

# Neural networks, probability, decision theory

I. A. Davidovich

<ivan.davidovich@ntnu.no>

P.G.L. Porta Mana

<piero.mana@ntnu.no>

Draft of 20 January 2019 (first drafted 9 August 2018)

Notes on Bayesian and decision-theoretic interpretations of neural networks

*Note: Dear Reader & Peer, this manuscript is being peer-reviewed by you. Thank you.*

## 1 Inference and decision

First a summary about inference and decision, with some remarks of import for the rest of these notes.

Let's call, somewhat improperly, *inference*, *prediction*, or *forecast* the process of assigning numerical values to our degrees of belief  $p(p \mid t, D, M)$  about a set of propositions  $\{p\}$  given an (in some problems optional) input  $t$ , data  $D$  and a model  $M$ . The propositions often concern the happening of events or the values of some physical quantity  $x$ . The model is the set of numerical degrees of belief that are sufficient to arrive at  $p(p \mid t, D, M)$  using the plausibility calculus.

By *decision* or *choice* we call the process of choosing an action among a set  $\{a\}$ . We must specify a gain function or matrix  $G(a \mid p, M)$  that tells our gain in performing action  $a$  when the proposition  $p$  holds. The gain function is part of our model. According to decision theory (Raiffa et al. 2000; Jaynes 2003 ch. 13) we should choose one of the actions that maximize the expected gain:

$$\text{choose } \arg \max_a \sum_p G(a \mid p, M) p(p \mid t, D, M). \quad (1)$$

This decision problem can be rephrased in term of the maximum, rather than the expectation, with respect to  $p$  ( $\arg \max_a \max_p$ ), just by using a different gain function. Decision problems always involve inference problems, but not vice versa.

In the statistics literature a lot of effort is often spent on the numerical choice and determination of degrees of belief, but little is spent on the choice of the gain function. Yet it's clear from formula (1) that the solution of the decision problem depends equally on both.

Two models  $M'$ ,  $M''$  can lead to different inferences and yet be identical for decision purposes if

$$\sum_p G(a | p, M') p(p | t, D, M') = \sum_p G(a | p, M'') p(p | t, D, M'') \quad \text{for all } a, D. \quad (2)$$

For example, when  $p(p | t, D, M') = p(p | t, D, M'') \times h(p)$  and  $G(a | p, M') = G(a | p, M'')/h(p)$ . It's therefore impossible to ascertain the degrees of belief determined by a model if we're only given the decisions ensuing from it. This point will be important for our Bayesian interpretations of machine-learning procedures.

Note the different nature of actions and propositions. Often the decision problem is metonymically represented by the choice of an event or of the value of a quantity, but the intended actions are different: we say 'choose  $p$ ' but we mean 'behave as if  $p$  were true'. For example: the uncertain events  $\{p\}$  are {'rainy', 'sunny'} and we may represent the actions  $\{a\}$  by the same terms, but what we mean is actually {'take the umbrella', 'don't take the umbrella'}. Often the literature presents problems that look like inferences but are actually decisions. The 'choice of a parameter' is an example of this ambiguity.

Another ambiguity appears in the literature in the discussion of inference problems. A model parameter is often presented as the object of our inference, but our ultimate inference actually regards some more concrete event. For example: we model our degree of belief about observations of a physical or biologic quantity  $x$  with a normal density, and the mean and variance of the density are often presented as the objects of uncertainty; but our ultimate uncertainty is about the value of  $x$  in a new observation. The parameters are secondary.

Both ambiguities above (see Kass et al. 1995 § 2 for some examples) will also be important in our discussion of Bayesian interpretations of machine-learning procedures.

## 2 Synopsis of interpretations

I think it's important to first make a distinction between two tasks:

- (A) the general task that the neural network is meant to solve, once ready for use,
- (B) the task of training the neural network.

My impression is that the heuristics used in neural-net training address both tasks. In fact, task (B) somehow involves task (A). This double, nested nature of neural-network training is probably what makes it difficult to understand from a Bayesian and decision-theoretical point of view.

The next question is whether each task above is an inference problem or a decision problem. I think both tasks are decision problems. For (A), even in those cases where the network is designed to report some measure of confidence, its main goal is to give a definite output out of a set of possible ones. Think of examples like optical-character recognition. For (B), the software engineer must in the end deliver a network, appropriately prepared, ready to be used in an automated way. One complication in this question is whether the ready-for-use neural network is meant to continue learning from data given to it or not. It seems to me this isn't the case (think again of optical-character recognition).

Let's now examine the two decision tasks above in more detail.

In task (A) the neural net must choose a specific output  $y$ , given an input  $t$ . Such choice is an action; thus  $y$  is a label for an action, more than for a quantity. In optical-character recognition, for example,  $t$  is an array of pixel intensities of an image of text; the uncertain proposition  $p$  is about the typographic character  $x$  that was intended by the original writer of the text; and  $y$  is the choice of a specific character, say as a unicode number. There may of course be an isomorphism between  $\{y\}$  and  $\{x\}$ . This is an example of the distinction between action and quantity mentioned in the previous section. In this decision problem no additional data  $D$  are used: the conditional plausibilities  $p(x | t, M)$  are given and fixed. The model  $M$  may of course have been constructed or chosen using data, but these do not enter the present task. For the neural-network user, this decision task reduces to a specific function  $f: t \mapsto y$ , or possibly  $f: t \mapsto (y, c)$  where  $c$  is a measure of confidence. This function can be obtained by many different pairs of plausibility distributions and gain functions; see the remarks in the previous section.

Task (B) can be seen from two points of view:

(1) If we see task (A) as summarized by a function  $f: t \mapsto y$ , then the task of the software engineer is to choose such a function from a set  $\{f\}$  of possible ones. To make this choice she needs a gain function  $G(f | p, E)$  and the set of plausibility distributions  $p(p | t, D, E)$ , chosen

according to some model  $E$ . But what are the propositions  $\{p\}$  about? that is, what is the engineer uncertain about? And what is the gain about?

The engineer is uncertain about which pairs  $(t, x)$  the neural network will encounter during its actual use (see Jaynes 2003 § 22.3 for an enlightening discussion about this). We may distinguish two cases of such uncertainty: either there is a deterministic relation between the two quantities, that is, the engineer knows that only one  $x$  corresponds to each  $t$ , although she doesn't know which; or there isn't a deterministic relation, and the engineer is uncertain about which values of  $x$  will appear for multiple occurrences of the same  $t$ .

As for the gain in choosing  $f$ , my intuition is that it is the average gain that the end user of the neural net will obtain during its use. If  $G(y | x)$  is the gain function for the end user, and  $p(x, t | E)$  is the engineer's degree of belief about the  $(t, x)$  pairs the end user will encounter, the expected gain of choosing a function  $f: t \mapsto y$  is

$$\sum_{x,t} G[f(t) | x] p(x, t | E), \quad (3)$$

so the engineer should choose an  $f$  which maximizes this expectation. If the engineer knows that there's a deterministic dependence between  $t$  and  $x$ , say  $r: t \mapsto x$ , but she's uncertain as to this function, then the expectation above may be rewritten

$$\sum_{r,t} G[f(t) | r(t)] p(r | E) p(t | E), \quad (4)$$

assuming that her knowledge of  $t$  is irrelevant for her uncertainty about  $r$ .

Data  $D$  are useful to the engineer for arriving at better-informed degree of beliefs about  $p(x, t | D, E)$ . So, considering the expectation (3), the software engineer must

$$\text{choose } \arg \max_f \sum_{x,t} G[f(t) | x] p(x, t | D, E). \quad (5)$$

The heuristic procedures typical neural-net learning could be doing the following:

- calculating the posterior  $p(x, t | D, E)$ ,
- searching for the  $f$  that maximizes the expectation above,

- both.

An important feature of the point of view currently under discussion is that there the space of functions  $\{f\}$  can be distinct from the space of plausibility models  $p(x, t \mid D, E)$  or the space of functions  $\{r\}$ . This is in contrast with the Bayesian interpretation most often discussed in the literature (Tishby et al. 1989; Levin et al. 1990; MacKay 1992a,b; Neal 1996; Barber et al. 1998; Bishop 2006 esp. §§ 1.2–3, 3.3–4, 5.7), which assumes a neural network to be a particular plausibility model.

✠ To be continued

---

old text

---

The purpose of a neural network and some other machine-learning methods is to choose an output  $y$  given an input  $x$  and data relating the possible inputs and outputs. The way this choice is made and the way it is influenced by already-observed data make intuitive sense. As stated, the task above involves uncertainty and decision; it can therefore also be approached methodically using the plausibility calculus – the ‘gold standard’ as Srivastava, Hinton, et al. (2014) called it – and the principles of decision theory. It is interesting to see how the intuitive neural-net approach and the methodical one relate to each other.

This relation was first studied by Tishby, Levin, Solla (Tishby et al. 1989; Levin et al. 1990) and then brilliantly explored by MacKay (1992a,b) and others (1996; 1998).

### 3 \*\*\*

The study and grasp of the connection between neural nets and probability is beneficial to both. For neural nets, this connection *might* not lead to any practical improvements – because of scalability etc.\*\*\* – but it’s very useful for understanding open problems, like the quantification of confidence of neural-network outputs [refs\*\*\*], overtraining, or the choice of architectural constants. For example, the successful idea of using ‘dropout’ was also motivated (Hinton et al. 2012; Srivastava et al. 2014) by the ‘mixing’ of predictions typical of probability. For probability, this connection suggests new, efficient parametric families and approximation methods.

In this note I want to explain the connection between neural nets and some probability models. What’s new in this note, compared with the 1990s studies cited above?

- I want to try to show this connection in a *geometric, visual* way.
- I extend those studies in two directions:
  1. including the decision-theoretic character of the problem. This leads to a different probabilistic interpretation of the error function of neural nets;
  2. basing the connection on so-called ‘partially exchangeable models’, to be discussed later. They give us a better understanding of overtraining, and are also connected with generative models [refs\*\*\*].

- *Repetita iuvant*: knowledge of those studies seems to be forgotten or much diluted today.

This is what I do not purport to do:

- proving that the connection between neural nets and probability discussed here is unique. Even if it isn't unique, it's still insightful and useful;
- suggest that neural nets should be used 'more in accordance' with the probability calculus. The connection already shows that they are approximate probability calculations; their speed and power come from their approximate nature;
- [add here whatever else you'd like my intentions not to be.]

## 4 Introduction to partially exchangeable models

Consider a classification problem: given input  $x \in \{1, \dots, M\}$ , we want to predict output  $y \in \{1, \dots, N\}$ . Let's consider both discrete.

Two cases must be distinguished: whether every input has several possible outputs, or just one. The second case is a special instance of the first, and could be simply treated as the prediction of a 'function'  $x \mapsto y$ . But both cases are instances of inference with a *partially exchangeable* model.

Our assumptions are summarized in the proposition *I*. We assume that the probability of  $y$ , for each kind of input  $x$ , is exchangeable. Thus the probability for several pairs

$$p(y^{(T)}, \dots, y^{(1)} \mid x^{(T)}, \dots, x^{(1)}, I) \quad (6)$$

must have a partially exchangeable form [refs]: for every kind of input  $x$  we consider the  $N$  relative frequencies  $f_x := (f_{1|x}, \dots, f_{N|x})$  of the  $N$  possible outputs. The probability above is given by this integral:

$$p(y^{(T)}, \dots, y^{(1)} \mid x^{(T)}, \dots, x^{(1)}, I) = \int \dots \int \left[ \prod_{x=1}^M \prod_t^{x^{(t)}=x} f_{y^{(t)}|x} \right] p(f_1, \dots, f_M \mid I) df_1 \dots df_M. \quad (7)$$

The complicated look of this formula doesn't do justice to its simple and intuitive interpretation:

Choose a kind of input  $x$ . We're judging that the probability for a very large sequence of outputs generated with this input doesn't depend on their order. Suppose we knew the relative frequencies of each kind of output in this sequence,  $(f_{1|x}, \dots, f_{N|x}) =: f_x$ , and knew nothing else. Then out of symmetry reasons we should give a probability  $f_{y|x}$  of observing outcome  $y$  at the next observation. This is just a 'drawing without replacement' problem. At the subsequent observations with the same input we give again the same probabilities to each  $y$ , because the sequence is so large that we approximate 'drawing without replacement' with 'drawing with replacement'. The probability for a sequence of outputs  $y^{(1)}, \dots, y^{(T)}$  from the same input  $x$  is then

$$p(y^{(T)}, \dots, y^{(1)} \mid \text{same input } x, f_x, I) = f_{y^{(T)}|x} \times \dots \times f_{y^{(1)}|x}. \quad (8)$$

Within a larger sequence of outputs from all possible inputs, the outputs  $y^{(0)t}$  coming from input  $x$  are those for which  $x^{(0)t} = x$ . Thus we can write the formula above more generally as

$$\prod_t^{x^{(t)}=x} f_{y^{(t)}|x} \quad (9)$$

The above reasoning applies for each kind of input  $x$ . Thus the probability for a sequence of outputs coming from different inputs is the product of the probabilities above for all different  $x$ :

$$p(\text{all outputs} \mid \text{inputs}, f_1, \dots, f_M, I) = \prod_{x=1}^M \prod_t^{x^{(t)}=x} f_{y^{(t)}|x}. \quad (10)$$

Now what if we don't know the relative frequencies  $f_x$ , for any input? Then we assign a probability distribution over their possible values and use the law of total probability:

$$p(\text{all outputs} \mid \text{inputs}, I) = \int \dots \int p(\text{all outputs} \mid \text{inputs}, f_1, \dots, f_M, I) p(f_1, \dots, f_M \mid I) df_1, \dots, df_M. \quad (11)$$



Substituting the explicit expression (10) in this formula we obtain formula (11). In summary,

$$\int \cdots \int \left[ \prod_{x=1}^M \prod_{t=1}^{x^{(t)=x}} f_{y^{(t)}|x} \right] p(f_1, \dots, f_M | I) df_1 \cdots df_M. \quad (12)$$

product over all inputs
probability for outputs from same input
uncertainty over all frequencies

The possible frequencies give one input,  $f_x \equiv (f_{1|x}, \dots, f_{N|x})$ , belong to the  $(N - 1)$ -dimensional simplex

$$\Delta := \{(f_1, \dots, f_N) \mid f_i \geq 0, \sum_i f_i = 1\}, \quad (13)$$

and the collection of possible frequencies  $(f_1, \dots, f_M)$  belongs to the  $M$ -fold Cartesian product  $\Delta^M$ . From now on we denote  $\mathbf{f} := (f_1, \dots, f_M)$ .

The probability for a new sequence of  $T'$  outputs given their inputs and given that we've learned a previous sequence of  $T$  input-output pairs is determined by Bayes's theorem:

$$p(y^{(T')}, \dots, y^{(T+1)} \mid x^{(T')}, \dots, x^{(T+1)}, y^{(T)}, x^{(T)}, \dots, y^{(1)}, x^{(1)}, I) = \frac{\int \left[ \prod_{x=1}^M \prod_{t=T+1, \dots, T'}^{x^{(t)=x}} f_{y^{(t)}|x} \right] p(\mathbf{f} \mid I) d\mathbf{f}}{\int \left[ \prod_{x=1}^M \prod_{t=1, \dots, T}^{x^{(t)=x}} f_{y^{(t)}|x} \right] p(\mathbf{f} \mid I) d\mathbf{f}}. \quad (14)$$

This formula is equivalent to (7) with an updated distribution for the frequencies:

$$p(\mathbf{f} \mid y^{(T)}, x^{(T)}, \dots, y^{(1)}, x^{(1)}, I) d\mathbf{f} = \frac{\left[ \prod_{x=1}^M \prod_{t=1, \dots, T}^{x^{(t)=x}} f_{y^{(t)}|x} \right] p(\mathbf{f} \mid I)}{\int \left[ \prod_{x=1}^M \prod_{t=1, \dots, T}^{x^{(t)=x}} f_{y^{(t)}|x} \right] p(\mathbf{f} \mid I) d\mathbf{f}} d\mathbf{f}. \quad (15)$$

The form of this updated distribution has important consequences for our learning process.

If the number of learned data is enough large compared with the numbers  $M, N$  of possible inputs and outputs and with the magnitude

of the initial distribution for the frequencies, and if the latter is strictly positive, then the updated distribution becomes very peaked on the collection of relative frequencies  $(q_1, \dots, q_M)$  of the learned outputs for all input values. This can be seen from the asymptotic expression in terms of the relative entropy (Kullback-Leibler divergence)  $D$ ,

$$p(f | y^{(T)}, x^{(T)}, \dots, y^{(1)}, x^{(1)}, I) \propto \exp[-\sum_x T_x D(q_x | f_x)] p(f | I), \quad (16)$$

where  $T_x$  is the number of observations with input  $x$ , with  $\sum_x T_x = T$ .

If the number of learned data is small compared with the dimensions of the input and output spaces, then the initial distribution for the frequencies  $p(f | I) df$  greatly influence our inference (14). This distribution determines two important \*\*\*

## 5 Utility functions and probabilities

Utility of behaving as if proposition  $B \in X$  is true given that proposition  $A \in X$  is true:  $c(B | A)$ . Probability for  $A$  given  $D, I$ :  $P(A | D, I)$ . Optimal decision is  $B$  that maximizes

$$\sum_A c(B | A) P(A | D, I). \quad (17)$$

Now consider different probabilities for all  $A$  given the same data  $D$ :  $P(A | D, I')$ . The decision  $B$  will still be the same if we use a new utility

$$c'(B | A) := c(B | A) \frac{P(A | D, I)}{P(A | D, I')}. \quad (18)$$

So the same choice can be made with a different probability, if the utility is appropriately changed, provided  $p(A | D, I') > 0$  for all  $A$ .

This leads to a slightly more general view than Tishby et al.'s (1989; 1990) and Mackay's (1992a,b)

✂ Point out that the joint parameter density allows us to make inferences 1. from data about one output to data about the same output; 2. from data about one output to data about a different output – *this is generalization*

✂ Point out that generalization (from one output to another) is *fully* determined by the form of the joint parameter prior

## Bibliography

(‘de X’ is listed under D, ‘van X’ under V, and so on, regardless of national conventions.)

- Barber, D., Bishop, C. M. (1998): *Ensemble learning in Bayesian neural networks*. In: *Neural networks and machine learning*. Ed. by C. M. Bishop (Springer), 215–237. <https://www.microsoft.com/en-us/research/publication/ensemble-learning-in-bayesian-neural-networks/>.
- Bishop, C. M. (2006): *Pattern Recognition and Machine Learning*. (Springer, New York).
- Hinton, G., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R. (2012): *Improving neural networks by preventing co-adaptation of feature detectors*. [arXiv:1207.0580](https://arxiv.org/abs/1207.0580).
- Jaynes, E. T. (2003): *Probability Theory: The Logic of Science*. (Cambridge University Press, Cambridge). Ed. by G. Larry Bretthorst. First publ. 1994. <https://archive.org/details/XQUHIUXHIQUHIQUXUIHX2>, <http://www-biba.inrialpes.fr/Jaynes/prob.html>.
- Kass, R. E., Raftery, A. E. (1995): *Bayes factors*. J. Am. Stat. Assoc. **90**<sup>430</sup>, 773–795. <https://www.stat.washington.edu/raftery/Research/PDF/kass1995.pdf>; <https://www.andrew.cmu.edu/user/kk3n/simplicity/KassRaftery1995.pdf>.
- Levin, E., Tishby, N., Solla, S. A. (1990): *A statistical approach to learning and generalization in layered neural networks*. Proc. IEEE **78**<sup>10</sup>, 1568–1574.
- MacKay, D. J. C. (1992a): *Bayesian interpolation*. Neural Comp. **4**<sup>3</sup>, 415–447. <http://www.inference.phy.cam.ac.uk/mackay/PhD.html>.
- (1992b): *A practical Bayesian framework for backpropagation networks*. Neural Comp. **4**<sup>3</sup>, 448–472. <http://www.inference.phy.cam.ac.uk/mackay/PhD.html>.
- Neal, R. M. (1996): *Bayesian Learning for Neural Networks*. (Springer, New York). <https://www.cs.toronto.edu/~radford/bnn.book.html>.
- Raiffa, H., Schlaifer, R. (2000): *Applied Statistical Decision Theory*, reprint. (Wiley, New York). First publ. 1961.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R. (2014): *Dropout: a simple way to prevent neural networks from overfitting*. J. Mach. Learn. Res. **15**, 1929–1958.
- Tishby, N., Levin, E., Solla, S. A. (1989): *Consistent inference of probabilities in layered networks: predictions and generalizations*. Int. Joint Conf. Neural Networks **1989**, II-403–II-409.