# B

## Bayesian Learning

Bradley C. Love[1], Matt Jones[2]
[1]Department of Psychology, The University of Texas at Austin, Austin, TX, USA
[2]University of Colorado, Boulder, CO, USA

## Synonyms

Bayesian model; Normative; Probabilistic approaches; Rational

## Theoretical Background

Bayesian methods have undergone tremendous progress in recent years, due largely to mathematical advances in probability and estimation theory (Chater et al. 2006). These advances have allowed theorists to express and derive predictions from far more sophisticated models than previously possible. These models have generated a good deal of excitement for at least two reasons. First, they offer a new interpretation of the goals of cognitive systems, in terms of inductive probabilistic inference, which has revived attempts at rational explanation of human behavior (Oaksford and Chater 2007). Second, Bayesian models may have the potential to explain some of the most complex aspects of human cognition, such as language acquisition or reasoning under uncertainty, where structured information and incomplete knowledge combine in a way that has defied previous approaches (e.g., Kemp and Tenenbaum 2008).

Constructing a Bayesian model involves two steps. The first step is to specify the set of possibilities for the state of the world, which is referred to as the hypothesis space. Each hypothesis can be thought of as a prediction by the subject about what future sensory information will be encountered. However, the term hypothesis should not be confused with its more traditional usage in psychology, connoting explicit testing of rules or other symbolically represented propositions. In the context of Bayesian modeling, hypotheses need have nothing to do with explicit reasoning, and indeed the Bayesian framework makes no commitment whatsoever on this issue.

For example, in Bayesian models of visual processing, hypotheses can correspond to extremely low-level information, such as the presence of elementary visual features (contours, etc.) at various locations in the visual field (Geisler et al. 2001). There is also no commitment regarding where the hypotheses come from. Hypotheses could represent innate biases or knowledge, or they could have been learned previously by the individual. Thus, the framework has no position on nativist–empiricist debates. Furthermore, hypotheses representing very different types of information (e.g., a contour in a particular location, whether or not the image reminds you of your mother, whether the image is symmetrical, whether it spells a particular word, etc.) are all lumped together in a common hypothesis space and treated equally by the model. Thus, there is no distinction between different types of representations or knowledge systems within the brain. In general, a hypothesis is nothing more than a probability distribution. This distribution, referred to as the *likelihood function*, simply specifies how likely each possible pattern of observations is according to the hypothesis in question.

The second step in constructing a Bayesian model is to specify how strongly the subject believes in each hypothesis before observing data. This initial belief is expressed as a probability distribution over the hypothesis space, and is referred to as the *prior*. The prior can be thought of as an initial bias in favor of some hypotheses over others, in that it contributes extra "votes" (as elaborated below) that are independent of any actual data. This decisional bias allows the model's predictions to be shifted in arbitrary directions regardless of the data. As we discuss below, the prior can be a strong point of the model if it is derived independently, from empirical statistics of real environments. However, more commonly, the prior is chosen ad hoc, providing substantial unconstrained flexibility to models that are advocated as rational and assumption-free.

Together, the hypotheses and the prior fully determine a Bayesian model. The model's goal is to decide how strongly to believe in each hypothesis after data have been observed. This final belief is again expressed as a probability distribution over the hypothesis space and is referred to as the *posterior*. The statistical identity known

83  as Bayes' Rule is used to combine the prior with the
84  observed data to compute the posterior. Bayes' Rule can
85  be expressed in many ways, but here we explain how it can
86  be viewed as a simple vote-counting model. Specifically,
87  Bayesian inference is equivalent to tracking evidence for
88  each hypothesis, or votes for how strongly to believe in
89  each hypothesis. The prior provides the initial evidence
90  counts, $E_{\text{prior}}$, which are essentially made-up votes that
91  give some hypotheses a head start over others, before
92  observing any actual data. When data are observed, each
93  observation adds to the existing evidence according to
94  how consistent it is with each hypothesis. The evidence
95  contributed for a hypothesis that predicted the observa-
96  tion will be greater than the evidence for a hypothesis
97  under which the observation was unlikely. The evidence
98  contributed by the $i$th observation, $E_{data_i}$, is simply added
99  to the existing evidence to update each hypothesis' count.
100 Therefore the final evidence, $E_{\text{posterior}}$, is nothing more
101 than a sum of the votes from all of the observations, plus
102 the initial votes from the prior. (Formally, $E_{\text{posterior}}$ equals
103 the logarithm of the posterior distribution, $E_{\text{prior}}$ is the
104 logarithm of the prior, and $E_{data}(H)$ is the logarithm of the
105 likelihood of the data under hypothesis $H$. The model's
106 prediction for the probability that hypothesis $H$ is correct,
107 after data have been observed, is proportional to exp
108 $[E_{\text{posterior}}(H)]$).

$$E_{\text{posterior}}(H) = E_{\text{prior}}(H) + \sum_i E_{data_i}(H) \qquad (1)$$

109 This sum is computed for every hypothesis, $H$, in the
110 hypothesis space. The vote totals determine how strongly
111 the model believes in each hypothesis in the end. Thus, any
112 Bayesian model can be viewed as tracking evidence for
113 each hypothesis, with initial evidence coming from the
114 prior and additional evidence coming from each new
115 observation. At its core, this is all there is to Bayesian
116 modeling.
117 To illustrate these two steps and how inference pro-
118 ceeds in a Bayesian model, consider the problem of deter-
119 mining whether a fan entering a football stadium is
120 rooting for the University of Southern California (USC)
121 Trojans or the University of Texas (UT) Longhorns based
122 on three simple questions: (1) Do you live by the ocean?
123 (2) Do you own a cowboy hat? (3) Do you like Mexican
124 food? The first step is to specify the space of possibilities
125 (i.e., hypothesis space). In this case, the hypothesis space
126 consists of two possibilities: being a fan of either USC or
127 UT. Both of these hypotheses entail probabilities for the
128 data we could observe, for example, $P(\text{ocean}|\text{USC}) = .8$
129 and $P(\text{ocean}|\text{UT}) = .3$. Once these probabilities are given,
130 the two hypotheses are fully specified. The second step is

131 to specify the prior. In many applications, there is no
132 principled way of doing this, but in this example, the
133 prior corresponds to the probability that a randomly
134 selected person will be a USC or a UT fan, that is, one's
135 best guess as to the overall proportion of USC and UT fans
136 in attendance.
137 With the model now specified, inference proceeds by
138 starting with the prior and accumulating evidence as new
139 data are observed. For example, if the football game is
140 being played in Los Angeles, one might expect that most
141 people are USC fans, and hence the prior would provide
142 an initial evidence count in favor of USC. If our target
143 person responded that he lives near the ocean, this obser-
144 vation would add further evidence for USC. The magni-
145 tudes of these evidence values will depend on the specific
146 numbers assumed for the prior and for the likelihood
147 function for each hypothesis, but all that the model does
148 is take the evidence values and add them up. Each new
149 observation adds to the balance of evidence among the
150 hypotheses, strengthening those that predicted it relative
151 to those under which it was unlikely.
152 There are several ways in which real applications of
153 Bayesian modeling become more complex than the simple
154 example above. However, these all have to do with the
155 complexity of the hypothesis space rather than the Bayes-
156 ian framework itself. For example, many models have
157 a hierarchical structure, in which hypotheses are essen-
158 tially grouped into higher-level *overhypotheses*.
159 Overhypotheses are generally more abstract and require
160 more observations to discriminate among; thus
161 hierarchical models are useful for modeling learning or
162 change over developmental timescales (e.g., Kemp et al.
163 2007). However, each overhypothesis is just a weighted
164 sum of elementary hypotheses, and inference among
165 overhypotheses comes down to exactly the same vote-
166 counting scheme as described above. As a second example,
167 many models assume special mathematical functions for
168 the prior, such as conjugate priors, that simplify the com-
169 putations involved in updating evidence. However, such
170 assumptions are generally made solely for the convenience
171 of the modeler, rather than for any psychological reason
172 related to the likely initial bias of a human subject. Finally,
173 for models with especially complex hypothesis spaces,
174 computing exact predictions often becomes computation-
175 ally intractable. In these cases, sophisticated approxima-
176 tion schemes are used, such as Markov-chain Monte Carlo
177 (MCMC) or particle filtering (i.e., sequential Monte
178 Carlo). These algorithms yield good estimates of the
179 model's true predictions while requiring far less compu-
180 tational effort. However, once again they are used for the
181 convenience of the modeler and usually are not meant as

182 proposals for how human subjects might solve the same
183 computational problems.

184 To summarize: Hypotheses are probability distribu-
185 tions and have no necessary connection to explicit reason-
186 ing. The model's predictions depend on the initial biases
187 on the hypotheses (i.e., the prior). The heart of Bayesian
188 inference – combining the prior with observed data to
189 reach a final prediction – is formally equivalent to
190 a simple vote-counting scheme. Learning and one-off
191 decision-making both follow this scheme, and are identi-
192 cal except for timescale and specificity of hypotheses. Most
193 of the elaborate mathematics that often arises in Bayesian
194 models comes from the complexity of their hypothesis sets
195 or the tricks used to derive tractable predictions, which
196 generally have little to do with the psychological claims of
197 the researchers. Bayesian inference itself, aside from its
198 assumption of optimality and close relation to vote-
199 counting models, does not make psychological claims in
200 recards to representational format, encoding, retrieval,
201 attention, etc. However, the flexibility and power of the
202 Bayesian framework has allowed researchers to model
203 complex learning and decision-making behaviors that
204 have proven intractable or unwieldly under other
205 formulations.

## Important Scientific Research and Open Questions

206
207

208 The restriction to computational-level accounts (cf. Marr
209 1982) severely limits contact with process-level theory and
210 data. Rational approaches attempt to explain *why* cogni-
211 tion produces the patterns of behavior that it does, but
212 they offer no insight into *how* cognition is carried out.
213 Second, in general, there are multiple rational theories of
214 any given task, corresponding to different assumptions
215 about the environment and the learner's goals. Conse-
216 quently, there is insufficient acknowledgement of these
217 assumptions and their critical roles in determining
218 model predictions. It is extremely rare to find

219 a comparison among alternative Bayesian models of the
220 same task to determine which is most consistent with
221 empirical data. Likewise, there is little recognition when
222 the critical assumptions of a Bayesian model logically
223 overlap closely with those of other theories. These chal-
224 lenges are currently being addressed by members of the
225 Bayesian community. The end goal is to integrate Bayesian
226 approaches with what we know about the mental pro-
227 cesses that support learning and decision making (Jones
228 and Love 2011).

## Cross-References

## References

239 Chater, N., Tenenbaum, J. B., & Yuille, A. (2006). Probabilistic models of
240 cognition: Conceptual foundations. *Trends in Cognitive Sciences,*
241 *10*(7), 287–291.
242 Geisler, W. S., Perry, J. S., Super, B. J., & Gallogly, D. P. (2001). Edge
243 co-occurrence in natural images predicts contour grouping perfor-
244 mance. *Vision Research, 41*, 711–724.
245 Jones, M., & Love, B. C. (in press, 2011). Bayesian fundamentalism or
246 enlightenment? On the explanatory status and theoretical contribu-
247 tions of bayesian models of cognition. *Behavioral and Brain Sciences.*
248 Kemp, C., & Tenenbaum, J. B. (2008). The discovery of structural form.
249 *Proceedings of the National Academy of Sciences, 105*, 10687–10692.
250 Kemp, C., Perfors, A., & Tenenbaum, J. B. (2007). Learning
251 overhypotheses with hierarchical Bayesian models. *Developmental*
252 *Science, 10*, 307–321.
253 Marr, D. (1982). *Vision.* San Francisco: W.H. Freeman.
254 Oaksford, M., & Chater, N. (2007). *Bayesian rationality: The probabilistic*
255 *approach to human reasoning.* Oxford: Oxford University Press.