

Naïve Inference viewed as Computation

Chris Thornton
Informatics
University of Sussex
Brighton
BN1 9QH
UK
c.thornton@sussex.ac.uk

July 11, 2011

Abstract

Use of Bayesian models to explain both high- and low-level aspects of cognitive function promises better connections between cognitive science and cognitive neuroscience. But standing in the way are fundamental problems, such as the computational intractability of Bayesian inference, and the general difficulty of understanding how Bayesian calculation can deal with structural representation. Getting around the problem of intractability seems to involve devising effective methods for approximating optimal inference. But there is the alternative of simplifying the interpretation of how inference arises. While the process is normally taken to involve calculations over an implied joint distribution, it is possible to view it more simply as data-driven application of conditional assertions. This naive interpretation has several advantages with regard to tractability and representation. The paper formalizes the model and demonstrates some of its virtues.

1 Introduction

Methods of probabilistic inference are increasingly under the spotlight as cognitive science moves towards greater use of Bayesian approaches (e.g. Knill and Richards, 1996; Chater et al., 2006). Strongly facilitating this trend are the ‘graphical models’ for performing inference with complex probabilistic information (e.g. Pearl, 1988). Somewhat obstructing it is the knowledge that Bayesian inference is computationally intractable (Cooper, 1990). For some, this intractability does not vitiate the explanatory value of Bayesian inference viewed as an optimal solution for a cognitive or perceptual problem (e.g. Anderson, 1990). The point is made that such models can be viewed as theories at a functional level of abstraction, e.g., the *computational level* of Marr’s scheme

(Marr, 1982). For others, the intractability issue is more concerning (e.g. Danks, 2008).

A more basic problem with Bayesian approaches relates to the constraints that methods of probabilistic inference place on source information. All standard methods assume a probabilistic model that implicitly represents an underlying joint distribution (Russell and Norvig, 2010). Inference is progressed through marginalization, i.e., constrained summations of values appearing within this joint distribution. The difficulty is that the process is then infeasible with regard to probabilistic information that does not properly represent a joint distribution.

The problem particularly affects information embodying conditional cycles. Imagine we have boolean variables X , Y and Z representing basic features of the environment (e.g., *rain*, *mud* and *humidity*). These are known to be conditionally related as follows: $P(Y|X) = 0.8$, $P(Z|Y) = 0.4$, $P(X|Z) = 0.6$. Given $P(X) = 1.0$ (e.g., observed evidence of *rain*) what probability should we infer for Y ? Standard methods of probabilistic inference cannot produce any answer. The existence of a cycle among the conditional relationships means the model cannot be viewed as representing a joint distribution. Inference through marginalization is ruled out.

One way around the problem is to take an approach which dispenses with the underlying joint distribution altogether. Instead of treating this as the key reference, we treat the probabilities the model asserts as fundamental. Conditional values are not viewed as constraints on an implied joint distribution. Rather they are viewed as mandating acts of inference. On this basis, $P(X|Y)$ legitimates inference of an unconditional value of X whenever an unconditional value of Y is identified. Simultaneous inferences for the same variable can be accommodated by letting them contribute equally to the inferred value, subject to the constraint that probabilities sum to 1. Inferred values are then obtained by normalized summation of products, much as in marginalization. The regime accommodates cyclic conditionality. Application to the case above, for example, produces a final inferred probability of 0.8 for variable Y .

This regime, in which asserted probabilities are treated as inference mandates, can be viewed as a naive form of probabilistic inference (cf. Hansson et al., 2008) in which inference arises directly from the semantics of conditional assertions. Application of conditional probabilities to relevant unconditional values has the potential to identify new unconditional values. These can then be the basis for production of further values, and so on, in a potentially infinite sequence. Naïve inference becomes the behaviour of a data-driven machine—a *naïve inference machine* as it will be called.

Is such a primitive procedure likely to have any useful application? There are various ways it might do so. Where there is a need to deal with probabilistic information embodying cycles, standard forms of inference cannot be used. Dealing with the situation without requiring additional assumptions, naive inference may then have a use as a least-commitment approach to inference in the presence of conditional cycles.

The regime also makes connections with non-inferential models of mecha-

nism. Accommodation of cycles means that naive inference can exhibit looping. Naïve inference machines have the potential for infinite processing. In the case where extremal probabilities (1s and 0s) are deployed, the behaviour is that of a digital device with iterative behaviour; a connection can then be made between naive inference and *computation*. In fact, as Section 4 demonstrates, naive inference machines are Turing equivalent: they can model any form of computational behaviour. This suggests the potential for explanations unifying modeling of inferential behaviour with modeling of computational behaviour.

An appealing application of naive inference is in connection with conventional Bayesian accounts. A difficulty with these is the computational intractability of the process model. Some accounts cash this out by establishing the means by which processing is implemented (e.g. Pouget et al., 2003; Körding and Wolpert, 2006). Others reserve the right *not* to do so on grounds of explanatory abstraction (e.g. Chater and Oaksford, 2008). In general, there is a need to establish a better connection between the theoretical ideal of Bayesian inference, and practical mechanisms by which it can be pursued.

Naïve inference cannot be viewed as approximating Bayesian inference. But since it treats the semantics of the conditional probability assertion in the same way, situations can arise in which both methods produce the same result. There is then the potential to treat naive inference as modeling the ‘bounded rationality’ (Simon, 1957) that ideal rationality must express in practice. Instead of approximating the Bayesian ideal, this introduces a less sophisticated interpretation of what inference involves.

Taking bounded rationality to be modeled by naive inference does achieve some of the benefits of approximation, however. The process model is no longer computationally intractable. There is also the prospect of better connections with models of neural processing. Mediated by simple operations of summation and normalization, the machinery of naive inference is likely to be more easily related to known functionalities of the brain (cf. Dayan and Abbott, 2001)

The aim of the paper is to set out the proposed model of naive probabilistic inference in more detail, to formalize it mathematically, and to examine its potential uses. There are five main sections. The next section formalizes the inference model. The third section examines the degree to which naive inference can emulate ideal Bayesian estimation. Following this, there is an examination of the sense in which naive inference machines are Turing equivalent. The final section is a summary.

2 Formalization

In this approach, probabilistic information is assumed to take the form of probability values for random discrete variables. Conventional notation is used. Thus $P(X = X_i)$ denotes the probability that random variable X has value X_i . If X is boolean, $P(X)$ is a shorthand for $P(X = \text{true})$.

Given variables X and Y are both boolean, the conditional expression $P(X|Y)$ expresses the probability that X is true given Y is true. In naive inference, this

is understood to directly mandate acts of inference. Identification of an unconditional probability for Y , either by assumption or (prior) inference, establishes the possibility of inferring an unconditional value of X . Where inference using multiple conditionals is legitimated, derived products are assumed to contribute equally to the inferred value, subject to the constraint that probabilities in a distribution must sum to 1.

A *probabilistic model* is defined to be a set of conditional and unconditional probability values for random discrete variables. Letting M label such a model, $P_M(X)$ is the unconditional probability of X in model M , and $P_M(X|Y)$ is the conditional probability of X given Y in model M . Bold font is used to denote distributions. Thus $\mathbf{P}_M(X)$ is the distribution on variable X represented by model M . Given X is boolean, $P(X) = 1$ becomes a shorthand for $\mathbf{P}(X) = \langle 0, 1 \rangle$.

Defining the unnormalized inferred probability for X_i in model M to be

$$P_M(X_i) = \sum_{c \in C(X_i, M)} P_M(X_i|c)P_M(c) \quad (1)$$

the distribution inferred for X in model M is

$$\mathbf{P}'_M(X) = \alpha \langle P_M(X_1), \dots, P_M(X_n) \rangle \quad (2)$$

In Eq. 1, $C(x, M)$ is the set of conditions that figure in conditional probabilities asserted for values of X in model M , n is the number of values of X , and α is the normalization function. The inference of Eq. 2 is taken to be defined just in case the model provides evaluations for all conditions. That is to say, it is defined if the model provides unconditional values for all applicable conditions.

Building on this, we can define the complete set of inferences that can be obtained through application of the inference step to an existing model. Termed a *revision*, this is denoted by adding a prime to the model label. Thus M' denotes the naive-inferential revision of model M :

$$M' = \{ \mathbf{P}'_M(X) \mid X \in M \wedge P'_M(X) \neq P_M(X) \} \quad (3)$$

Here, $X \in M$ is true if and only if variable X features in model M .

Recursive evaluation of M' can then be the means of generating a sequence of inferential revisions of a particular model. M_i , the i 'th model in the sequence, must satisfy

$$M_i = M'_{i-1} \implies M_{i-1}$$

where the ' \implies ' operator denotes imposition of M'_{i-1} on M_{i-1} . Specifically, $M' \implies M$ represents addition of all unconditional values in M' to M , with preference given to values of M' where both sets give values for the same variable. Letting M_0 label the set of conditional values in model M , and M'_0 be the corresponding set of unconditional values, the sequence of revisions for model M then takes a well-defined form. This is denoted $N(M)$:

$$N(M) = (M'_0, \dots, M'_n)$$

$N(M)$ can also be viewed as labeling the naive inference machine defined by model M . The behaviour of the machine is production of revisions. The output is the revision sequence itself.

A simple illustration is provided by the model of Table 1.

$$\begin{array}{l|l} P(X = 1) = 1 & P(X = 0) = 0 \\ P(X = 1|X = 1) = 0 & P(X = 0|X = 1) = 1 \\ P(X = 1|X = 0) = 1 & P(X = 0|X = 0) = 0 \end{array}$$

Table 1: Naïve-inferential oscillator

In this probabilistic model, X is conditionally dependent on itself, but with the conditioned value always being the opposite of the conditioning value. Naïve inference then yields an *infinite* revision sequence within which the value of X continually changes between its two values. Given shorthand representations for binary-valued distributions, and vertical arrangement of the elements of the sequence, $N(M)$ evaluates as

$$N(M) = \left(\begin{array}{l} \{P(X) = 1\}, \\ \{P(X) = 0\}, \\ \{P(X) = 1\}, \\ \{P(X) = 0\}, \\ \dots \end{array} \right)$$

This is the behaviour of the naive inference machine defined by the model of Table 1. The (infinite) sequence of revisions generated is the output the machine produces.

3 Emulation of ideal Bayesian inference

In optimal Bayesian inference, a probabilistic model comprising conditional and unconditional probabilities (i.e. *priors*) is assumed to represent a joint distribution. The process of inference involves determining unobserved values in this distribution. Where unconditional values are considered to constitute *evidence*, derived values are *posteriors*. While the process is computationally intractable (Cooper, 1990), graphical models such as (Pearl, 1988) are often effective. These allow the process to be progressed in a way that maximally exploits independence relationships for factorising calculations.

The naive model of inference relinquishes the assumption of an underlying joint distribution. Inference is taken to involve data-driven application of conditional assertions. However, the two approaches place the same interpretation on conditional assertions. In both interpretations, it is axiomatic that

$$P(X) = P(X|Y)P(Y)$$

given known values for $P(Y)$ and $P(X|Y)$. Inference mediated solely by this rule is thus progressed identically under naive and ideal protocols.

This can be illustrated using the ‘sprinkler’ example, a popular scenario for illustrating the behaviour of Bayesian inference using Bayesian networks (e.g. Pearl, 1988, p. 56). In this example, variable *Rain* represents the occurrence of rain, variable *Sprinkler* represents a sprinkler being on overnight, and variable *GrassWet* represents the grass being wet. These are all boolean variables taking values T and F, representing *true* and *false* respectively. Conditional and unconditional probabilities for these variables are illustrated schematically in Figure 1.

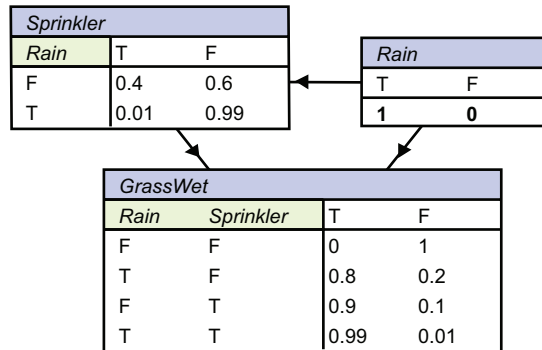


Figure 1: Probabilistic model for ‘sprinkler’ example.

In this diagram, each variable is represented by a table. Values of the variable correspond to columns, while rows represent conditions. Where unconditional values are given, they appear in the bottom row of a table. Thus, the unconditional probability of *Rain* is here shown to be 1. The conditional probability of *GrassWet* given *Rain* = T and *Sprinkler* = F is 0.8, and so on. The values shown can be viewed as comprising a probabilistic model in the present sense of the term. They can also be regarded as comprising a probabilistic model in the conventional sense of the term. Indeed, given the simplicity of the model guarantees conditional independence of *GrassWet* given *Sprinkler* and *Rain*, it also represents a Bayesian network. On this view, the tables are the conditional probability tables (CPTs) of a standard Bayesian network.

Say we discover that $P(Rain) = 1.0$, and wish to infer the effect on $P(Sprinkler)$. We must decide whether we wish to treat $P(Rain) = 1.0$ as the unconditional probability of *Rain*, or as observed evidence. This makes a difference in the case of the Bayesian network: in one case the network will calculate new values for *Sprinkler* and *GrassWet* through probabilistic inference. In the other, priors for these two variables will become implicitly defined. Derived probabilities are the same however. The emerging prior (or inferred probability) for *Sprinkler*

is $P(\textit{Sprinkler}) = 0.01$ and the emerging prior (or inferred probability) for $\textit{GrassWet}$ is $P(\textit{GrassWet}) = 0.802$.

In this case, optimal Bayesian inference relies purely on derivation (and normalization) of products. There is no application of Bayes' rule. Naïve inference is then able to emulate the process to the same effect. Applying Eq. 2 to the model of Figure 1, the initial revision is determined to contain $P(\textit{Sprinkler}) = 0.01$. It will not contain a value for $\textit{GrassWet}$ however, since all the conditions for that variable require unconditional values for both \textit{Rain} and $\textit{Sprinkler}$. Establishment of an unconditional probability for $\textit{Sprinkler}$ then prompts a second revision comprising $P(\textit{GrassWet}) = 0.802$.

The Bayesian network generates the same posterior value (or emergent prior) for $\textit{GrassWet}$ as does naive inference. Indeed, given the assumption that $\textit{Sprinkler}$ is summed-out in the Bayesian network before $\textit{GrassWet}$, the two regimes produce the derivations in the same order. In simple cases like this, ideal Bayesian inference and naive inference can produce the same result. The Bayesian network has a more extensive behavioural repertoire, of course. Utilizing Bayes' rule for inverting conditional probabilities, it could be the means of calculating an unconditional value for \textit{Rain} given evidence involving $\textit{GrassWet}$, for example.

The key difference between naive and Bayesian inference is that the former makes no direct use of Bayes' rule for inverting conditional probabilities. However, this does not necessarily mean that naive inference is unable to reproduce classical Bayesian hypothesis selection. In this scenario, inference is used to determine the hypothesis that optimally explains certain data, given priors on the hypotheses and the data, and conditional values for data given hypotheses (i.e., relevant *likelihoods*). The functionality applied, however, involves deriving products of conditional and unconditional values in the usual way. Naïve inference can thus reproduce the effect provided variables are provided whose unconditional values are those that *would* be obtained through application of Bayes' rule. On the assumption that such proxies are introduced (or assumed to exist), naive inference has the potential to reproduce hypothesis-selection involving application of Bayes' rule. On this basis, naive inference can reproduce the classic inferential scenario of Bayesian estimation.

3.1 Introducing a conditional cycle

Naïve inference has the advantage of being able to accommodate models representing conditional cycles. This effect can be illustrated using a modification of the 'sprinkler' example. In the original example, there are no cycles among the conditions. The conditional structure takes the form of a directed acyclic graph (DAG), as required for a Bayesian network. Consider now the variation of Figure 2. Here variables \textit{Rain} and $\textit{GrassWet}$ have the same conditional relationship. But we now have a $\textit{Humidity}$ variable, which is conditionally dependent on $\textit{GrassWet}$. This produces a conditional cycle in which $\textit{GrassWet}$ is made more probable by \textit{Rain} , $\textit{Humidity}$ is made more probable by $\textit{GrassWet}$, and \textit{Rain} is made more probable by $\textit{Humidity}$.

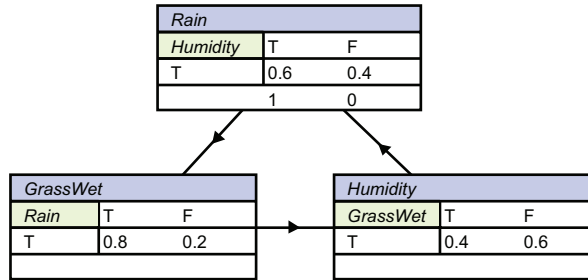


Figure 2: Model for Rain/Grass/Humidity example.

This cycle violates the conditional-independence requirements of the Bayesian network, and thus the assumption of an underlying joint distribution. Standard methods cannot be applied but the naive procedure is unaffected. Revisions are identified in the usual way. The presence of the cycle creates the potential for an infinite sequence. But in this case, inference rapidly converges on a particular set of unconditional values. Taking M to be model of Figure 2, we obtain the following finite sequence:

$$N(M) = \left(\begin{array}{l} \{P(Rain) = 1\}, \\ \{P(GrassWet) = 0.8\}, \\ \{P(Humidity) = 0.4\}, \\ \{P(Rain) = 0.6\}, \\ \{P(GrassWet) = 0.8\}, \end{array} \right)$$

The original unconditional probability for *Rain* appears here as the zeroth element of the sequence. Derivation of distributions by Eq. 2 then yields four revisions, the last of which makes no changes to the model. Inference terminates at this point, with a final probability of 0.8 for *GrassWet*. This inferred value may be viewed as reflecting the cyclical dependency between the three variables. Alternatively, the inferential process may be viewed as a dynamic projection of the asserted conditional relationships.

4 Emulation of Turing Machines

Attention now given turns to other interpretations that can be applied to naive inference. This section examines the sense in which naive inference is Turing equivalent. By showing that naive inference can model any Turing-machine computation, the procedure is shown to have computational power equivalent to that of a digital computer, or any other device for effective computation.

A Turing machine is defined in terms of a state-transition table, a ‘tape’ containing a sequence of symbols, and an initial tape position. In each cycle, the machine reads the symbol from the current position on the tape and responds by writing a symbol at that position, moving the tape one position left or

right, and entering a new state. The behaviour of the machine is the result of repeatedly applying such transitions, until the halt state is reached. The final output obtained is the revised contents of the tape.

To translate a Turing machine into an equivalent naive inference machine we can proceed as follows. For each cell of the Turing machine's tape, we introduce a named variable whose values are the symbols used by the Turing machine. We introduce variables to represent the current symbol read, and the current symbol to be written. We also introduce variables to represent the current state, the current move, and the current tape position. Finally, we establish a clocking variable that cycles through a sequence of values representing the read/write/move cycle of execution. This functionality is achieved through use of self-referential conditions, as in the model of Table 1.

Probabilities are extremal (i.e., 1s and 0s) in all cases. Conditional values are configured so that values change in accordance with the transitions of the Turing machine. States of the clocking variable are referenced for purposes of sequencing the individual steps of each transition. Execution of the machine then produces extremal distributions over variables that precisely replicate the read/write/move states of the Turing machine. Tape-cell variables are updated exactly as the Turing machine updates its tape.

As an illustration, consider the Turing machine defined by the state-transition table of Table 2. The function of this machine is to increment whatever binary number is represented on its tape. Each entry in the table specifies a single transition. For example, the first entry, says that in state 0 reading symbol #, the machine should write a 1, move right (R), and enter state 1. The symbol # represents an empty tape cell.

State	Read	Write	Move	New state
0	#	1	R	1
0	0	1	R	1
0	1	0	L	0
1	#	#	L	h
1	0	0	R	1
1	1	1	R	1

Table 2: Incrementing Turing Machine

Running the machine with a tape representing a binary number has the effect of producing a binary number on the tape that is one greater than the initial value. Given an initial tape with contents [# # 1 1 #], and initial read position at index 4 (i.e., over the final 1), the machine executes a series of transitions eventually producing the tape state [# 1 0 0 #].

The equivalent naive inference machine appears in Figure 3. In this translation, variable R represents the current read, W the current write, S the current state and M the current move. Variable I is the current tape index and variable K is the three-phase clock. Variables $T1$, $T2$, $T3$ etc. represent the tape

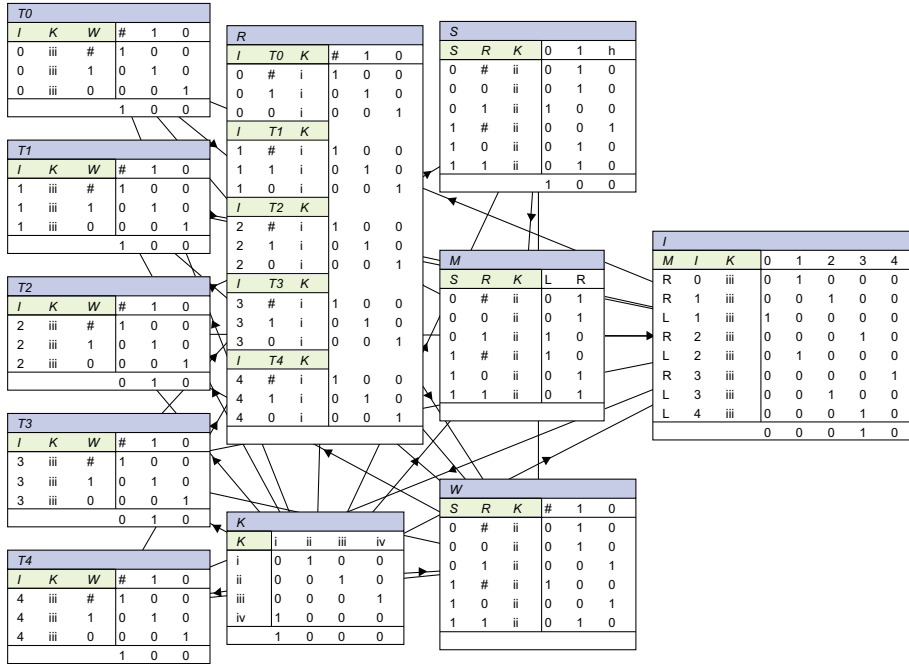


Figure 3: Naïve-inferential Turing machine.

contents at positions 1, 2, 3 etc. To correctly initialize the machine, we include unconditional distributions which have the effect of setting the clock variable to phase 1, the tape position variable to the appropriate index, and the tape variables to values corresponding to the initial contents of the tape. The ensuing behaviour then perfectly emulates the original Turing machine, terminating once the halt state is reached. Tape-state variables at that point correctly represent the final tape state of the corresponding Turing machine.

This is not the only way to translate a Turing machine into an equivalent naive inference machine. The demonstration suffices to show that the behaviour of a Turing machine can be obtained from naive inference, however. On this basis, computational behaviour is contained within naive inference, and naive inference has the capacity to ‘compute’ any computable function.

A better understanding is then obtained of the sense in which naive inference addresses the issue of inferential intractability. The protocol does not invoke an intractable algorithm, due to the fact that it does not invoke an algorithm of any sort. Rather, it invokes the concept of computation. Rather than *solving* the problem of intractability, then, naive inference provides an interpretation in which the problem does not apply. The demonstration that naive inference is able to compute any function also shows it can be the medium for representing any formal structure. On that basis, a naive inference machine can be the means

of applying inferential processes to structural representations.

5 Discussion

The paper has proposed a ‘naive’ method of probabilistic inference that deals with the problem of conditional cycles. While it cannot be viewed as approximating Bayesian estimation, it does offer some of the advantages of that approach. It avoids the problem of intractability, by linking inference to computation in general rather than to the action of a specific algorithm. There may then be implications for ‘the challenge of applying probabilistic methods over structured symbolic representations’ (Chater and Oaksford, 2008, p. 510). The present method addresses that challenge to some extent by demonstrating a way in which ‘it is possible to integrate probability with logic’ (ibid.)

Naïve inference offers a novel way to bridge the gap between the ideal of Bayesian calculation and the realities of innately constrained behaviour. Whereas the tradition of rational analysis involves modeling ‘cognitive abilities using sophisticated forms of probabilistic inference’ (Chater et al., 2006, p. 287), this approach allows them to be modeled using naive forms. Rather than assuming this necessarily involves applying some heuristic approximation of ideal Bayesian calculation, moreover, it can be taken to involve a non-heuristic procedure operating under a less sophisticated interpretation of inference.

References

- Anderson, J. R. (1990). *The Adaptive Character of Thought*. Erlbaum.
- Chater, N. and Oaksford, M. (2008). The probabilistic mind: where next? In Chater and Oaksford (Eds.), *The Probabilistic Mind: Prospects for Bayesian Cognitive Science* (pp. 501-514). Oxford: Oxford University Press.
- Chater, N. and Oaksford, M. (2008). The probabilistic mind: prospects for a bayesian cognitive science. In Chater and Oaksford (Eds.), *The Probabilistic Mind: Prospects for Bayesian Cognitive Science* (pp. 1-32). Oxford: Oxford University Press.
- Chater, N., Tenenbaum, J. and Yuille, A. (2006). Probabilistic models of cognition: conceptual foundations. *Trends in Cognitive Sciences (Special issue on Probabilistic Models of Cognition)*, 10, No. 7 (pp. 287-291).
- Cooper, G. F. (1990). The computational complexity of probabilistic inference using bayesian belief networks. *Artificial Intelligence*, 42 (pp. 393-405w).
- Danks, D. (2008). Rational analyses, instrumentalism, and implementations. In Chater and Oaksford (Eds.), *The Probabilistic Mind: Prospects for Bayesian Cognitive Science* (pp. 59-75). Oxford: Oxford University Press.

- Dayan, P. and Abbott, L. (2001). *Theoretical Neuroscience: Computational and Mathematical Modelling of Neural Systems*. MIT Press.
- Hansson, P., Juslin, P. and Winman, A. (2008). The [naive] intuitive statistician: organism—environment relations from yet another angle. In Chater and Oaksford (Eds.), *The Probabilistic Mind: Prospects for Bayesian Cognitive Science*. Oxford: Oxford University Press.
- Knill, D. and Richards, W. (1996). *Perception as Bayesian Inference*. Cambridge University Press.
- Körding, K. O. and Wolpert, D. (2006). Bayesian decision theory in sensory motor control. *Trends in Cognitive Sciences*, 10 (pp. 319-326).
- Marr, D. (1982). *Vision*. New York: W.H. Freeman.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Mateo: Morgan and Kaufman.
- Pouget, A., Dayan, P. and Zemel, R. S. (2003). Inference and computation with population codes. *Annual Review of Neuroscience*, 26 (pp. 381-410).
- Russell, S. and Norvig, P. (2010). *Artificial Intelligence: A Modern Approach* (Third Edition). Boston: Pearson.
- Simon, H. A. (1957). *Models of Man, Social and Rational: Mathematical Essays on Rational Human Behavior in a Social Setting*. New York: Wiley.