminative’ and ‘generative’ modes⁷, even if the algorithm was not designed for a generative use. Second, its computation automatically also yields a ‘probability of the probability’, or more precisely a quantification of how much the probability would change if we had further training data.

The probability thus obtained is also combined with utilities to perform the final classification task.

We show on a concrete dataset that this solution always leads to an improvement in classification, in some cases with a relative increase as high as 30%.

The main idea is explained in § 2. Section ✂ ... explains how the resulting probabilities are combined with utilities to perform classification. In § ✂ ... we apply the transducer to ✂ ...

✂ finish synopsis

2 An output-to-probability transducer

In the present section we explain the main idea in informal and intuitive terms, and give its mathematical essentials only. Stricter logical derivation and more mathematical details are left to appendix ✂ ...

⁷ Russell & Norvig 2022 § 21.2.3; Murphy 2012 § 8.6.

2.1 Main idea: algorithm output as a proxy for the features

Let us consider first the essentials behind a classification (or regression) problem. We have the following quantities:

- the *feature* values of a set of known units,
- the *classes* of the same set of units,

which together form our *training data*; and

- the feature value of a new unit,

where the ‘units’ can be widgets, images, patients, drug compounds, and so on, depending on the classification problem. From these quantities we would like to infer the class of the new unit. This inference consists in probability values (we omit ‘new’ in equations for brevity)

$$P(\text{class of unit} \mid \text{feature of unit, classes \& features of training data}) \quad (1)$$

for each possible class.

These probabilities are obtained through the rules of the probability calculus⁸; in this case specifically through the so-called de Finetti theorem⁹ which connects training data and new unit.

Combined with a set of utilities, these probabilities allow us to determine an optimal, further decision to be made among a set of alternatives. Note that the inference (1) includes deterministic interpolation, i.e. the assessment of a function $\text{class} = f(\text{feature})$, as a special case, when the probabilities are essentially 0s and 1s.

A trained classifier should ideally output the probabilities above when applied to the new unit. Machine-learning classifiers trade this capability for computational speed – with an increase in the latter of several orders of magnitude¹⁰. Thus their output cannot be considered a probability, but *it still carries information about both class and feature variables*.

Now we acknowledge the fact that the information contained in the feature and in the training data, relevant to the class of the new unit, is simply inaccessible to us because of computational limitations. We do have access to the output for the new unit, however, which does

⁸ Jaynes 2003; Russell & Norvig 2022 chs 12–13; Gregory 2005; Hailperin 2011; Jeffreys 1983. ⁹ Bernardo & Smith 2000 ch. 4; Dawid 2013. ¹⁰ to understand this trade-off in the case of neural-network classifiers see e.g. MacKay 1992b,c,a; Murphy 2012 § 16.5 esp. 16.5.7; see also the discussion by Self & Cheeseman 1987.

carry relevant information. Thus want to calculate a probability such as $P(\text{class of unit} \mid \text{output for unit})$.

To calculate this probability, however, it is necessary to have examples of further pairs (class of unit, output for unit), of which the new unit's pair can be considered a representative sample and vice versa – exactly for the same reason why we need training data in the first place.

For this purpose, can we use the pairs (class of unit, output for unit) of the training data? This would be very convenient, as those pairs are readily available. But answer is no. The reason is that the outputs of the training data are produced from the features *and the classes* jointly; this is the very point of the training phase. There is therefore a direct informational dependence between the classes and the outputs of the training data. For the new unit, on the other hand, the classifier produces its output from the feature alone. The new unit is not a representative sample of the training data.

We need a data set where the outputs are generated by simple application of the algorithm to the feature, as it would occur in its concrete use, and the classes are known. The *test data* of standard machine-learning procedures are exactly what we need. The new unit can be considered a representative sample of the test data, and vice versa. We call such data 'transducer-calibration data', owing to its specific purpose.

The probability we arrive at is therefore

$$P(\text{class of unit} \mid \text{output for unit, classes \& outputs for calibration data}) . \quad (2)$$

For algorithms that yield an output much simpler than the features, such as a vector of few real components, the probability above can be exactly calculated. Thus, once we obtain the classifier's output for the new unit, we can calculate a probability for the new unit's class.

The idea above can also be informally understood in two ways. First: the classifier's output is regarded as a proxy for the feature. Second: the classifier is regarded as something analogous to a *diagnostic test*, such as any common diagnostic or prognostic test used in medicine for example. We do not take diagnostic-test results at face value – if a flu test is 'positive' we do not conclude that the patient has the flu – but rather arrive at a probability that the patient has the flu, given some statistics

about results of tests performed on ‘verified samples’ of true-positive and true-negative patients¹¹.

The probability values (4) for a fixed class and variable output give us a sort of ‘calibration curve’ (or calibration hypersurface for multidimensional outputs) of the output-to-probability transducer for the classifier. See fig. ... as an example. It must be stressed that such curve needs to be calculated only once, and it can be used for all further applications of the classifier to new units.

What is the relation between the ideal incomputable probability (1) and the probability (4) obtained by proxy? It can be shown that the two are related by convex combination or mixing:

$P(\text{class} \mid \text{output, classes \& outputs for calibration data}) =$

$$\sum_{\substack{\text{feature} \\ \text{training data}}} w(\text{feature, training data}) \times P(\text{class} \mid \text{feature, training data}) \quad (3)$$

where w are positive and normalized functions. This means that the proxy probability is generally less sharp, that is, farther away from 0 and 1, than the ideal one. This is an obvious consequence of the loss of information about the feature value. Such mixing, on the other hand, only leads to more conservative probabilities, not to over-confident ones.

 add note about the fact that the goodness of the results still greatly depends on the goodness of the classifier.

2.2 Calculation of the probabilities

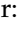
Let us denote by c the class value of a new unit, by y the output of the classifier for the new unit, and by $D := \{c_i, y_i\}$ the classes and classifier outputs for the transducer-calibration data.

Instead of the conditional probability (4), of the class given the output and calibration data, we focus instead on the joint probability of class and output given the data, which we write

$$p(c, y \mid D). \quad (4)$$

This probability is calculated using standard non-parametric Bayesian methods¹². ‘Non-parametric’ in this case means that we do not make any

¹¹ Sox et al. 2013 ch. 5; Hunink et al. 2014 ch. 5; see also Jenny et al. 2018. ¹² for introductions and reviews see e.g. Walker 2013; Müller & Quintana 2004; Hjort 1996.

assumptions about the shape of the probability curve as a function of c, y (contrast this with logistic regression, for instance), or about special independence between the variables (contrast this with naive-Bayes). The only assumption made – and we believe it is quite realistic – is that the curve must have some minimal degree of smoothness. This assumption allows for much leeway, however: fig.  ... for instance shows that the probability curve can still have very sharp bends, as long as they are not cusps.

Non-parametric methods differ from one another in the kind of ‘coordinate system’ they select on the infinite-dimensional space of all possible probability curves, that is, in the way they represent a general positive normalized function.

We choose the representation discussed by Dunson & Bhattacharya¹³. The end result of interest in the present section is that the probability density $p(c, y | D)$, with c discrete and y continuous and possibly multi-dimensional, is expressed as a sum

$$p(c, y | D) = \sum_k q_k A(c | \alpha_k) B(y | \beta_k) \quad (5)$$

of a finite but large number of terms. Each term is the product of a positive weight q_k , a probability distribution $A(c | \alpha_k)$ for c depending on parameters α_k , and a probability density $B(y | \beta_k)$ for y depending on parameters β_k . The parameter values can be different from term to term, as indicated by the index k . The weights $\{q_k\}$ are normalized.

This mathematical representation can approximate (under some norm) any smooth probability density in c and y . It has the advantages of being automatically positive and normalized, and of readily producing the marginal distributions for c and for y :

$$p(c) = \sum_k q_k A(c | \alpha_k), \quad p(y) = \sum_k q_k B(y | \beta_k), \quad (6)$$

from which also the conditional distributions are easily obtained:

$$p(c | y) = \sum_k \frac{q_k B(y | \beta_k)}{\sum_l q_l B(y | \beta_l)} A(c | \alpha_k), \quad (7a)$$

$$p(y | c) = \sum_k \frac{q_k A(c | \alpha_k)}{\sum_l q_l A(c | \alpha_l)} B(y | \beta_k). \quad (7b)$$

¹³ Dunson & Bhattacharya 2011; see also the special case presented by Rasmussen 1999.

These will be important in the following discussion.

The parametric distributions $A(c \mid \alpha)$ and $B(y \mid \beta)$ are chosen by us according to convenience. For the first we use a simple categorical distribution with parameters α . In a binary-classification case, where the class variable c assumes conventional values $\{0, 1\}$, it is

$$A(c \mid \alpha) = c \alpha + (1 - c) (1 - \alpha) \equiv \begin{cases} \alpha & \text{if } c = 1, \\ 1 - \alpha & \text{if } c = 0. \end{cases} \quad (8)$$

For the second we use a gaussian distribution, $\beta \equiv (\mu, \sigma)$ being its mean and standard deviation:

$$B(y \mid \beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(y - \mu)^2}{2\sigma^2} \right], \quad (9)$$

or a product of such gaussians if y is multidimensional. If the classifier's output y is a bounded variable we perform a transformation of its range onto the real line first (taking due care of Jacobian determinants for density transformations).

The weights $\{q_k\}$ and the parameters $\{\alpha_k\}, \{\beta_k\}$ are the heart of this representation, because the shape of the probability curve $p(c \mid y, D)$ depends on their values. They are determined by the test data D . Their calculation is done via Markon-chain Monte Carlo sampling, which we discuss in appendix [A](#).... For low-dimensional y and discrete c (or even continuous, low-dimensional c , which means we are working with a regression algorithm), this calculation can be done in a matter of hours, and *it only needs to be done once*.

Once calculated, these parameters are saved in memory and can be used to compute any of the probabilities (5), (6), (7) as needed, as discussed in the next subsection. Such computations take less than a second.

Note that the role of the classifier in this calculation is simply to produce the outputs y for the calibration data, after having been trained in any standard way on a training data set. No changes in its architecture or in its training procedure have been made, nor are any required.

2.3 Effecting the probability transduction: discriminative and generative modes

Let us now discuss the application of the various probability functions presented above. For simplicity we shall omit the dependence ‘ \dots, D' ’ on the calibration data, leaving it implicitly understood.

We have a new unit, which could be part of an evaluation test-set or coming from a real application after deployment. Its features are fed to the classifier, which outputs the real value y . We can then proceed in two ways:

Discriminative mode: *probability of class given output*

We calculate $p(c | y)$ from formula (7a), for each value of c , say $c = 0$ and $c = 1$ in a binary-classification case. These are the probabilities of the classes.

The discriminative mode is the standard way of proceeding in classification and does not need further discussion. We have simply translated the classifier’s raw output into a more sensible probability. From this point of view the function $p(c | y)$ can be considered as a more appropriate substitute of the softmax function, for instance, at the output of a neural network.

Generative mode: *probability of class given output and base rates*

We calculate $p(y | c)$ from formula (7b), again for each value of c . The probability of each class c is then obtained via Bayes’s theorem, supplying the actual *population prevalence* r_c of the class:

$$p(c | y, \text{base rates}) = \frac{p(y | c) r_c}{\sum_c p(y | c) r_c} . \quad (10)$$

The population prevalences¹⁴, also called *base rates*¹⁵, are the relative frequencies of occurrence of the various classes in the population whence our unit originates.

The generative mode, which is generally not available for classifiers with a discriminative design, is the correct way

3 Sensible probabilities for classifiers

¹⁴ Sox et al. 2013 ch. 3; Hunink et al. 2014 § 5.1. ¹⁵ Bar-Hillel 1980; Axelsson 2000.

🔧 only pieces will be taken from this section

Why are probability values important? As we argue in a companion work 🔧🔧, our ultimate purpose in classification is seldom only to guess a class; most often it is to choose a specific course of action, or to make a decision, among several available ones. A clinician, for example, does not simply tell a patient “you will probably not contract the disease”, but has to decide among dismissal or different kinds of preventive treatment¹⁶. Said otherwise, in classification we must choose the *optimal* class, not the probably true one. Making optimal choices in situations of uncertainty is the domain of Decision Theory¹⁷. In order to make an optimal choice, decision theory requires the use of probability values that properly reflect our state of uncertainty.

Determining class probabilities conditional on the input features is unfortunately computationally unfeasible at present for problems that involve very high-dimensional spaces, such as image classification; in fact if an exact probabilistic analysis were possible we would not be developing machine-learning classifiers in the first place¹⁸. 🔧🔧 Maybe useful to add a reminder that probability theory is the *learning* theory par excellence (even if there’s no ‘learning’ in its name)? Its rules are all about making logical updates given new data.

In the present work we propose an alternative solution that has a low computational cost and that can be applied to all commonly used classifiers, even those that only output class labels.

The essential idea comes from seeing an analogy between a classifier and a diagnostic test, such as any common diagnostic or prognostic test used in medicine for example. There are many parallels in the way machine-learning classifiers and diagnostic tests, a flu test for instance, are devised and work. Our basic motivation in using either is that we would like to assess some situational variable – class, pathological condition – by means of its correlation (in the general sense of the word, not the linear Pearson one; and including deterministic dependence as a particular case) with a set of ‘difficult’ variables that are either too complex or hidden – image pixels, presence of replicating viral agents –:

situational variable \longleftrightarrow difficult variables

We devise an auxiliary variable – algorithm output, test result – to be correlated with the difficult variables:

situational variable \longleftrightarrow difficult variables \longleftrightarrow aux variable

¹⁶ Sox et al. 2013; Hunink et al. 2014. ¹⁷ Russell & Norvig 2022 ch. 15; Jeffrey 1965; North 1968; Raiffa 1970. ¹⁸ Russell & Norvig 2022 chs 2, 12; Pearl 1988.

We can now assess the situational variable by observing the more easily accessible auxiliary variable instead of the difficult ones. In probability language we are *marginalizing* over the difficult variables. This is the procedure dictated by the probability calculus whenever we do not have informational access to a set of variables. The correlation of the auxiliary variable is achieved by the training process in the case of the machine-learning algorithm, and by the exploitation of biochemical processes or reactions in the case of the flu test.

The situational variable is *informationally screened* from the auxiliary variable by the difficult variables. That is, the auxiliary variable does not – in fact, cannot – contain any more information about the situational variable than that contained in the difficult variables. This means that the probability relationship between the three variables is as follows:

$$p\left(\begin{array}{c} \text{situational} \\ \text{variable} \end{array} \middle| \begin{array}{c} \text{aux} \\ \text{variable} \end{array}\right) = \sum_{\text{difficult variables}} p\left(\begin{array}{c} \text{situational} \\ \text{variable} \end{array} \middle| \begin{array}{c} \text{difficult} \\ \text{variables} \end{array}\right) \times p\left(\begin{array}{c} \text{difficult} \\ \text{variables} \end{array} \middle| \begin{array}{c} \text{aux} \\ \text{variable} \end{array}\right), \quad (11)$$

the sum running over all possible values of the difficult variables.

In the case of the diagnostic test we determine the probability $p\left(\begin{array}{c} \text{situational} \\ \text{variable} \end{array} \middle| \begin{array}{c} \text{aux} \\ \text{variable} \end{array}\right)$ by carrying out the test on a representative sample of cases and collecting joint statistics between the test's output and the true situation, the presence of the flu in our example. These statistics are typically displayed in a so-called contingency table¹⁹, akin to a confusion matrix.

Unlike the case of a diagnostic test, the output of a machine-learning classifier is usually taken at face value: if the output is a class label, that label is regarded as the true class; if the output is a unity-normalized tuple of positive numbers, that tuple is regarded as the probability distribution for the classes.

We instead propose *to treat the classifier's output just like a diagnostic test's result*. This output, discrete or continuous, is regarded as a quantity that has some correlation with the true class. This correlation can be analysed in a set of representative samples and used to calculate a sensible probability for the class given the classifier's output. This analysis only

¹⁹ Fienberg 2007; Mosteller et al. 2013.

needs to be made once and is computationally cheap, because the classifier output takes values in a discrete or low-dimensional space.

This approach differs from the computationally infeasible one discussed above in that we are calculating the computationally easier probability

$$p(\text{class} \mid \text{output}) \quad (12)$$

rather than

$$p(\text{class} \mid \text{feature}) . \quad (13)$$

The former probability, as we saw in eq. (11), is the marginal

$$p(\text{class} \mid \text{output}) = \sum_{\text{feature}} p(\text{class} \mid \text{feature}) \times p(\text{feature} \mid \text{output}) . \quad (14)$$


We can thus think of this approach as a marginalization over the possible features, which is a necessary operation as we have no effective access to them.

A hallmark of this approach is that we are calculating exact probabilities conditional on reduced information, rather than approximate probabilities conditional on full information. This protects us from biases that are typically present in the approximation method. The price of using reduced information is that the probabilities may be open to more variability as we collect more representative data. But as we shall see this variability is actually quite low, and moreover it can be exactly assessed.


This approach also offers the following advantages:

- It does not require any changes of the standard training procedures.
- It is easily implemented as an additional low-cost computation of a function at the end of the classifier's output, or as a replacement of a softmax-like computation.
- It does not make any assumptions such as linearity or gaussianity.
- It yields not only the probability distribution for the classes, but also a measure of how much this distribution could change if we collected more test data (the 'probability of the probability', so to speak).
- It allows us to use the classifier both in a discriminative and generative way. That is, we can use either $p(\text{class} \mid \text{output})$, or $p(\text{output} \mid \text{class})$ in conjunction with Bayes's theorem. The latter approach enables us to avoid possible base-rate fallacies²⁰.

²⁰ Russell & Norvig 2022 § 12.5; Axelsson 2000; Jenny et al. 2018.

- It can be seamlessly integrated with a utility matrix to compute the optimal class, as shown in the companion work .

In § 4 we present some notation and the general procedure for the calculation of the probabilities; more technical details are given in appendix 10. Section 5 explains how to augment a classifier's output with the probability calculation. Results of numerical experiments are presented in § 6.

 We could show that even if we used a biased test set, the method corrects the bias (provided we know what the bias is).

4 Calculation of the probabilities: general procedure

Let us denote the class variable by C and the classifier-output variable by X . We assume that C is discrete and finite, its values can be simply renamed $1, 2, 3, \dots$. We assume that X is either discrete and finite (it is isomorphic to C for many classifying algorithms) or a low-dimensional tuple of real variables; a combination of both cases can also be easily accommodated. We assume to have a sample of M such data pairs denoted collectively by D :

$$D := \{(c_1, x_1), (c_2, x_2), \dots, (c_M, x_M)\}. \quad (15)$$

We call them *calibration data*. Let us emphasize that these are not pairs of class & feature values, but pairs of class & classifier-output values, obtained as described in § 5.


Instead of the conditional probability $p(\text{class} \mid \text{output})$, that is, $p(C \mid X)$, we can actually calculate the joint probability

$$p(C, X) \quad (16)$$

given the sample data. The computational cost is the same, but from the joint probability we can easily derive both conditionals

$$p(C \mid X) = \frac{p(C, X)}{\sum_C p(C, X)}, \quad (17)$$

$$p(X \mid C) = \frac{p(C, X)}{\sum_X p(C, X)}. \quad (18)$$

It is advantageous to have both, as we shall see in § : if one of them is biased owing to the way the test samples were obtained, we can still use the other.

In our specific inference problem, where no time trends are assumed to exist in future data (the probability distribution for future data is exchangeable), probability theory dictates that the joint probability (16) for a new datapoint (c_0, x_0) is equal to the *expected* frequency of that datapoint in a hypothetically infinite run of observations, that is, the average

$$p(c_0, x_0) = \int F(c_0, x_0) w(F) dF . \quad (19)$$


This formula is the so-called de Finetti's theorem²¹. It is derived from first principles but can be intuitively interpreted: We are considering every possible long-run frequency distribution $F(\cdot, \cdot)$, giving it a weight $w(F)$, and then taking the weighted sum of all such distributions. The result is still a distribution, and its value at (c_0, x_0) is the probability of this datapoint.

The weight $w(F)$ – a probability density – given to a frequency distribution F is proportional to two factors:

$$w(F) \propto F(D) w_g(F) . \quad (20)$$

- The first factor ('likelihood') $F(D)$ quantifies how well F fits known data of the same kind, the sample data D in our case. It is simply proportional to how frequent the known data would be, according to F :

$$F(D) = F(c_1, x_1) F(c_2, x_2) \cdots F(c_M, x_M) . \quad (21)$$

- The second factor ('prior') $w_g(F)$ quantifies how well F generalizes beyond the data we have seen, owing to reasons such as physical or biological constraints for example. In our case we expect F to be somewhat smooth in X when this variable is continuous²²  add a picture – a sample from the prior over F – to illustrate the expected range of smoothness. No assumptions are made about F when X is discrete.

Formula (20) is just Bayes's theorem. Its normalization factor is the integral $\int F(D) w_g(F) dF$, which ensures that $w(F)$ is normalized.

The first factor becomes larger as the number of known data increases. Thus a large amount of data indicating a non-smooth distribution F will override any smoothness preferences embodied in the second factor.

²¹ Bernardo & Smith 2000 ch. 4; Dawid 2013; de Finetti 1929; 1937. ²² Cf. Good & Gaskins 1971.

Note that no assumptions about the shape of F – no gaussians, logistic curves, sigmoids, and so on – are made in this approach.

The integral in (19) is calculated in either of two ways, depending on whether X is discrete or continuous. For X discrete and taking on the same values as the class variable C , the integral is over \mathbf{R}^{n_c} where n_c is the number of possible classes, and can be done analytically. For X continuous, the integral is numerically approximated by a sum over a representative sample, obtained by Markov-chain Monte Carlo, of distributions F according to the weights (20). The error of this approximation can be calculated and made as small as required by increasing the number of Monte Carlo samples.

The expected value (19) is calculated for all possible values of (c_0, x_0) , obtaining the full joint probability distribution $p(C, X)$. From this joint distribution we calculate the direct and inverse conditional distributions

$$p(C | X) = \frac{p(C, X)}{\sum_C p(C, X)} , \quad (22)$$

$$p(X | C) = \frac{p(C, X)}{\sum_X p(C, X)} . \quad (23)$$

It is very convenient to have both, as discussed in § 7.

The conditional distributions above are just matrices when X is discrete. For continuous X they can be regarded as n_c tuples of functions in X . We can find convenient approximate expressions, such as polynomial interpolants, for faster numerical implementations of these functions.

The integration procedure for (19) also tells us how much the probability distribution $p(C, X)$ would change if we acquired new data (a sort of ‘probability of the probability’).

For further mathematical details see appendix 10.

5 Implementation in the classifier output

The implementation of our approach takes place after the training of the classifier has been carried out in the usual way. We assume that a collection of M test data were set aside as usual:

$$T := \{(c_1, z_1), (c_2, z_2), \dots, (c_M, z_M)\} , \quad (24)$$

where the c_i are the true classes and z_i the corresponding feature values.

The M feature values z_i are given as inputs to the classifier, which produces M corresponding outputs x_i . We now consider data pairs consisting in the true classes c_i and the outputs x_i : these are the *calibration data* discussed in § 4:

$$D := \{(c_1, x_1), (c_2, x_2), \dots, (c_M, x_M)\}.$$


They are used to find the direct and inverse conditional probability distributions $p(C | X)$ and $p(X | C)$ as described in § 4.

We can finally augment our classifier either in a ‘direct’ or ‘discriminative’ way, or an ‘inverse’ or ‘generative’ way, by adding one computation step at the end of the classifier’s operation:

Direct: from its output x_0 we obtain the probability for each class, $p(c | x_0)$.

Inverse: from its output x_0 we obtain the probability of the output itself, conditional on each class, $p(x_0 | c)$.

These n_c probabilities are the final output of the augmented classifier.

In the direct or discriminative case, at each new use of the classifier the output probabilities can be used together with a utility matrix to choose the *optimal* class for that case, as discussed in the companion paper .

In the inverse or generative case, at each new use of the classifier the probabilities for the classes are obtained via Bayes’s theorem:


$$p(c) = \frac{p(x_0 | c) B(c)}{\sum_c p(x_0 | c) B(c)}, \quad (25)$$

where $B(c)$ is the base rate of class c . The probabilities $p(c)$ can finally be used together with a utility matrix to choose the *optimal* class.

6 Numerical experiments and results

7 Circumventing biases

8 Alternative to ensembling

 The idea is to implement this output-to-probability conversion in several classifiers, in a *generative way*. These probabilities can then be multiplied together and with a base rate, to obtain the ‘ensembled’ probabilities of the classes. This way of doing ensembling would

be a rigorous application of the probability calculus; should be superior to a majority vote or similar.

9 Methods and materials

9.1 Data

The data comes from a publicly open bioactivity database called ChEMBL **ChEMBL**. The dataset that is used in this paper was introduced in the paper by Koutsoukas et. al. **DL-investigating-bioactivity**. The dataset is a source of structure-activity relationship from version 20 of ChEMBL with protein target Carbonic Anhydrase II (ChEMBL205).

9.2 Pre-processing

For our pre-processing pipeline, we use two different methods to represent the molecule, one for the RF and one for the DNN. The first method turns the molecule into a hashed bit vector of circular fingerprints called ECFP (Extended Connectivity Fingerprints) **Rogers-ECFP**. From our experiments, there was little to no improvement using a 2048-bit vector over a 1024-bit vector.

For our DNN model, the data is represented by converting the molecule into images in 224x224 size. This is done by using the SMILES string from the dataset and using Rdkit **rdkit** to generate a canonical graph structure of the molecule. This differs from ECFP in that it represents the actual structure of the molecule rather than properties generated from the molecule.

The dataset has in total 1631 active compounds and 16310 non-active molecules which act as decoys. In training, the active molecules are oversampled to match the same number of non-active molecules.

9.3 Prediction

Virtual screening is the process of finding chemical activity in the interaction between a compound (molecule) and a target (protein). The goal of the machine learning algorithms is to find structural features or chemical properties that show that the molecule is active against the protein. DNNs have previously been shown to outperform RF and other

linear models in VHTS and quantitative structure-activity relationship (QSAR) problems **DL-investigating-bioactivity**.

9.4 Models

The algorithms and methods used to create the models have previously been shown to give great results for all different fields.

The second model is a pre-trained residual network (ResNet) **resnet** with 18 hidden layers trained on the well-known ImageNet dataset **ImageNet** by using the PyTorch framework **PyTorch**. ResNet has shown to outperform other pre-trained CNN models **resnet**. A ResNet with 34 hidden layers showed little to no performance gain, so we chose to go with the simpler model. The model is trained with the following hyperparameters:

- **Learning rate:** 0.003
- **Optimization technique:** Stochastic Gradient Descent (SGD)
- **Activation Function:** ReLU
- **Dropout:** 50%
- **Number of epochs:** 20
- **Loss function:** Cross Entropy Loss


9.5 Train, Validation, Test split

The data is split into four parts:

- Test set 1: Uses 20% of the total dataset for the Bayesian inference.
- Test set 2: Uses 20% of the dataset for evaluation.
- Validation set: 25% of the remaining data for validating the model after each epoch.
- Training set: The rest of the dataset is used to train the model.

10 Mathematical details of the nonparametric density regression

The joint probability (16) is calculated nonparametrically, that is, without making any assumptions such as linearity or gaussianity, besides very mild and reasonable assumptions of continuity.

using the versatile computational approach by Dunson & Bhattacharya (2011), and obtain the probability (12) by conditionalization. The calculation requires Monte Carlo sampling  refs [here](#) but needs to be made only once.

11 Summary and discussion

 Add note about how (sequential) decision theory was used during World War I; see Good (1950) around § 6.2

PIECES OF TEXT

12 An example

 will delete this section

An example. Imagine you’re a clinician and must attend to a patient with a particular disease. The disease may appear in two variants *I* and *II*; you are not sure which type affects your patient. For this disease there are three kinds of treatment *A*, *B*, *C* available at present, and you must choose one of them. Their efficacies, measured on some scale, depend on the disease type according to the following table:

		disease type	
		<i>I</i>	<i>II</i>
treatment	<i>A</i>	2	−2
	<i>B</i>	1	1
	<i>C</i>	−2	2

Treatment *A* is very effective for disease type *I* but causes harm (hence the negative value) if administered to a patient of type *II*; vice versa for treatment *C*. Treatment *B* never causes harm but is less effective on either disease type.

Which treatment do you choose?

You could say “the answer depends on the patient’s disease type”. This would be correct if you knew the disease type: type *I*, choose treatment *A*; type *II*, choose treatment *C*. But you do *not* know the disease type, so cannot make your answer depend on unavailable knowledge.

Rather, the answer depends on *how sure you are about the disease type*, given whatever evidence you have. This is clear in the case where you are completely uncertain about the type (and have no way to dispel your uncertainty), which could equally be either way. The rational, safest choice in this case is treatment *B*: the patient will have some relief for sure, albeit not the highest, and no harm will be done. If you are *extremely sure* that the disease type is *I* instead, then you recommend your patient to go for treatment *A*, as there’s low risk for harm and greater benefits to be reaped (surgical treatments are typical examples of this case: there’s always some minimal risk that a surgery goes awry, but such risk may be offset by the benefits of the more likely successful surgery). Analogously if you are sure about type *II* instead, recommending treatment *C*.

Now suppose that a shiny new machine-learning algorithm has been acquired by your clinic to help diagnosing the disease type quickly and with minimal expenses. This algorithm takes as input the results of various cheap clinical tests performed on a patient, and outputs the disease type. You use it for your patient, it outputs '*I*', so you go for treatment *A*, and the patient indeed gets better. You use it for two more patients. The outputs are *II* and *I*, so you administer treatments *C* and *A*. Both patients get better. Bliss! You use the algorithm for a third patient. Clinical tests are again made and their results fed into the algorithm. It outputs '*II*'. You therefore administer treatment *C*. But the patient is badly harmed and gets worse (and possibly sues you). To check for human error, the tests are repeated; the algorithm consistently outputs '*II*'. Eventually you perform an expensive and invasive but reliable test. It turns out the disease type is *I*, not *II*. The algorithm is simply wrong, and you chose the worst treatment.

You ask the algorithm's vendor or developers for an explanation: "isn't the algorithm supposed to tell me the actual disease type?". They tell you that no, there's always a chance that the algorithm is wrong, depending on the case. For some input values (test results) the prediction is virtually certain, but for others there may be a larger margin of error. You tell them:

Sorry, but there has been a misunderstanding then. I don't need an algorithm that gives me a best guess, making me sometimes almost kill my patients. *I need an algorithm that tells me how sure a disease type is.* Because if there's much uncertainty I can give my patients a treatment that's guaranteed harm-free even if not the best!

This example exposes several misconceptions and issues in how classification problems and also some regression problems are often stated and faced in machine learning:

(i) The ultimate goal in a classification is often not the guess a class, but the choice among a set of alternative decisions. In a 'cat vs dog' image classification the classes are 'cat' and 'dog', and the decisions might be 'put into folder Cats' vs 'put into folder Dogs'. A clinician, however, does

not simply tell a patient “you probably have such-and-such disease”, but has to decide among different kinds of treatment²³.

(ii) The alternative decisions, for example treatments in medicine, often do not correspond to the unknown classes in a one-one way; they may be more numerous than the classes.

(iii) The target is not what’s most probable, but what’s *optimal*. The two may be different. In the clinical example above, the clinician and patients would opt for treatment B , rather than A , even if type I had slightly higher probability than II – say 60% vs 40% (we shall show later that the threshold for this case is 70%). This is also true when the decisions have some natural one-one correspondence with the classes. Consider for example two other treatments A' , B' with this efficiency table:

	disease type	
	I	II
treatment A'	2	-2
treatment B'	1	1

Clearly A' is best for type I , and B' for type II , but we would choose A' only if type I had quite a higher probability than II (the threshold is again 70%).

12.1 Actual utility yield

The utility matrix is not only the basis for making optimal decisions by means of expected-utility maximization. It also provides the metric to rank a set of decisions already made – for example on a test set – by some algorithm, if we know the corresponding true classes. Suppose we have N test instances, in which each class c occurs N_c times, so that $\sum_c N_c = N$. A decision algorithm made decision d when the true class was c a number M_{dc} of times. These numbers form the confusion matrix (M_{dc}) of the algorithm’s output. The numbers M_{dc} are must satisfy the constraints $\sum_d M_{dc} = N_c$ for each c .

For given decision d and class c , in each of the M_{dc} instances the algorithm yielded a utility U_{dc} . The actual average utility yield in the test set is then

$$\frac{1}{N} \sum_{dc} U_{dc} M_{dc} . \quad (26)$$

²³ Sox et al. 2013; Hunink et al. 2014.

It is convenient to consider the average utility yield, rather than the total utility yield (without the $1/N$ factor), because if we shift the zero or change the measurement unity of our utilities then the yield changes in the same way.

The summary of decision theory just given suffices to address issues **i1–i4**.

In comparing, evaluating, and using machine-learning classifiers we face a number of questions and issues; some are well-known, others are rarely discussed:

- i1 Choice of valuation metric.** When we have to evaluate and compare different classifying algorithms or different hyperparameter values for one algorithm, we are avalanched by a choice of possible evaluation metrics: accuracy, area under curve, F_1 -measure, mean square contingency²⁴ also known as Matthews correlation coefficient²⁵, precision, recall, sensitivity, specificity, and many others²⁶. Only vague guidelines are usually given to face this choice. Typically one computes several of such scores and hopes that they will lead to similar ranking.
- i2 Rationale and consistency.** Most or all of such metrics were proposed only on intuitive grounds, from the exploration of specific problems and relying on tacit assumptions, then heedlessly applied to new problems. The Matthews correlation coefficient, for example, relies on several assumptions of gaussianity²⁷, which for instance do not apply to skewed population distributions²⁸. The area under the receiver-operating-characteristic curve is heavily affected by values of false-positive and false-negative frequencies, as well as by misclassification costs, that have nothing to do with those of the specific application of the classifier²⁹. The F_1 -measure implicitly gives correct classifications a weight that depends on their frequency or probability³⁰; such dependence amounts to saying, for

²⁴ Yule 1912 denoted ' r ' there. ²⁵ Matthews 1975; Fisher 1963 § 31 p. 183. ²⁶ Sammut & Webb 2017; see also the analysis in Goodman & Kruskal 1954; 1959; 1963; 1972.

²⁷ Fisher 1963 § 31 p. 183 first paragraph. ²⁸ Jeni et al. 2013; Zhu 2020. ²⁹ Baker & Pinsky 2001; Lobo et al. 2008. ³⁰ Hand & Christen 2018.

example, “this class is rare, *therefore* its correct classification leads to high gains”, which is a form of scarcity cognitive bias³¹.

We are therefore led to ask: are there valuation metrics that can be proven, from first principles, to be free from biases and unnecessary assumptions?

- i3 Class imbalance.** If our sample data are more numerous for one class than for another – a common predicament in medical applications – we must face the ‘class-imbalance problem’: the classifier ends up classifying all data as belonging to the more numerous class³², which may be an undesirable action if the misclassification of cases from the less numerous class entails high losses. 🛠️ [discussion and refs about cost-sensitive learning](#)

i4 Optimality vs truth.

All the issues above are manifestly connected: they involve considerations of importance, gain, loss, and of uncertainty.

In the present work we show how issues **i1–i4** are all solved at once by using the principles of *Decision Theory*. Decision theory gives a logically and mathematically self-consistent procedure to catalogue all possible valuation metrics, to make optimal choices under uncertainty, and to evaluate and compare the performance of several decision algorithms. Most important, we show that implementing decision-theoretic procedures in a machine-learning classifier does not require any changes in current training practices 🛠️ (possibly it may even make procedures like under- or over-sampling unnecessary!), is computationally inexpensive, and takes place downstream after the output of the classifier.

The use of decision theory requires sensible probabilities for the possible classes, which brings us to issue ?? above. In the present work we also present and use a computationally inexpensive way of calculating these probabilities from the ordinary output of a machine-learning classifier, both for classifiers such as 🛠️ [example here](#) that can only output a class label, and for classifiers that can output some sort of continuous score.

🛠️ Write here a summary or outlook of the rest of the paper and a summary of results:

- The admissible valuation metrics for a binary classifier form a two-dimensional family;

³¹ Camerer & Kunreuther 1989; Kim & Markus 1999; Mittone & Savadori 2009. ³² Sammut & Webb 2017; Provost 2000.

that is, the choice of a specific metric corresponds to the choice of two numbers. Such choice is problem-dependent and cannot be given a priori. • Admissible metrics are only those that can be expressed as a linear function of the elements of the population-normalized confusion matrix. Metrics such as the F_1 -measure or the Matthews correlation coefficient are therefore inadmissible

13 Classification from the point of view of decision theory

In using machine-learning classifiers one typically considers situations where the set of available decisions and the set of possible classes have some kind of natural correspondence and equal in number. In a ‘cat vs dog’ image classification, for example, the classes are ‘cat’ and ‘dog’, and the decisions could be ‘put into folder Cats’ vs ‘put into folder Dogs’. In a medical application the classes could be ‘ill’ and ‘healthy’ and the decisions ‘treat’ vs ‘dismiss’. In the following when we speak of ‘classification’ we mean a *decision* problem of this kind. The number of decisions thus equals that of classes: $n_d = n_c$.

✚ For simplicity we will focus on binary classification, $n_d = n_c = 2$, but the discussion generalizes to multi-class problems in an obvious way.

13.1 Choice of valuation metric, rationale and consistency (issues i1, i2)

13.2 Optimality vs truth (issue i4)

According to decision theory a classification algorithm should, at each application, calculate the probabilities ($p_{c|z}$) for the possible classes, given the feature z provided as input; calculate the expected utility of the available decisions according to eq. (??), using the probabilities and the utility matrix; and finally output the decision d^* having maximal expected utility:

$$d^* = \arg \max_d \sum_c U_{dc} p_{c|z} . \quad (27)$$


We assume here that the utilities are given and the same at each application – the latter assumption could be dropped, however; see the discussion in § 11.

Current common practice with algorithms capable of outputting some sort of probability-like score is simply to output the class c^* having highest probability:

$$c^* = \arg \max_c p_{c|z} . \quad (28)$$

As discussed in § 3, issue i4, this means choosing the *most probable* class, not the *optimal* class, and the two are often different, the second being what we typically want. This choice is also the one that would be made with an identity utility matrix $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$.

How can we amend current practice for this kind of classifiers, so that they look for optimality rather than truth?

A first idea could be to simply modify the standard output step (28) into (27). It is an easily implementable and computationally cheap modification: we just multiply the probability tuple by a matrix. Such simple modification, however, has a profound implication for the training procedure: we are modifying the algorithm to output the optimal class, and therefore it should also *learn what is optimal*, not what is true: *the targets in the training and validation phases should be the optimal classes, not the true classes*. But optimality depends on the value of sensible probabilities for the specific situation of uncertainty, in this case conditional on the input features. Determining the optimal classes would thus require a probabilistic analysis that is computationally unfeasible at present for problems that involve very high-dimensional spaces, such as image classification – if an exact probabilistic analysis were possible we would not be developing machine-learning classifiers in the first place³³.  Maybe useful to add a reminder that probability theory is the *learning* theory par excellence (even if there's no 'learning' in its name)? Its rules are all about making logical updates given new data.

Appendix: broader overview of binary classification

Let us consider our binary-classification problem from a general perspective and summarize how it would be approached and solved from first principles³⁴ if our computational resources had no constraints.

In our long-term task we will receive 'units' of a specific kind; the units for example could be gadgets, individuals, or investment portfolios. Each new unit will belong to one of two classes, which we can denote $X=0$ and $X=1$; for example they could be 'defective' vs 'non-defective', 'ill' vs 'healthy'. The class will be unknown to us. For each new unit we shall need to decide among two possible actions, which we can denote $A=\hat{0}$ and $A=\hat{1}$; for example 'discard' vs 'keep', or 'treat' vs 'dismiss'.

³³ Russell & Norvig 2022 chs 2, 12; Pearl 1988. ³⁴ Russell & Norvig 2022 part IV.

The utility of each action depends on the unknown class of the unit; we denote these utilities by $U(A | X)$. For each new unit we will be able to measure a ‘feature’ Z of a specific kind common to all units; for example Z could be a set of categorical and real quantities, or an image such as a brain scan. We have a set of units – our ‘sample units’ or ‘sample data’ – that are somehow “representative” of the units we will receive in our long-term task³⁵. we know both the class and the feature of each of these sample units. Let us denote this sample information by D .

According to the principles of decision theory and probability theory, for each new unit we would proceed as follows:

1. Assign probabilities to the two possible values of the unit’s class, given the value of the unit’s feature $Z=z$, our sample data D , and any other available information:

$$p(X=0 | Z=z, D), \quad p(X=1 | Z=z, D) \equiv 1 - p(X=0 | Z=z, D), \quad (29)$$

according to the rules of the probability calculus.

2. Calculate the expected utilities \bar{u} of the two possible actions:

$$\begin{aligned} \bar{u}(\hat{0}) &:= U(\hat{0} | X=0) p(X=0 | Z=z, D) + U(\hat{0} | X=1) p(X=1 | Z=z, D) \\ \bar{u}(\hat{1}) &:= U(\hat{1} | X=0) p(X=0 | Z=z, D) + U(\hat{1} | X=1) p(X=1 | Z=z, D) \end{aligned} \quad (30)$$

and choose the action having maximal expected utility.

How is the probability $p(X | Z=z, D)$ determined by the probability calculus? Here is a simplified, intuitive picture. First consider the case where the feature Z can only assume a small number of possible values, so that many units can in principle have the same value of Z .

Consider the collection of all units having $Z=z$ that we received in the past and will receive in the future. Among them, a proportion $F(X=0 | Z=z)$ belong to class 0, and a proportion $1 - F(X=0 | Z=z) \equiv F(X=1 | Z=z)$ to class 1. For example these two proportions could be 74% and 26%. Our present unit with $Z=z$ is a member of this collection. The probability $p(X=0 | Z=z)$ that our unit belongs to class 0, given that its feature has value z , is then intuitively equal to the proportion $F(X=0 | Z=z)$. Analogously for $X=1$.

³⁵ for a critical analysis of the sometimes hollow term ‘representative sample’ see Kruskal & Mosteller 1979a,b,c; 1980.

The problem is that we do not know the proportion $F(X=0 | Z=z)$. However, we expect it to be roughly equal to the analogous proportion seen in our sample data; let us denote the latter by $F_s(X=0 | Z=z)$:

$$F(X=0 | Z=z) \sim F_s(X=0 | Z=z) . \quad (31)$$

this is indeed what we mean by saying that our sample data are ‘representative’ of the future units. Later we shall discuss the case in which such representativeness is of different kinds. We expect the discrepancy between $F(X=0 | Z=z)$ and $F_s(X=0 | Z=z)$ to be smaller, the larger the number of sample data. Vice versa we expect it to be larger, the smaller the number of sample data.

If Z can assume a continuum of values, as is the case for a brain scan for example, then the collection of units having $Z=z$ is more difficult to imagine. In this case each unit will be unique in its feature value – no two brains are exactly alike.



old text below

Given the unit’s feature Z we will assign probabilities to the possible values of the unit’s class: according to the rules of the probability calculus.

As mentioned in § ??, a decision problem under uncertainty is conceptually divided into two steps

The Suppose we have a population of units or individuals characterized by a possibly multidimensional variable Z and a binary variable $X \in \{0, 1\}$. Different joint combinations of (X, Z) values can appear in this population. Denote by $F(X=x, Z=z)$, or more simply $F(x, z)$ when there is no confusion, the number of individuals having specific joint values $(X=x, Z=z)$. This is the absolute frequency of the values (x, z) . We can also count the number of individuals having a specific value of $Z=z$, regardless of X ; this is the marginal absolute frequency $F(z)$. It is easy to see that

$$F(z) = F(X=0, z) + F(X=1, z) \equiv \sum_x F(x, z) . \quad (32)$$

Analogously for $F(x)$.

Select only the subpopulation of individuals that have a specific value $Z=z$. In this subpopulation, the *proportion* of individuals having

a specific value $X = x$ is $f(x | Z = z)$. This is the conditional relative frequency of x given that z . It is easy to see that

$$f(x | z) = \frac{F(x, z)}{F(z)}. \quad (33)$$

Now suppose that we know all these statistics about this population. An individual coming from this population is presented to us. We measure its Z and obtain the value z . What could be the value of X for this individual? We know that among all individuals having $Z = z$ (and the individual before us is one of them) a proportion $f(x | z)$ has $X = x$. Thus we can say that there is a probability $f(x | z)$ that our individual has $X = x$. And this is all we can say if we only know Z .

For this individual we must choose among two actions $\{a, b\}$. The utility of performing action a if the individual has $X = x$, and given any other known circumstances, is $U(a | x)$; similarly for b . If we knew the value of X , say $X = 0$, we would simply choose the action leading to maximal utility:

$$\begin{aligned} &\text{if } U(a | X=0) > U(b | X=0) \text{ then choose action } a, \\ &\text{if } U(a | X=0) < U(b | X=0) \text{ then choose action } b, \\ &\text{else it does not matter which action is chosen.} \end{aligned} \quad (34)$$

But we do not know the actual value of X . We have probabilities for the possible values of X given that $Z = z$ for our individual. Since X is uncertain, the final utilities of the two actions are also uncertain; but we can calculate their *expected* values $\bar{U}(a | Z = z)$ and $\bar{U}(b | Z = z)$:

$$\begin{aligned} \bar{U}(a | z) &:= U(a | X=0) f(X=0 | z) + U(a | X=1) f(X=1 | z), \\ \bar{U}(b | z) &:= U(b | X=0) f(X=0 | z) + U(b | X=1) f(X=1 | z). \end{aligned} \quad (35)$$

Decision theory shows that the optimal action is the one having the maximal expected utility. Our choice therefore proceeds as follows:

$$\begin{aligned} &\text{if } \bar{U}(a | z) > \bar{U}(b | z) \text{ then choose action } a, \\ &\text{if } \bar{U}(a | z) < \bar{U}(b | z) \text{ then choose action } b, \\ &\text{else it does not matter which action is chosen.} \end{aligned} \quad (36)$$

The decision procedure just discussed is very simple and does not need any machine-learning algorithms. It could be implemented in a

simple algorithm that takes as input the full statistics $F(X, Z)$ of the population, the utilities, and yields an output according to (36).

Our main problem is that the full statistics $F(X, Z)$ is almost universally not known. Typically we only have the statistics $F_s(X, Z)$ of a sample of individuals that come from the population of interest or from populations that are somewhat related to the one of interest. This is where probability theory steps in. It allows us to assign probabilities to all the possible statistics $F(X, Z)$. From these probabilities we can calculate the *expected* value $\bar{f}(x | z)$ of the conditional frequencies $f(x | z)$. Decision theory says that the expected value $\bar{f}(x | z)$ should then be used, in this uncertain case, in eq. (35) in place of the unknown $f(x | z)$. The decision procedure (36) can then be used again.

Probability theory says that in this particular situation the probability of a particular possible statistics $F(X, Z)$ is the product of two factors having intuitive interpretations:

- the probability of observing the statistics $F_s(X, Z)$ of our data sample, assuming the full statistics to be $F(X, Z)$. With some combinatorics it can be shown that this probability is proportional to

$$\exp \left[\sum_{X, Z} F_s(X, Z) \ln F(X, Z) \right] \quad (37)$$

The argument of the exponential is the cross-entropy between $F_s(X, Z)$ and $F(X, Z)$; this is the reason of its appearance in the loss function used for classifiers³⁶.

This factor tells us how much the possible statistics *fit* the sample data; it gives more weight to statistics with a better fit.

- the probability of the full statistics $F(X, Z)$ for reasons not present in the data, for example because of physical laws, biological plausibility, or similar.

This factor tells us whether the possible statistics should be favourably considered, or maybe even discarded instead, for reasons that go beyond the data we have seen; in other words, whether the hypothetical statistics would *generalize* well beyond the sample data.

³⁶ Bridle 1990; MacKay 1992a.

The final probability comes from the balance between these ‘fit’ and ‘generalization’ factors. Note that the first factor becomes more important as the sample size and therefore $F_s(X, Z)$ increases; the sample data eventually determine what the most probable statistics is, if the sample is large enough.

A similar probabilistic reasoning applies if our sample data come not from the population of interest but from a population having at least the same *conditional* frequencies of as the one of interest, either $f(X | Z)$ or $f(Z | X)$. The latter case must be examined with care when our purpose is to guess X from Z . In this case we cannot use the conditional frequencies $f_s(X | Z)$ that appear in the data to obtain the expected value $\bar{f}(X | Z)$: they could be completely different from the ones of the population of interest. We must instead use the sample conditional frequencies $f_s(Z | X)$ to obtain the expected value $\bar{f}(Z | X)$, and then combine the latter with an appropriate probability $P(X)$ through Bayes’s theorem:

$$\frac{\bar{f}(Z | X) P(X)}{\sum_X \bar{f}(Z | X) P(X)} . \quad (38)$$

The probability $P(X)$ cannot be obtained from the data, but requires a separate study or survey. In medical applications, where X represents for example the presence or absence of a disease, the probability $P(X)$ is the base rate of the disease. Direct use of $f_s(X | Z)$ from the data instead of (38) is the ‘base-rate fallacy’³⁷.

In supervised learning the classifier is trained to learn the most probable $f(X | Z)$ from the data. The training finds the $f(X | Z)$ that most closely fits the conditional frequency $f_s(X | Z)$ of the sampled data; this roughly corresponds to maximizing the first factor (37) described above. The architecture and the parameter regularizer of the classifier play the role of the second factor.

Bibliography

(‘de X ’ is listed under D, ‘van X ’ under V, and so on, regardless of national conventions.)
 Axelsson, S. (2000): *The base-rate fallacy and the difficulty of intrusion detection*. ACM Trans. Inf. Syst. Secur. 3³, 186–205. DOI:10.1145/357830.357849, <http://www.scs.carleton.ca/~soma/id-2007w/readings/axelsson-base-rate.pdf>.

³⁷ Russell & Norvig 2022 § 12.5; Axelsson 2000; Jenny et al. 2018.

- Baker, S. G., Pinsky, P. F. (2001): *A proposed design and analysis for comparing digital and analog mammography special receiver operating characteristic methods for cancer screening*. J. Am. Stat. Assoc. **96**⁴⁵⁴, 421–428. DOI:10.1198/016214501753168136.
- Bar-Hillel, M. (1980): *The base-rate fallacy in probability judgments*. Acta Psychol. **44**³, 211–233. DOI:10.1016/0001-6918(80)90046-3.
- Barber, D. (2020): *Bayesian Reasoning and Machine Learning*, online update. (Cambridge University Press, Cambridge). <http://www.cs.ucl.ac.uk/staff/d.barber/brml>. First publ. 2007.
- Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P., Heckerman, D., Smith, A. F. M., West, M., eds. (2011): *Bayesian Statistics 9*. (Oxford University Press, Oxford). DOI:10.1093/acprof:oso/9780199694587.001.0001.
- Bernardo, J.-M., Berger, J. O., Dawid, A. P., Smith, A. F. M., eds. (1996): *Bayesian Statistics 5*. (Oxford University Press, Oxford).
- Bernardo, J.-M., Smith, A. F. (2000): *Bayesian Theory*, repr. (Wiley, New York). DOI:10.1002/9780470316870. First publ. 1994.
- Bishop, C. M. (2006): *Pattern Recognition and Machine Learning*. (Springer, New York). <https://www.microsoft.com/en-us/research/people/cmbishop/prml-book>.
- Bridle, J. S. (1990): *Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition*. Neurocomputing **68**, 227–236. DOI:10.1007/978-3-642-76153-9_28.
- Camerer, C. F., Kunreuther, H. (1989): *Decision processes for low probability events: policy implications*. J. Policy Anal. Manag. **8**⁴, 565–592. DOI:10.2307/3325045.
- Cifarelli, D. M., Regazzini, E. (1979): *Considerazioni generali sull'impostazione bayesiana di problemi non parametrici. Le medie associative nel contesto del processo aleatorio di Dirichlet*. Riv. mat. sci. econ. soc. **2**^{1,2}, 39–52, 95–111.
- Damien, P., Dellaportas, P., Polson, N. G., Stephens, D. A., eds. (2013): *Bayesian Theory and Applications*. (Oxford University Press, Oxford). DOI:10.1093/acprof:oso/9780199695607.001.0001.
- Dawid, A. P. (2013): *Exchangeability and its ramifications*. In: Damien, Dellaportas, Polson, Stephens (2013): ch. 2:19–29. DOI:10.1093/acprof:oso/9780199695607.003.0002.
- de Finetti, B. (1929): *Funzione caratteristica di un fenomeno aleatorio*. In: *Atti del Congresso Internazionale dei Matematici*: ed. by S. Pincherle (Zanichelli, Bologna): 179–190. <https://www.mathunion.org/icm/proceedings>, <http://www.brunodefinetti.it/Opere.htm>. Transl. in Cifarelli, Regazzini (1979). See also de Finetti (1930).
- (1930): *Funzione caratteristica di un fenomeno aleatorio*. Atti Accad. Lincei: Sc. Fis. Mat. Nat. **IV**⁵, 86–133. <http://www.brunodefinetti.it/Opere.htm>. Summary in de Finetti (1929).
- (1937): *La prévision: ses lois logiques, ses sources subjectives*. Ann. Inst. Henri Poincaré **7**¹, 1–68. http://www.numdam.org/item/AIHP_1937__7_1_1_0. Transl. in Kyburg, Smokler (1980), pp. 53–118, by Henry E. Kyburg, Jr.
- Dunson, D. B., Bhattacharya, A. (2011): *Nonparametric Bayes regression and classification through mixtures of product kernels*. In: Bernardo, Bayarri, Berger, Dawid, Heckerman, Smith, West (2011): 145–158. <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.178.1521>, DOI:10.1093/acprof:oso/9780199694587.003.0005, older version at https://www.researchgate.net/publication/228447342_Nonparametric_Bayes_Regression_and_Classification_Through_Mixtures_of_Product_Kernels.
- Dunson, D. B., Pillai, N., Park, J.-H. (2007): *Bayesian density regression*. J. R. Stat. Soc. B **69**², 163–183.

- Ferguson, T. S. (1983): *Bayesian density estimation by mixtures of normal distributions*. In: Rizvi, Rustagi, Siegmund (1983): 287–302.
- Fienberg, S. E. (2007): *The Analysis of Cross-Classified Categorical Data*, 2nd ed. (Springer, New York). DOI:10.1007/978-0-387-72825-4. First publ. 1980.
- Fisher, R. A. (1963): *Statistical Methods for Research Workers*, rev. 13th ed. (Hafner, New York). First publ. 1925.
- (1967): *Statistical Methods and Scientific Inference*, repr. of 2nd rev. ed. (Oliver and Boyd, Edinburgh). First publ. 1956.
- Gal, Y., Ghahramani, Z. (2016): *Dropout as a Bayesian approximation: representing model uncertainty in deep learning*. Proc. Mach. Learn. Res. **48**, 1050–1059. See also Appendix at arXiv DOI:10.48550/arXiv.1506.02157.
- Good, I. J. (1950): *Probability and the Weighing of Evidence*. (Griffin, London).
- Good, I. J., Gaskins, R. A. (1971): *Nonparametric roughness penalties for probability densities*. Biometrika **58**², 255–277. DOI:10.1093/biomet/58.2.255.
- Goodman, L. A., Kruskal, W. H. (1954): *Measures of association for cross classifications*. J. Am. Stat. Assoc. **49**²⁶⁸, 732–764. DOI:10.1080/01621459.1954.10501231. See corrections Goodman, Kruskal (1957; 1958) and also Goodman, Kruskal (1959; 1963; 1972).
- (1957): *Corrigenda: Measures of association for cross classifications*. J. Am. Stat. Assoc. **52**²⁸⁰, 578. DOI:10.1080/01621459.1957.10501415. See Goodman, Kruskal (1954).
- (1958): *Corrigenda: Measures of association for cross classifications*. J. Am. Stat. Assoc. **53**²⁸⁴, 1031. DOI:10.1080/01621459.1958.10501492. See Goodman, Kruskal (1954).
- (1959): *Measures of association for cross classifications. II: Further discussion and references*. J. Am. Stat. Assoc. **54**²⁸⁵, 123–163. DOI:10.1080/01621459.1959.10501503. See also Goodman, Kruskal (1954; 1963; 1972).
- (1963): *Measures of association for cross classifications. III: Approximate sampling theory*. J. Am. Stat. Assoc. **58**³⁰², 310–364. DOI:10.1080/01621459.1963.10500850. See correction Goodman, Kruskal (1970) and also Goodman, Kruskal (1954; 1959; 1972).
- (1970): *Corrigenda: Measures of association for cross classifications. III: Approximate sampling theory*. J. Am. Stat. Assoc. **65**³³⁰, 1011. DOI:10.1080/01621459.1970.10481142. See Goodman, Kruskal (1963).
- (1972): *Measures of association for cross classifications, IV: Simplification of asymptotic variances*. J. Am. Stat. Assoc. **67**³³⁸, 415–421. DOI:10.1080/01621459.1972.10482401. See also Goodman, Kruskal (1954; 1959; 1963).
- Gregory, P. C. (2005): *Bayesian Logical Data Analysis for the Physical Sciences: A Comparative Approach with Mathematica Support*. (Cambridge University Press, Cambridge). DOI: 10.1017/CB09780511791277.
- Guyon, I., Gunn, S., Nikravesh, M., Zadeh, L. A., eds. (2006): *Feature Extraction: Foundations and Applications*. (Springer, Berlin). DOI:10.1007/978-3-540-35488-8.
- Hailperin, T. (2011): *Logic with a Probability Semantics: Including Solutions to Some Philosophical Problems*. (Lehigh University Press, Plymouth, UK).
- Hand, D., Christen, P. (2018): *A note on using the F-measure for evaluating record linkage algorithms*. Stat. Comput. **28**³, 539–547. DOI:10.1007/s11222-017-9746-6.
- Hjort, N. L. (1996): *Bayesian approaches to non- and semiparametric density estimation*. In: Bernardo, Berger, Dawid, Smith (1996): 223–253. With discussion by M. Lavine, M. Gasparini, and reply.
- Hunink, M. G. M., Weinstein, M. C., Wittenberg, E., Drummond, M. F., Pliskin, J. S., Wong, J. B., Glasziou, P. P. (2014): *Decision Making in Health and Medicine: Integrating*

- Evidence and Values*, 2nd ed. (Cambridge University Press, Cambridge). DOI:10.1017/CB09781139506779. First publ. 2001.
- Jaynes, E. T. (2003): *Probability Theory: The Logic of Science*. (Cambridge University Press, Cambridge). Ed. by G. Larry Bretthorst. First publ. 1994. DOI:10.1017/CB09780511790423, <https://archive.org/details/XQUHIUXHIQUHIQUXUIHX2>, <http://www-biba.inrialpes.fr/Jaynes/prob.html>.
- Jeffrey, R. C. (1965): *The Logic of Decision*. (McGraw-Hill, New York).
- Jeffreys, H. (1983): *Theory of Probability*, 3rd ed. with corrections. (Oxford University Press, London). First publ. 1939.
- Jeni, L. A., Cohn, J. F., De La Torre, F. (2013): *Facing imbalanced data: recommendations for the use of performance metrics*. Proc. Int. Conf. Affect. Comput. Intell. Interact. 2013, 245–251. DOI:10.1109/ACII.2013.47.
- Jenny, M. A., Keller, N., Gigerenzer, G. (2018): *Assessing minimal medical statistical literacy using the Quick Risk Test: a prospective observational study in Germany*. BMJ Open 8, e020847, e020847corr2. DOI:10.1136/bmjopen-2017-020847, DOI:10.1136/bmjopen-2017-020847corr2.
- Kim, H., Markus, H. R. (1999): *Deviance or uniqueness, harmony or conformity? A cultural analysis*. J. Pers. Soc. Psychol. 77⁴, 785–800. DOI:10.1037/0022-3514.77.4.785.
- Kruskal, W., Mosteller, F. (1979a): *Representative sampling, I: Non-scientific literature*. Int. Stat. Rev. 47¹, 13–24. See also Kruskal, Mosteller (1979b,c; 1980).
- (1979b): *Representative sampling, II: Scientific literature, excluding statistics*. Int. Stat. Rev. 47², 111–127. See also Kruskal, Mosteller (1979a,c; 1980).
- (1979c): *Representative sampling, III: The current statistical literature*. Int. Stat. Rev. 47³, 245–265. See also Kruskal, Mosteller (1979a,b; 1980).
- (1980): *Representative sampling, IV: The history of the concept in statistics, 1895–1939*. Int. Stat. Rev. 48², 169–195. See also Kruskal, Mosteller (1979a,b,c).
- Kyburg Jr., H. E., Smokler, H. E., eds. (1980): *Studies in Subjective Probability*, 2nd ed. (Robert E. Krieger, Huntington, USA). First publ. 1964.
- Lindley, D. V., Novick, M. R. (1981): *The role of exchangeability in inference*. Ann. Stat. 9¹, 45–58. DOI:10.1214/aos/1176345331.
- Lobo, J. M., Jiménez-Valverde, A., Real, R. (2008): *AUC: a misleading measure of the performance of predictive distribution models*. Glob. Ecol. Biogeogr. 17², 145–151. DOI: 10.1111/j.1466-8238.2007.00358.x, <https://www2.unil.ch/biomapper/Download/Lobo-GloEcoBioGeo-2007.pdf>.
- MacKay, D. J. C. (1992a): *The evidence framework applied to classification networks*. Neural Comput. 4⁵, 720–736. <http://www.inference.phy.cam.ac.uk/mackay/PhD.html>, DOI: 10.1162/neco.1992.4.5.720.
- (1992b): *Bayesian interpolation*. Neural Comput. 4³, 415–447. <http://www.inference.phy.cam.ac.uk/mackay/PhD.html>, DOI:10.1162/neco.1992.4.3.415.
- (1992c): *A practical Bayesian framework for backpropagation networks*. Neural Comput. 4³, 448–472. <http://www.inference.phy.cam.ac.uk/mackay/PhD.html>, DOI: 10.1162/neco.1992.4.3.448.
- Matthews, B. W. (1975): *Comparison of the predicted and observed secondary structure of T4 phage lysozyme*. Biochim. Biophys. Acta 405², 442–451. DOI:10.1016/0005-2795(75)90109-9.
- Mittone, L., Savadori, L. (2009): *The scarcity bias*. Appl. Psychol. 58³, 453–468. DOI: 10.1111/j.1464-0597.2009.00401.x.
- Mosteller, F., Fienberg, S. E., Rourke, R. E. K. (2013): *Beginning statistics with data analysis*, repr. (Dover, Mineola, USA). First publ. 1983.

- Müller, P., Quintana, F. A. (2004): *Nonparametric Bayesian data analysis*. Stat. Sci. **19**¹, 95–110. <http://www.mat.puc.cl/~quintana/publications/publications.html>.
- Murphy, K. P. (2012): *Machine Learning: A Probabilistic Perspective*. (MIT Press, Cambridge, USA). <https://problml.github.io/pml-book/book0.html>.
- Neal, R. M., Zhang, J. (2006): *High dimensional classification with Bayesian neural networks and Dirichlet diffusion trees*. In: Guyon, Gunn, Nikraves, Zadeh (2006): ch. 10:265–296. DOI: 10.1007/978-3-540-35488-8_11.
- North, D. W. (1968): *A tutorial introduction to decision theory*. IEEE Trans. Syst. Sci. Cybern. **4**³, 200–210. DOI:10.1109/TSSC.1968.300114, <https://stat.duke.edu/~scs/Courses/STAT102/DecisionTheoryTutorial.pdf>.
- Pearl, J. (1988): *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, rev. 2nd pr. (Kaufmann, San Francisco). DOI:10.1016/C2009-0-27609-4.
- Provost, F. (2000): *Machine learning from imbalanced data sets 101*. Tech. rep. WS-00-05-001. (AAAI, Menlo Park, USA). <https://aaai.org/Library/Workshops/2000/ws00-05-001.php>.
- Raiffa, H. (1970): *Decision Analysis: Introductory Lectures on Choices under Uncertainty*, 2nd pr. (Addison-Wesley, Reading, USA). First publ. 1968.
- Rasmussen, C. E. (1999): *The infinite Gaussian mixture model*. Adv. Neural Inf. Process. Syst. (NIPS) **12**, 554–560. <https://www.seas.harvard.edu/courses/cs281/papers/rasmussen-1999a.pdf>.
- Rizvi, M. H., Rustagi, J. S., Siegmund, D., eds. (1983): *Recent Advances in Statistics: Papers in Honor of Herman Chernoff on His Sixtieth Birthday*. (Academic Press, New York).
- Russell, S. J., Norvig, P. (2022): *Artificial Intelligence: A Modern Approach*, Fourth Global ed. (Pearson, Harlow, UK). First publ. 1995.
- Sammut, C., Webb, G. I., eds. (2017): *Encyclopedia of Machine Learning and Data Mining*, 2nd ed. (Springer, Boston). DOI:10.1007/978-1-4899-7687-1. First publ. 2011.
- Self, M., Cheeseman, P. C. (1987): *Bayesian prediction for artificial intelligence*. In: *Proceedings of the third conference on uncertainty in artificial intelligence (uai'87)*, ed. by J. Lemmer, T. Levitt, L. Kanal (AUAI Press, Arlington, USA): 61–69. Repr. in arXiv DOI:10.48550/arXiv.1304.2717.
- Sox, H. C., Higgins, M. C., Owens, D. K. (2013): *Medical Decision Making*, 2nd ed. (Wiley, New York). DOI:10.1002/9781118341544. First publ. 1988.
- Thorburn, D. (1986): *A Bayesian approach to density estimation*. Biometrika **73**¹, 65–75.
- Walker, S. G. (2013): *Bayesian nonparametrics*. In: Damien, Dellaportas, Polson, Stephens (2013): ch. 13:249–270.
- Yule, G. U. (1912): *On the methods of measuring association between two attributes*. J. R. Stat. Soc. **75**⁶, 579–652. DOI:10.1111/j.2397-2335.1912.tb00463.x.
- Zhu, Q. (2020): *On the performance of Matthews correlation coefficient (MCC) for imbalanced dataset*. Pattern Recognit. Lett. **136**, 71–80. DOI:10.1016/j.patrec.2020.03.030.