
THEORETICAL AND REVIEW ARTICLES

Three case studies in the Bayesian analysis of cognitive models

MICHAEL D. LEE

University of California, Irvine, California

Bayesian statistical inference offers a principled and comprehensive approach for relating psychological models to data. This article presents Bayesian analyses of three influential psychological models: multidimensional scaling models of stimulus representation, the generalized context model of category learning, and a signal detection theory model of decision making. In each case, the model is recast as a probabilistic graphical model and is evaluated in relation to a previously considered data set. In each case, it is shown that Bayesian inference is able to provide answers to important theoretical and empirical questions easily and coherently. The generality of the Bayesian approach and its potential for the understanding of models and data in psychology are discussed.

Psychology as an empirical science progresses through the development of formal models incorporating theoretical ideas designed to explain and predict observations of psychological phenomena. This means that progress in psychology relies upon the quality and completeness of the methods it uses to relate models and data. There is little point in developing theories and models, on the one hand, and collecting data in the laboratory or the field, on the other, if the two cannot be brought into contact in useful ways.

In most empirical sciences, Bayesian methods have been or are rapidly being adopted as the most complete and coherent available way to relate models and data. Psychology has long been aware of problems with traditional frequentist and null hypothesis significance-testing approaches to parameter estimation and model selection, and recognition of the Bayesian alternative has followed from a number of recent articles and special volumes addressing the general issues (e.g., Lee & Wagenmakers, 2005; Myung, Forster, & Browne, 2000; Myung & Pitt, 1997; Pitt, Myung, & Zhang, 2002). Beyond the illustrative applications provided in these general treatments, however, there are few worked examples of Bayesian methods being applied to models at the forefront of modern psychological theorizing. Perhaps one reason is that there has been too great a focus on model selection defined in a narrow sense—particularly through the evaluation of Bayes factors—rather than a full Bayesian analysis. The perception that all that Bayesian methods have to offer for the evaluation of psychological models is a number that quantifies how much more likely one model is than another is dangerously limiting.

In this article, three previous cognitive-modeling studies are revisited, in an attempt to demonstrate the generality

and usefulness of the Bayesian approach. The three applications involve the multidimensional scaling (MDS) representation of stimulus similarity (Shepard, 1962, 1980), the generalized context model (GCM) account of category learning (Nosofsky, 1984, 1986), and a signal detection theory (SDT) account of inductive and deductive reasoning (Heit & Rotello, 2005). These applications were chosen in order to span a range of cognitive phenomena, to involve well-known and influential theories, and to put a focus on the ability of Bayesian methods to provide useful answers to important theoretical and empirical questions.

METRIC MULTIDIMENSIONAL SCALING

Theoretical Background

MDS representations of stimuli use a low-dimensional metric space in which points correspond to stimuli and the distance between points models the (dis)similarity between stimuli (Shepard, 1957, 1962, 1987, 1994). Nonmetric varieties of MDS algorithms for inferring these representations from pairwise similarity data (e.g., Kruskal, 1964) make only weak assumptions about the form of the relationship between distance in the MDS space and stimulus similarity. However, Shepard's (1987) *universal law of generalization* provides a compelling case that similarity decays exponentially with distance, at least for relatively low-level perceptual stimulus domains. We make this assumption¹ and, so, will consider the form of metric MDS that uses an exponential decay function to relate distances to similarities.

A classic issue in all MDS modeling has involved the interpretation of different metric assumptions for the representational space. Typically, consideration is restricted to

M. D. Lee, mdlee@uci.edu

the Minkowskian family of distance metrics. For points $\mathbf{p}_i = (p_{i1}, \dots, p_{iD})$ and $\mathbf{p}_j = (p_{j1}, \dots, p_{jD})$ in a D -dimensional space, the Minkowski r -metric distance is given by

$$d_{ij} = \left[\sum_{x=1}^D |p_{ix} - p_{jx}|^r \right]^{1/r}. \quad (1)$$

The $r = 1$ (city block) and $r = 2$ (Euclidean) cases are usually associated with so-called separable and integral stimulus domains, respectively (Garner, 1974; Shepard, 1991). The basic idea is that many stimulus domains, such as different shapes or different sizes, have component dimensions that can be attended to separately. These are termed separable and are well modeled by the distance metric that treats each dimension independently in accruing distance. Other stimulus domains, such as color, however, have component dimensions that are fused and not easily distinguished, and so the comparison of stimuli involves all of the dimensions simultaneously. These are termed integral and are well modeled by the familiar Euclidean distance metric. In addition, metrics with $r < 1$ have been given a psychological justification (e.g., Gati & Tversky, 1982; Shepard, 1987, 1991) in terms of modeling stimuli with component dimensions that “compete” for attention.²

Despite the theoretical elegance of this framework for relating the Minkowskian metric family to core psychological properties of stimulus domains, there have been few attempts to infer r from similarity data by using MDS modeling. Shepard (1991) has presented a focused attack on the problem that gives a good account of the capabilities and pitfalls of using standard methods. The basic approach is a brute force one of applying standard nonmetric MDS algorithms assuming a large number of different r values and comparing the solutions on the basis of a measure of goodness of fit.

Besides the set of computational problems that are noted by Shepard (1991), which are severe enough to preclude even considering the theoretically interesting possibilities with $r < 1$, this approach suffers from failing to account for the functional form effects of model complexity inherent in varying the metric parameter. Since the value of r dictates how the coordinate location parameters interact, different values of r will certainly change the functional form of parametric interaction and, hence, the complexity of the metric space representational model being considered. One of the great attractions of Bayesian inference is that, through its basis in a coherent and axiomatized probabilistic framework for inference, model complexity issues such as these are automatically handled in a principled way.

Graphical Model for MDS

All of the Bayesian models in this article rely on posterior sampling from graphical models (see Griffiths, Kemp, & Tenenbaum, in press, and Jordan, 2004, for psychological and statistical introductions, respectively). In these models, nodes represent variables of interest, and the graph structure is used to indicate dependencies between the variables, with children depending on their parents. The conventions of representing continuous variables with circular nodes

and discrete variables with square nodes and of representing unobserved variables without shading and observed variables with shading are used. Stochastic and deterministic unobserved variables are distinguished by using single and double borders, respectively. Plate notation, enclosing with square boundaries subsets of the graph that have independent replications in the model, is also used.

Figure 1 presents a graphical model interpretation of metric MDS. At the top is the coordinate representation of the points corresponding to stimuli. The p_{ix} node corresponds to the single coordinate value of the i th stimulus on the x th dimension, and the surrounding plates repeat these coordinates over the $i = 1, \dots, N$ stimuli and $x = 1, \dots, D$ dimensions. The node is shown as single-bordered, without shading, and circular because each coordinate dimension is stochastic, unknown, and continuous, respectively. Under the Bayesian approach, a prior distribution for these coordinate location parameters must be specified. We make the obvious prior assumption that all of the coordinates have equal prior probability of being anywhere in a sufficiently large (hyper)cube with bounds $(-\delta, +\delta)$:

$$p_{ix} \sim \text{Uniform}(-\delta, \delta); \delta > 0, \quad (2)$$

where “sufficiently large” means large enough that increasing δ does not alter the posterior distribution over the coordinate point parameters.

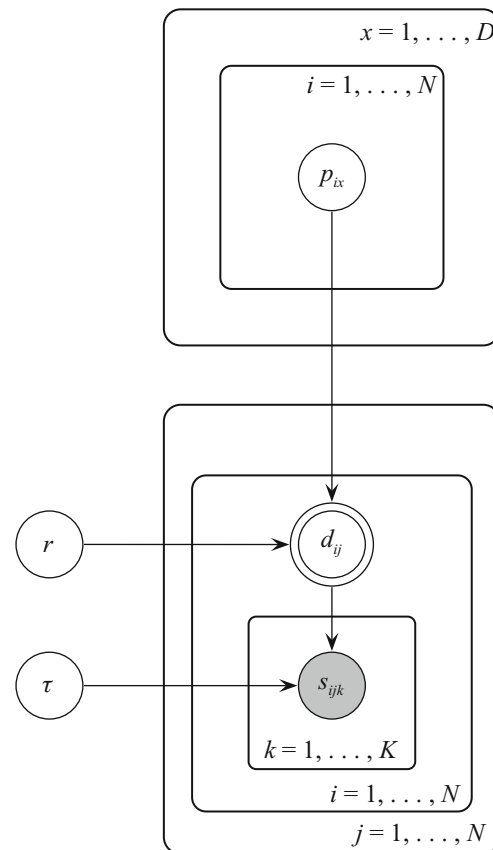


Figure 1. Graphical model for metric multidimensional scaling.

The metric parameter r that is the focus of this application is also a stochastic, unobserved, and continuous node. Following the earlier discussion, a prior distribution is used that is uniform over the theoretically interpretable interval between zero and two:

$$r \sim \text{Uniform}(0, 2). \quad (3)$$

Given the value of r and the coordinate locations p_{ix} , the pairwise distances d_{ij} are automatically given by Equation 1. In the graphical model in Figure 1, this means that the d_{ij} node is double-bordered, to indicate that it is deterministic, and has as parents the r and p_{ix} nodes. The d_{ij} node is encompassed by two plates, $i = 1, \dots, N$ and $j = 1, \dots, N$, to express the repetition over all pairs of N stimuli.

The similarity data considered here provide similarity ratings for each pair of stimuli as generated independently by K participants. The observed similarity between the i th and the j th stimuli given by the k th participant is denoted s_{ijk} and, so, is enclosed by an additional plate representing the $k = 1, \dots, K$ participants. These similarities are assumed to be generated as the exponential decay of the distance between these points but are subject to noise and, so, are stochastic, observed, and continuous. The noise process is assumed to be a zero-mean Gaussian with common variance across all participants and stimulus pairs. The precision (i.e., the reciprocal of the variance) is represented by the stochastic, unobserved, and continuous parameter τ , so that

$$s_{ijk} \sim \text{Gaussian}\left(\exp(-d_{ij}), \tau\right), \quad (4)$$

with the standard (see Spiegelhalter, Thomas, Best, & Gilks, 1996) near noninformative prior distribution for the precision

$$\tau \sim \text{Gamma}(\varepsilon, \varepsilon), \quad (5)$$

where $\varepsilon = .001$ is set near zero.³

Posterior Inference in Graphical Models

The graphical model in Figure 1 defines a precise and complete probabilistic relationship between the MDS parameters—the coordinate locations of stimulus points and the metric parameter—and the observed similarity data. Bayesian inference uses this relationship to update what is known about the parameters, converting prior distributions to posterior distributions on the basis of the evidence provided by data.

The graphical model is a generative one, specifying how stimulus points and a metric combine to produce similarity data. Once similarity data are observed, inference is the conceptually easy process of reversing the generative process and working out what stimulus points and metric are likely to have produced the data. The posterior probability distribution represents this information, specifying the relative probability that each possible combination of stimulus points and metric is the one that generated the data.

Although conceptually straightforward, for most interesting cognitive models, it will not be possible to find the posterior distribution analytically, and it is also unlikely that standard approximations will be very useful. Mod-

ern Bayesian inference for complicated models proceeds computationally by drawing samples from the posterior distribution. We implement our graphical models using WinBUGS (Spiegelhalter, Thomas, & Best, 2004), which uses a range of Markov chain Monte Carlo computational methods, including adaptive rejection sampling, slice sampling, and Metropolis–Hastings (see, e.g., Chen, Shao, & Ibrahim, 2000; Gilks, Richardson, & Spiegelhalter, 1996; Mackay, 2003), to perform posterior sampling.

For the MDS model in Figure 1, each posterior sample lists values for the unobserved variables

$$(r, \tau, p_{11}, \dots, p_{ND}, d_{11}, \dots, d_{NN}). \quad (6)$$

The basic principle of posterior sampling is that, over a large number of samples, the relative frequency of a particular combination of parameter values appearing corresponds to the relative probability of those values in the posterior distribution. This correspondence allows the information that is conceptually in the exact joint posterior distribution to be accessed approximately by simple computations across the posterior samples. For example, a histogram of the sampled values of a variable approximates its marginal posterior distribution, and the arithmetic average over these values approximates its expected posterior value. Considering the sampled values of one variable, for only those samples where another variable takes a specific value, corresponds to considering a conditional probability. Considering the combination of values taken by two or more variables corresponds to considering their joint distribution, and so on.

Inference From MDS Data

Our MDS applications consider three sets of individual-participant similarity data. Initial investigations with averaged data, of the type considered by Shepard (1991), showed clearly that the repeated measures nature of individual-participant data was important for making sound inferences about the metric structure of the representational space. This is consistent with results showing that averaging similarity data with individual differences can systematically affect the metric structure of MDS spaces (see Ashby, Maddox, & Lee, 1994; Lee & Pope, 2003). Only three data sets could be found for which raw individual-participant data were available and for which reasonable predictions about the separability or integrality of the stimulus domain could be made.

The first of these related to rectangles of different height with interior line segments in different positions, using eight of the possibilities in a 4×4 factorial design, as reported by Kruschke (1993). The second related to circles of different sizes with radial lines at different angles, following (essentially) a 3×3 factorial design, as reported by Treat, MacKay, and Nosofsky (1999). The third related to 10 spectral colors, as reported by Helm (1959). Previous results would strongly suggest the first two of these domains are separable, whereas the colors are integral. Also, on the basis of previous analyses (e.g., Lee, 2001; Shepard, 1962) and the explicit two-factor combinatorial designs for two of these stimulus domains, a two-dimensional representational space was assumed for all three stimulus domains.

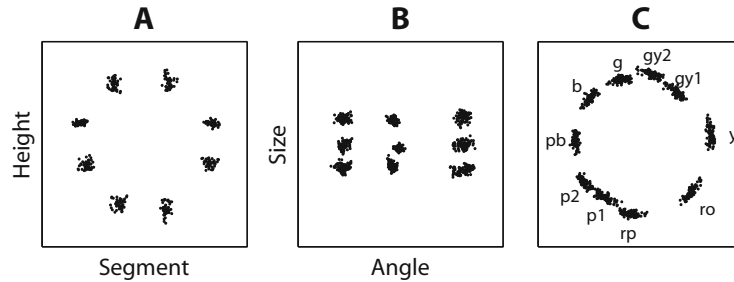


Figure 2. Multidimensional scaling representations for three stimulus domains—relating to (A) rectangles with interior lines, (B) circles with radial lines, and (C) spectral colors—showing samples from the posterior distributions for the representational points.

For each data set, we calculated 5,000 such samples after a 1,000-sample “burn-in” period (i.e., a period of sampling that is not recorded but allows the Markov chain to converge to sampling from the posterior distribution). We used multiple chains to check convergence and observed a small proportion of these chains showing a degenerate behavior, with r becoming trapped near zero. Although this behavior requires an explanation in the future, these chains were removed from the present analysis. Postprocessing of the posterior samples for the coordinate location parameters was also required, to accommodate natural translation, reflection, and axis permutation invariances inherent in the MDS model. We achieved this by translating to center at the origin, reflecting where necessary so that both coordinate values for the first stimulus were positive, and permuting the axes where necessary so that the first coordinate value was larger than the second.

Figure 2 shows 50 randomly selected posterior samples for the stimulus points, displaying the representations that have been inferred from each data set. In panel A, the eight stimuli are appropriately located within the 4×4 factorial structure. In panel B, the nine stimuli follow the exhausted 3×3 factorial structure. In panel C, the stimuli follow the standard *color circle* representation. In each case, showing samples from the distribution also gives a natural visual representation of the uncertainty in the coordinate locations.

Figure 3 shows the posterior distribution over the metric parameter r for each of the three data sets. It is clear that the color stimulus domain, which is expected to be integral, is distributed between about 1.6 and 2.0. It is not clear whether the mode is slightly below 2.0 because, consistent with previous theorizing, full integrality is not achieved or as a consequence of the theoretically driven restriction⁴ that r not exceed 2.0. The posterior distribution of r for the radial lines and circles domains is centered about the value of 1.0 associated with separability, as would be expected. The rectangle with interior lines stimulus domain has a posterior that lies a little below 1.0. One plausible interpretation of this, again consistent with previous theorizing (e.g., Gati & Tversky, 1982; Shepard, 1991), is that the rectangles and lines compete for attention and constitute a “highly separable” stimulus domain.

Summary

This application considered a Bayesian formulation of metric MDS modeling for similarity-based representation. The formulation was not intended to be definitive: It neglected issues of dimensionality determination and made plausible, but contestable, assumptions about the form of the generalization gradient and the distribution of empirical similarity. The present model also assumed that there are no individual differences in the stimulus representations for different participants. All of these issues await fuller exploration within a Bayesian graphical model framework.

What the application does show, however, is that even with this simple formulation, a Bayesian approach automatically provides useful information not available in previous analyses. It provides a full posterior distribution over the location of the points representing stimuli, without the need to make distributional assumptions about these posteriors, as with many probabilistic MDS methods (e.g., Ramsay, 1982). Under the sampling approach to Bayesian inference, posterior distributions are not constrained to follow any particular distribution but are free to take the form

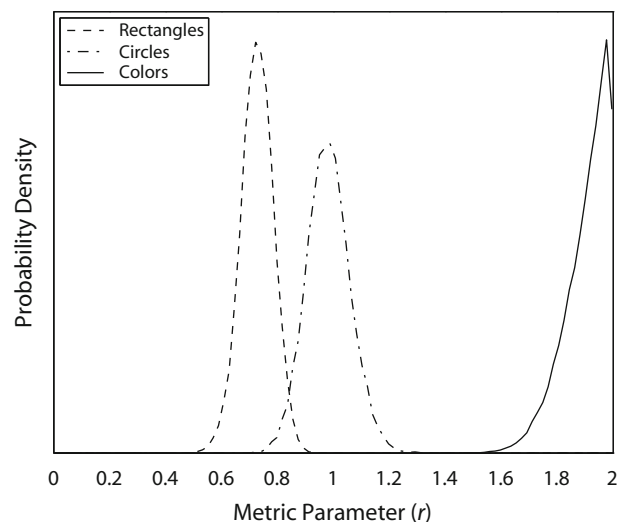


Figure 3. Posterior distributions over the metric parameter r for three stimulus domains.

that follows from the specification of the model and the information provided by data. The present application also provides a full posterior distribution over the parameter indexing the metric structure of the space. This posterior is sensitive to differences in the functional form complexity of parameter interaction, unlike the approach based on goodness of fit originally considered by Shepard (1991).

CATEGORY LEARNING

Theoretical Background

The GCM is a highly influential model of exemplar-based category learning (Nosofsky, 1984, 1986). The model assumes that stimuli are represented as exemplars in memory according to a previously derived MDS representation, which is subject to a selective attention process that weights the dimensions of the representation. The similarity between stimuli is modeled as an exponentially decaying function of distance in this transformed space, using a generalization gradient parameter. Category decisions are made probabilistically according to the ratio of similarity between the presented stimulus and those in the different categories, using bias parameters that weight the different category responses.

Nosofsky (1986) presented a thorough and impressive study of the performance of the GCM on individual-participant data in related identification and two-category learning tasks. Of particular interest are the categorization tasks, which involved four different two-category structures, termed *dimensional*, *criss-cross*, *interior–exterior*, and *diagonal*. These category structures are shown in Figure 4. The stimuli are arranged in a 4×4 grid, corresponding to their MDS representation. For each structure, the eight stimuli assigned to the two categories are shown as four black and four white squares in the grid, and the unassigned stimuli are shown as gray squares.

Of the many modeling issues Nosofsky (1986) addressed using these category structures, our focus will be on two. The first is an estimation issue and relates to the values of the attention, generalization, and bias parameters. In several places, theoretical questions are directly addressed by knowledge of the values that these parameters take, and Nosofsky (1986) reported standard tests of significance to decide, for example, whether attention was equally distributed over the two components of the stimuli. The posterior distribution over these parameters obtained automatically by Bayesian analysis contains the relevant information for addressing these sorts of inferences.

The second issue is a model selection issue and relates to the augmented version of the GCM proposed by Nosofsky (1986). In this augmented version, not only are the stimuli presented in the category-learning task used in assessing similarity, but the other stimuli, shown in gray in Figure 4, from the domain encountered in the earlier identification task are also used. These additional stimuli are assumed to have a latent assignment to one of the categories. What inferences can be made about these assignments is not readily amenable to standard statistical testing, and so Nosofsky (1986) considered every possible pattern of latent assignment to draw conclusions. Whether the improved fit of this augmented GCM over that of the original version warrants the additional model complexity is also a difficult question to answer using standard methods. Nosofsky (1986, p. 48) acknowledged this difficulty and argued for the appropriateness of the augmented model for just the criss-cross and interior–exterior category structures on the basis of unspecified “computer simulations.”

Graphical Model

Figure 5 presents a graphical model interpretation of the augmented GCM, as applied to the two-dimensional stimulus domain in Nosofsky (1986). The x_i and θ nodes relate only to model comparison applications and will be explained in that section. At the top of Figure 5 are the observed MDS coordinate locations for the $i = 1, \dots, N$ stimuli in $x = 1, 2$ dimensions. The attention weight parameter gives the relative emphasis given to the first stimulus dimension over the second. This weight is given a uniform prior distribution over the interval between zero and one:

$$w \sim \text{Uniform}(0, 1). \quad (7)$$

The version of the GCM used in Nosofsky (1986) models similarity as an exponentially decaying function of the squared distance between the representative points. These squared distances are represented by the d_{ij}^2 node, which is deterministically defined in terms of the attention weight and coordinates:

$$d_{ij}^2 = w(p_{i1} - p_{j1})^2 + (1-w)(p_{i2} - p_{j2})^2. \quad (8)$$

Given these squared distances, the generalization gradient parameter c determines the similarities between each pair of stimuli:

$$s_{ij} = \exp\left[-(cd_{ij})^2\right]. \quad (9)$$

Dimensional	Criss-Cross	Interior–Exterior	Diagonal
13 14 15 16	13 14 15 16	13 14 15 16	13 14 15 16
9 10 11 12	9 10 11 12	9 10 11 12	9 10 11 12
5 6 7 8	5 6 7 8	5 6 7 8	5 6 7 8
1 2 3 4	1 2 3 4	1 2 3 4	1 2 3 4

Figure 4. The four category structures used in the Nosofsky (1986) study. Based on Nosofsky (1986, Figure 5).

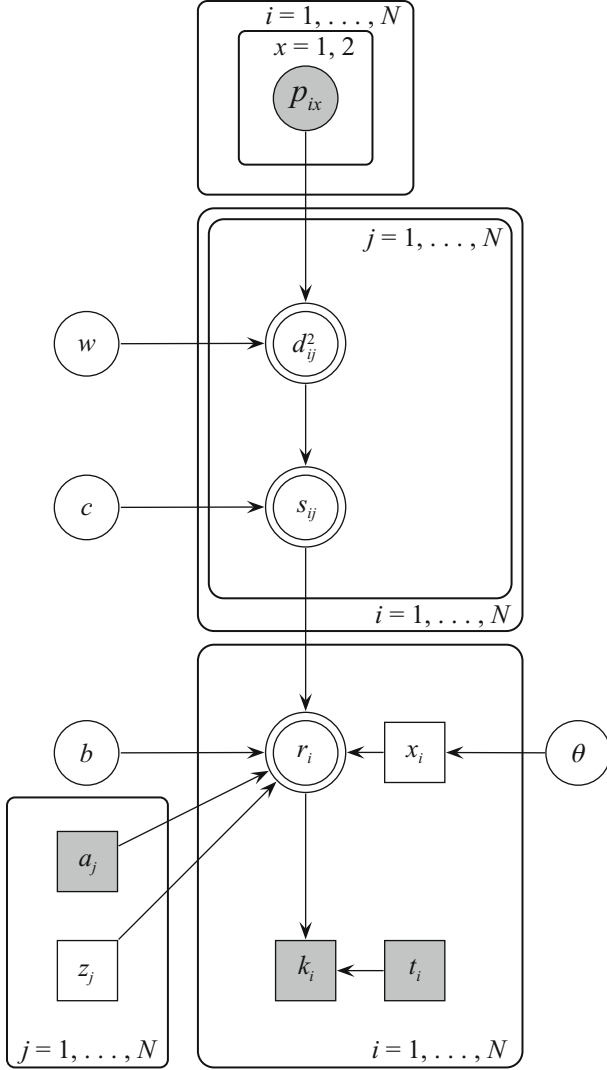


Figure 5. Graphical model for the augmented generalized context model.

The c parameter functions as an inverse scale (i.e., $1/c$ scales the distances), implying c^2 functions as a precision, and is given the standard near noninformative prior:

$$c^2 \sim \text{Gamma}(\varepsilon, \varepsilon), \quad (10)$$

where $\varepsilon = .001$ is set near zero. Both the distances and the similarities are repeated across all pairs of stimuli and, so, are enclosed in two plates.

The probability of responding to the i th stimulus is determined by the similarities between the stimuli, the response bias, and the assignment of the stimuli to the categories. The response bias b is stochastic, unobserved, and continuous and is given a uniform prior distribution over the interval between zero and one:

$$b \sim \text{Uniform}(0, 1). \quad (11)$$

The assignment of stimuli to the two categories is determined by two indicator variables. The indicator variables a_j represent the known assignment of the j th presented

stimulus, ranging over the $N/2$ such stimuli in each category structure. The indicator variables z_j represent the latent assignment of the j th unassigned stimulus, ranging over the $N/2$ such stimuli in each category structure. The latent variables are stochastic and are given a Bernoulli prior distribution with a rate of $1/2$, making the stimuli equally likely a priori to be assigned to each category:

$$z_j \sim \text{Bernoulli}(1/2). \quad (12)$$

From the similarities, bias, and assignments, the response probability for the i th stimulus being chosen as a member of the first category (“Category A”) is

$$r_i = \frac{b \left[\sum_{a \in A} s_{ia} + \sum_{z \in A} s_{iz} \right]}{b \left[\sum_{a \in A} s_{ia} + \sum_{z \in A} s_{iz} \right] + (1-b) \left[\sum_{a \in B} s_{ia} + \sum_{z \in B} s_{iz} \right]}. \quad (13)$$

The GCM uses these response probabilities to account for the observed data, which are the counts, k_i , of the number of times the i th stimulus was chosen in Category A out of the t_i trials on which it was presented. Accordingly, the counts k_i follow a binomial distribution:

$$k_i \sim \text{Binomial}(r_i, t_i). \quad (14)$$

Inference From Data

For each of the four category-learning data sets for Subject 1 in Nosofsky (1986), 100 chains were run collecting 1,000 posterior samples drawn after a burn-in of 1,000 samples. Each of the chains used a different random initial pattern of assignment. For two of the category structures—the dimensional and interior–exterior ones—a single pattern of latent assignment was observed to dominate the posterior. These patterns are shown, together with the original category structures, in Figure 6 and match exactly those reported by Nosofsky (1986, Table 5).

Given the consistency in latent assignments, it is straightforward to interpret the posterior distributions for the attention, generalization, and bias parameters for each category structure, as is shown in Figure 7. These distributions are entirely consistent with the maximum-likelihood estimates reported by Nosofsky (1986, Table 5). The posterior distributions in Figure 7, however, carry useful additional information, since they provide a complete characterization of the uncertainty in knowledge of each parameter. It is clear, for example, that attention in the interior–exterior condition has significant density at the theoretically important value of 0.5. It is also clear that the dimensional and interior–exterior bias and generalization parameters are very likely to be different, since their posterior distributions do not significantly overlap.

Finally, it is worth noting that the posterior of the attention parameter for the dimensional condition shows how Bayesian methods naturally handle the theoretical restriction of their range. Frequentist confidence intervals based on asymptotic assumptions are unlikely to be suitable for inference in cases such as these, and more difficult and ad hoc methods, such as bootstrapping, would probably

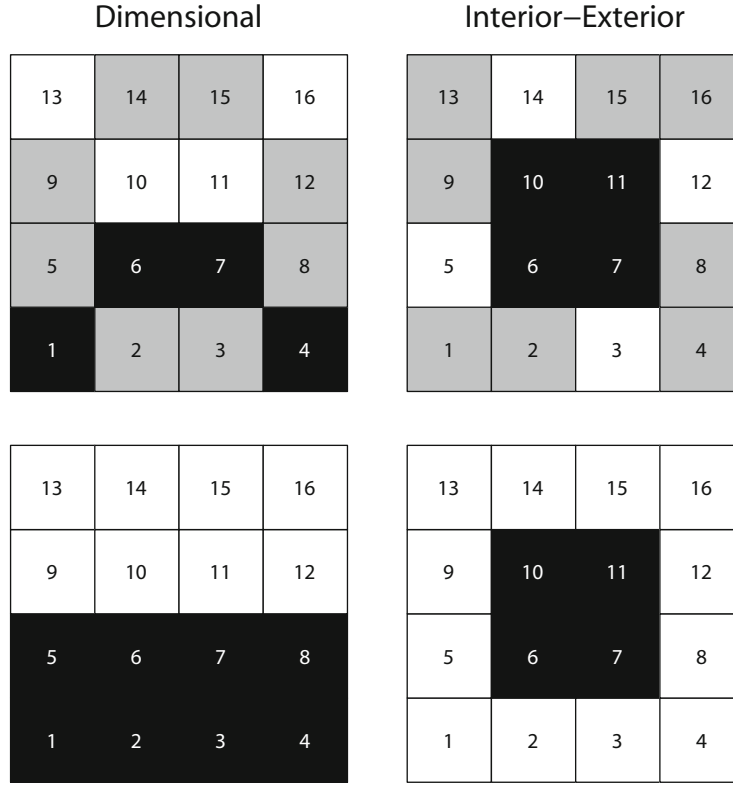


Figure 6. The augmented generalized context model latent assignments for the dimensional and interior–exterior structures.

be required. In contrast, the Bayesian approach handles the constraint automatically and naturally.

Across the 100 chains, the posterior samples for the criss-cross and diagonal category structures revealed two and four different patterns, respectively, of latent stimulus assignment.⁵ These latent assignments are shown in Figure 8.

Only the latent assignments CI, DII, and DIV were identified by Nosofsky (1986). Considering the additional latent assignments for the diagonal structure is particularly satisfying, because the total of four patterns exhausts the possibilities for Stimuli 7 and 10. In this way, the two newly found assignments complement and complete those already established. Of course, it is possible that by lowering the threshold of goodness of fit used to find latent assignments, Nosofsky (1986) could also have found these new assignments. But it is important to understand that unlike the Bayesian results reported here, such an analysis would not be sensitive to differences in the complexity of different assignments. Formally, Nosofsky (1986) considered patterns of latent assignment on the basis of their *maximum* likelihood, $p(z|w^*, c^*, b^*, D)$, whereas the

posterior samples in Figure 8 come from the *marginal* likelihood $p(z|D)$, where D is the category-learning data. Only the marginal density accounts for model complexity, because it takes into consideration how likely a category representation is, averaged across all of the different possible values for attention, generalization, and bias.

Model Comparison

To address the model comparison issue of whether the additional complexity involved in allowing latent assignments in the augmented GCM is justified by the data, the x_i and θ nodes in Figure 5 are used. The x_i nodes are latent indicator variables for each of the i stimuli being categorized. These indicators determine whether or not the associated response probability r_i uses the latent stimulus assignments z_j as per the augmented GCM, or simply relies on the fixed assignments from the category-learning task given by a_j .

Formally, this extension can be expressed by updating Equation 13 to Equation 15, as shown at the bottom of the page.

$$r_i = \begin{cases} \frac{b \sum_{a \in A} S_{ia}}{b \sum_{a \in A} S_{ia} + (1-b) \sum_{a \in B} S_{ia}} & \text{if } x_i \text{ is 0} \\ \frac{b \left[\sum_{a \in A} S_{ia} + \sum_{z \in A} S_{iz} \right]}{b \left[\sum_{a \in A} S_{ia} + \sum_{z \in A} S_{iz} \right] + (1-b) \left[\sum_{a \in B} S_{ia} + \sum_{z \in B} S_{iz} \right]} & \text{if } x_i \text{ is 1.} \end{cases} \quad (15)$$

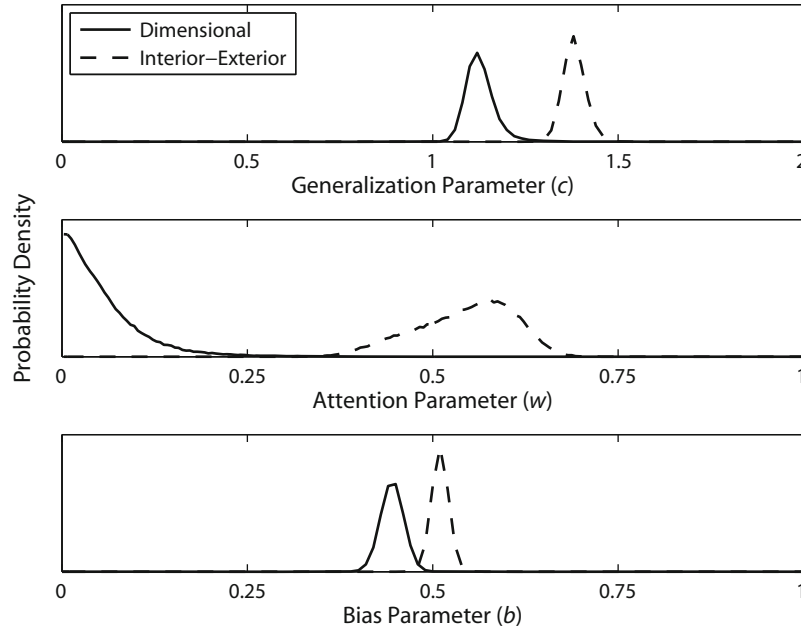


Figure 7. Posterior distributions for the augmented generalized context model parameters for two category-learning structures.

All of the indicators x_i are assumed to support the standard and augmented GCM accounts according to a fixed underlying rate of use, θ . This rate of use is given a uniform prior distribution,

$$\theta \sim \text{Uniform}(0,1), \quad (16)$$

and its posterior provides a measure of the relative usefulness of the standard and augmented GCM accounts.

The posterior rate of use provides a measure of the relative importance of the two models in accounting for the way the participant categorized all of the stimuli. The nature of the measure is best understood by noting its relationship to the standard Bayes factor (see Kass & Raftery, 1995). If θ were given a prior that only allowed the possibility that *every* stimulus was categorized by the standard GCM, or *every* stimulus was categorized by the augmented GCM, the posterior distribution would naturally allow the estimation of the Bayes factor. That is, the Bayes factor is a form of mixture estimation, when the only possible mixing rates are zero and one, because exactly one of the models is true. The assumption of a uniform prior employed here corresponds to allowing the possibility that neither model is exactly and exclusively true but both might be useful, and the issue of relative merit is the issue of what mixture of standard to augmented GCM can be inferred from the data.

This is the information provided by the posteriors for rate of use shown in Figure 9 for the four category structures. It is clear that the augmented GCM is rarely used for the dimensional category structure but is used significantly often for the other three structures, particularly in the case of the criss-cross and interior–exterior structures. In general, exactly what rate of use is required before a model is declared necessary, or superior to a competitor, is a question of the standards of scientific evidence needed and must be made

by researchers in each specific context. For the present application, we would conclude from Figure 9 that the augmented GCM is a useful and justified theoretical extension for all but the dimensional category structure.

Summary

This application demonstrated the ability of Bayesian methods to improve both parameter estimation and model comparison for the GCM account of category learning. The parameter posterior distributions provide a complete representation of what is known and unknown about the psychological variables—selective attention, stimulus generalization, and response bias—used by the model to explain the observed category-learning behavior. Under the posterior sampling approach to Bayesian inference, these distributions are again not constrained to follow any particular distribution but are free to take the form that follows from the specification of the model and the information provided by data.

A different sort of parameter estimation is demonstrated by the patterns of latent assignment in Figure 6. The augmented version of the GCM involves an additional set of membership parameters, which indicate the assignment of untrained stimuli to the two categories. In contrast to the difficulties encountered by Nosofsky (1986) with standard methods, Bayesian inference applies exactly the same principles to estimating these discrete parameters as it does for the continuous attention, generalization, and bias parameters. For two of the category structures, the Bayesian analysis showed the same augmented assignments as those originally found by Nosofsky (1986), but for the criss-cross and diagonal structures, it showed additional and intuitively satisfying patterns of associating untrained stimuli with the categories, and made these inferences with sensitivity to model complexity.

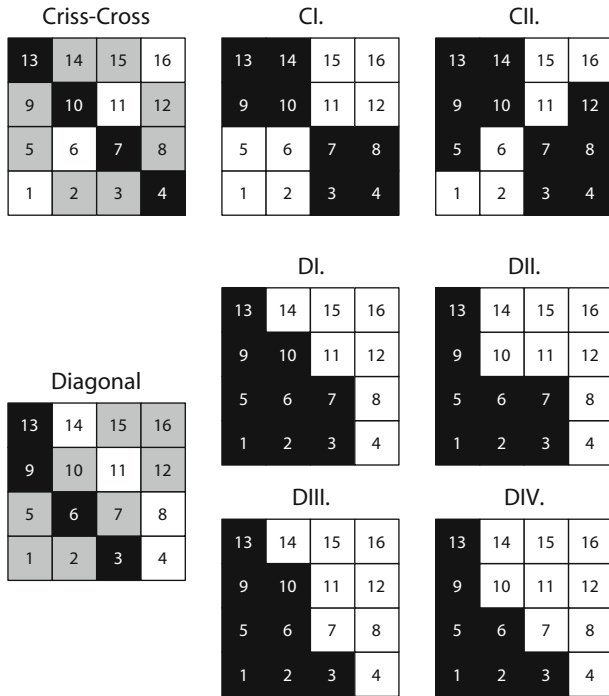


Figure 8. The augmented generalized context model latent assignments for the criss-cross and diagonal category structures.

Finally, this application shows how Bayesian inference can provide answers for a difficult model selection problem that was not addressed in any formal way by Nosofsky (1986). Using a mixture modeling approach to compare the standard and the augmented GCM accounts, strong evidence was found for the additional complexity of the augmented account for three of the four category structures.

SIGNAL DETECTION MODEL OF REASONING

Theoretical Background

Heit and Rotello (2005) have presented a clever model-based evaluation of the conjecture that both inductive and deductive reasoning involve the same single psychological dimension of *argument strength* (Rips, 2001). Heit and Rotello used SDT (for detailed treatments, see Green & Swets, 1966; Macmillan & Creelman, 2005) to model this conjecture. The basic idea is to assume that the strength of an argument is unidimensional but that different decision criteria control inductive and deductive reasoning. In particular, a relatively lesser criterion of argument strength is assumed to decide between *weak* and *strong* arguments for induction, whereas a relatively greater criterion decides between *invalid* and *valid* arguments for deduction. Under this conception, deduction is simply a more stringent form of induction. Accordingly, empirical evidence for or against the SDT model has strong implications for the many-threaded contemporary debate over the existence of different kinds of reasoning systems or processes (e.g., Chater & Oaksford, 2000; Heit, 2000; Parsons & Osherson, 2001; Sloman, 1998).

To obtain empirical evidence for evaluating the SDT model, Heit and Rotello (2005) tested the inductive and deductive judgments of 80 participants on eight arguments. They used a between-subjects design, so that 40 participants were asked induction questions about the arguments (i.e., whether or not the conclusion was “plausible”), whereas the other 40 participants were asked deduction questions (i.e., whether or not the conclusion was “necessarily true”). For each participant, there were four *signal* questions, where the conclusions were plausible or necessarily true, and four *noise* questions, where the conclusions were not plausible or necessarily true. Accordingly, the decisions made by the participants had a natural characterization in terms of hit and false alarm rates, which could then be converted to standard measures of discriminability (or synonymously, sensitivity) and bias using SDT.

In one of the key analyses of Heit and Rotello (2005), standard significance testing was used to reject the null hypothesis that there was no difference between discriminability for induction and deduction conditions. Their analysis involved calculating the mean discriminabilities for each participant, using edge corrections where perfect performance was observed. These sets of discriminabilities gave means of 0.93 for the induction condition and 1.68 for the deduction condition. By calculating via the *t* statistic—and so, assuming associated Gaussian sampling distributions—and observing that the *p* value was less than .01, Heit and Rotello rejected the null hypothesis of equal means. According to Heit and Rotello, this finding of different discriminabilities provided evidence against the criterion-shifting unidimensional account offered by SDT.

Although the statistical inference methods used by Heit and Rotello (2005) are widely used and accepted, they explicitly or implicitly make a number of problematic assumptions that can be dealt with effectively by using the Bayesian approach. First, the uncertainty about the

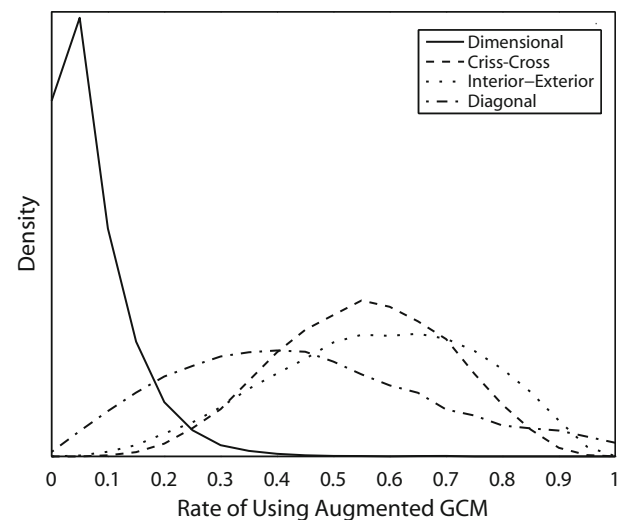


Figure 9. Posterior distribution of the rate at which stimuli are assigned to the augmented generalized context model (GCM), rather than to the standard GCM, for each of the four category structures.

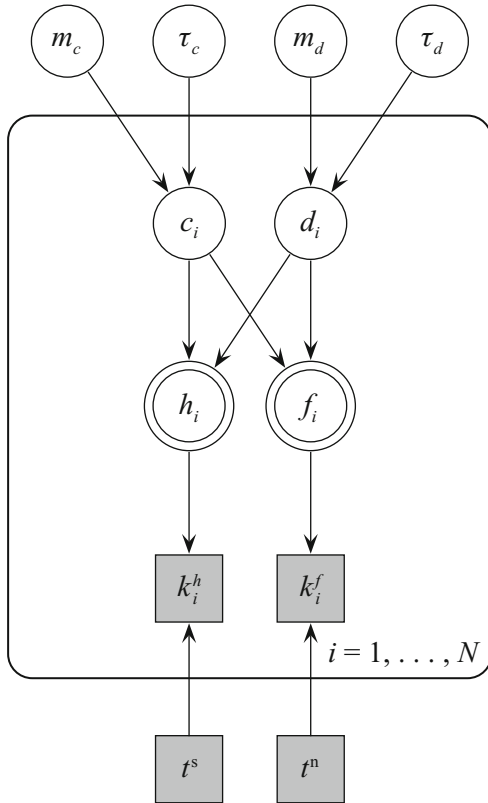


Figure 10. Graphical model for signal detection theory analysis allowing for Gaussian variation in discriminability and bias across participants.

discriminability of each individual is ignored, since it is represented by a single point estimate. Intuitively, making decisions corresponding, for example, to three hits and one false alarm is consistent, to varying degrees, with a range of possible hit and false alarm rates and, hence, to varying degrees, with a range of discriminabilities. The Bayesian approach naturally represents this uncertainty by making prior assumptions about hit and false alarm rates and then using the evidence provided by the decisions to calculate posterior distributions. These posterior distributions are naturally mapped into posterior distributions for discriminability and bias according to SDT, which avoids the need for ad hoc edge corrections.

In addition, and perhaps more important, the statistical analyses undertaken by Heit and Rotello (2005) implicitly assume that there are no individual differences across participants within each condition. The mean discriminabilities for each group they calculate are based on the statistical assumption that there is exactly one underlying point that generates the behavior of every participant in that group. That is, all of the individual-participant data are used to estimate a single discriminability, with a standard error representing only the uncertainty about this single point. However, it seems psychologically implausible that there are not some individual differences in higher order cognitive abilities such as reasoning. Ideally, what ought

to be estimated is a *distribution* of individual-participant discriminabilities, with the parameters of this distribution becoming more certain as additional data become available. Bayesian methods naturally achieve this extension to accommodate individual differences using hierarchical models.

Graphical Model

Figure 10 shows a graphical model for a hierarchical version of SDT that allows for individual differences in discriminability and bias across participants and is very similar to that developed by Rouder and Lu (2005). The plate represents repetitions over participants. Within the plate, the graphical model shows the relationships for the i th participant between their discriminability d_i , bias c_i , hit rate h_i , and false alarm rate f_i , and their observed counts of hit k_i^h and false alarm k_i^f decisions.

The discriminability of each participant is assumed to be a value drawn from an overarching Gaussian distribution with mean m_d and precision τ_d . Similarly, the bias of each participant is drawn from a Gaussian with mean m_c and precision τ_c . This means that

$$\begin{aligned} d_i &\sim \text{Gaussian}(m_d, \tau_d) \\ c_i &\sim \text{Gaussian}(m_c, \tau_c). \end{aligned} \quad (17)$$

These overarching Gaussians represent the individual differences in discriminability and bias over participants. Their mean and precision parameters are given standard near noninformative priors:

$$\begin{aligned} m_d &\sim \text{Gaussian}(0, \varepsilon) \\ m_c &\sim \text{Gaussian}(0, \varepsilon) \\ \tau_d &\sim \text{Gamma}(\varepsilon, \varepsilon) \\ \tau_c &\sim \text{Gamma}(\varepsilon, \varepsilon), \end{aligned} \quad (18)$$

where $\varepsilon = .001$ is set near zero.

The discriminability and bias variables for each participant can be reparameterized according to equal-variance SDT into hit and false alarm rates, according to

$$\begin{aligned} h_i &= \Phi\left(\frac{1}{2}d_i - c_i\right) \\ f_i &= \Phi\left(-\frac{1}{2}d_i - c_i\right), \end{aligned} \quad (19)$$

where $\Phi(\cdot)$ is the standard cumulative Gaussian function. Finally, the counts of hit and false alarm decisions follow a binomial distribution with respect to the hit and false alarm rates, and the number of *signal* t^s (i.e., valid or strong) and *noise* t^n (i.e., invalid or weak) arguments presented, so that

$$\begin{aligned} k_i^h &\sim \text{Binomial}(h_i, t^s) \\ k_i^f &\sim \text{Binomial}(f_i, t^n). \end{aligned} \quad (20)$$

Inference

The graphical model for SDT with individual differences was applied to both the induction and the deduction condition data of Heit and Rotello (2005), drawing 100,000

samples after a burn-in period of 1,000 samples. One useful analysis of the full joint posterior distribution, concentrating on the group-level means of discriminability and bias for each condition, is shown in Figure 11. The main panel shows 500 random samples from the joint posterior of the means m_d and m_c , shown as circles for the induction condition and crosses for the deduction condition. The side panels show the marginal distribution for each of these means.

Figure 11 shows that the two conditions have different patterns of mean discriminability and bias. In particular, the induction condition seems to have worse mean discriminability than does the deduction condition. It is also clear that there is a large negative bias for the induction condition, indicating a tendency to overrespond *strong*, whereas the deduction condition shows little if any bias toward overresponding *valid*.

Summary

The conclusion from the Bayesian analyses is that, in complete agreement with Heit and Rotello (2005), it is important to allow discriminability in the induction condition to be different from that in the deduction condition. The contribution of the Bayesian analysis is that this conclusion has been reached, unlike in the Heit and Rotello analysis, allowing for the possibility of individual differences in discriminability and bias across participants and accommodating the clear limitations in how accurately hit and false alarm rates can be estimated from only four observations per participant. In this way, the application demonstrates the ability of Bayesian methods to implement more realistic theoretical and methodological assumptions in drawing inferences from data.

DISCUSSION

The aim of this article was to demonstrate that Bayesian methods can be applied generally and usefully to aid in the understanding and evaluation of psychological models. The three applications tried to span a range of cognitive models and demonstrate a range of Bayesian analyses for addressing interesting theoretical questions informed by the available empirical data. In each case, the idea was to learn something useful through the Bayesian approach that would be difficult to achieve with the ways of relating models to data traditionally used in psychological modeling.

We concede that it is probably the case that the collection of ad hoc methods dominating current practice could be enlarged further with specific methods to achieve the outcomes reported here (after all, that is what *ad hoc* means). But the conceptual insight and technical skills needed to develop new methods stands in stark contrast to the conceptual simplicity and ease of implementation and analysis for the Bayesian graphical modeling approach.

To the extent that the applications succeeded in encouraging the use of Bayesian methods, a number of obvious questions arise. One has to do with the extent to which Bayesian methods can be applied to diverse types of cognitive models and cognitive-modeling approaches. Another involves the scalability of computational forms of Bayesian inference to large-scale models and data sets. A final question involves the extent to which Bayesian inference is being used “just for data analysis,” rather than as a model of human cognition. Some tentative answers will be attempted.

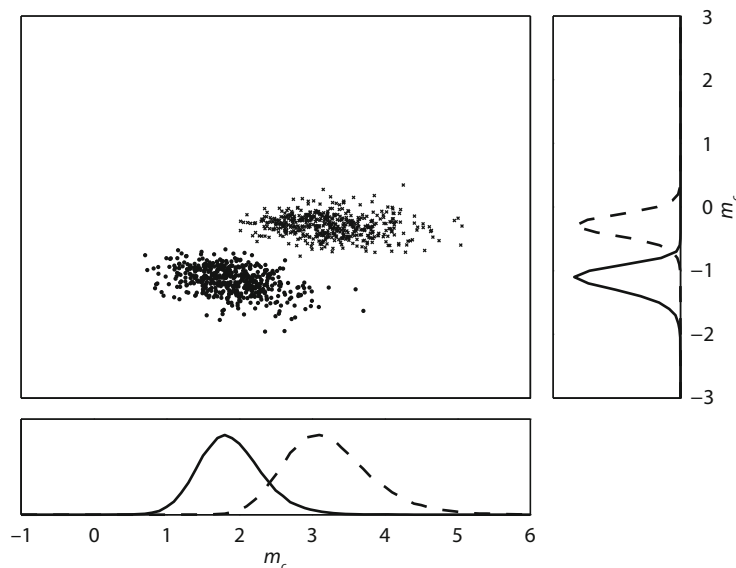


Figure 11. The main panel shows samples from the joint distribution of mean discriminability and mean bias, using circles for the induction condition and crosses for the deduction condition. The side panels show the corresponding marginal distributions, using solid lines for the induction condition and broken lines for the deduction condition.

Generality of Bayesian Methods

One thing that is clear from the three applications presented is that Bayesian analysis can be possible for models not originally developed in Bayesian terms. But a natural question is how generally psychological models can be accommodated in the structured probabilistic framework needed for graphical model interpretation. Deterministic (or qualitative) models that do not have an error theory—algorithmic models of decision making such as take-the-best (Gigerenzer & Goldstein, 1996) constitute one prominent example—clearly need some additional assumptions before being amenable to probabilistic treatment. One promising source for providing principled error theories to such models are *entropification* methods arising from the minimum description length coding approach to model evaluation (Grünwald, 1998, 1999), which have already been applied successfully to a number of psychological models (e.g., Lee, 2006; Lee & Cummins, 2004).

Another class of psychological models that present a challenge are those that do not have a fixed set of parameters. Examples include the original version of the ALCOVE model (Kruschke, 1992) or the SUSTAIN model (Love, Medin, & Gureckis, 2004) of category learning, which specify processes for introducing additional representation nodes within a task, so that their parameter sets change as a function of the data. The mechanics of the applications presented here, with their reliance on WinBUGS and its associated standard Markov chain Monte Carlo methods, do not extend automatically to these types of models. Instead, a more general Bayesian approach is required, using Bayesian nonparametric (also known as *infinite dimensional*) methods (Escobar & West, 1995; Ferguson, 1973; Ghosh & Ramamoorthi, 2003; Neal, 2000). Navarro, Griffiths, Steyvers, and Lee (2006) provide a general introduction to many Bayesian nonparametric ideas and a specific application to modeling individual differences in psychology.

Scalability of Bayesian Methods

Bayesian inference relies upon the full joint posterior distribution over the model parameters as the basis for understanding and evaluating models against data. This is powerful, because the posterior represents everything that is known and unknown about the psychologically interesting variables represented by the parameters. Analytic power, however, comes with a computational burden, and it is reasonable to ask how well the approach scales to large models or data sets. There seem to be grounds for optimism on this front. The topics model (Griffiths & Steyvers, 2002, 2004) of language processing, for example, has been applied to a text corpus with about 2,000,000 words and, hence, has successfully made inferences from data about the joint distribution of about 2,000,000 latent variables. Pioneering hierarchical and generative Bayesian models in vision have also succeeded at impressive scales (e.g., Yuille & Kersten, 2006). The application of Bayesian methods in other fields, such as biology and machine learning, give successful examples at very large scales (e.g., Ridgeway & Madigan, 2003).

More generally, there is a natural tension in psychological modeling between building models and addressing data that have the scale and complexity needed to account for what goes on in the real world, on the one hand, and maintaining the ability to evaluate those models and data in rigorous ways, on the other. The advantage of the Bayesian approach is that it guarantees a principled relationship between models and data. The potential of the Bayesian approach is that it will be able to accommodate progressively larger and more sophisticated models and data.

Bayesian Modeling Versus Data Analysis

It is possible to draw a distinction between two ways that Bayesian ideas can be applied to the modeling of human cognition. One is to assume that the mind solves the inference problems it faces in a Bayesian way. That is, a theoretical assumption is made that the mind does Bayesian inference. Good recent examples of this approach include models of concept and category learning (e.g., Anderson, 1991; Sanborn, Griffiths, & Navarro, 2006; Tenenbaum & Griffiths, 2001) and models of inductive inference and decision making (e.g., Griffiths & Tenenbaum, 2005). These are impressive models and have significantly increased our understanding of the basic abilities of human cognition that they address.

The second way Bayesian ideas can improve our understanding of cognition is to use them to relate model to data, improving the ability to make inferences about parameters and models using the incomplete and uncertain information provided by empirical data. The applications in this article are largely of this type and are part of a more general enterprise that has addressed diverse areas from similarity modeling and structure learning (e.g., Navarro & Lee, 2004), response time distributions (Lee, Fuss, & Navarro, 2007; Rouder, Lu, Speckman, Sun, & Jiang, 2005), and individual differences (e.g., Lee & Webb, 2005; Navarro et al., 2006; Rouder & Lu, 2005).

Although the distinction between *Bayesian models of cognition* and *Bayesian analyses of models of cognition* is an intuitively appealing and practically useful one, it can mask a number of important issues. One issue is that the mere act of analyzing a model from the Bayesian perspective almost always requires making additional theoretical assumptions and, so, changes the model itself to some degree. Most obviously, this happens in specifying prior distributions for parameters, as in all of the applications presented here. Occasionally, existing theory will suggest a form for these priors, but more often the goal will be to specify priors that affect the posteriors following from data as little as possible. In either case, the introduction of priors makes new theoretical assumptions about the psychologically meaningful variables used by a model. In this sense, the adoption of Bayesian methods can never amount to “just data analysis.”

In some cases, applying Bayesian methods can have more dramatic theoretical consequences. One example is the Rescorla–Wagner model of classical conditioning, which, under a non-Bayesian treatment, does not predict backward-blocking effects (Rescorla & Wagner, 1972).

Historically, this failure has been remedied by making additional theoretical assumptions and augmenting the basic model (e.g., Van Hamme & Wasserman, 1994). Dayan and Kakade (2001) have shown, however, that a Bayesian treatment of the basic learning mechanism underlying the Rescorla–Wagner model automatically predicts backward blocking. In this way, the Bayesian analysis lessens the appeal of the more complicated models that have been developed, and so the Bayesian analysis makes a strong contribution to the development of theory.

A final point is that Bayesian inference, by itself, will often not provide all of the ideas needed to model any significant part of human cognition and, so, will often require additional theory to be applied. This means it will be rare to have a purely Bayesian model of some aspect of cognition. As Griffiths et al. (in press) argued:

Bayesian inference stipulates how rational learners should update their beliefs in the light of evidence. The principles behind Bayesian inference can be applied whenever we are making inferences from data, whether the hypotheses involved are discrete or continuous, or have one or more unspecified free parameters. However, developing probabilistic models that can capture the richness and complexity of human cognition requires going beyond these basic ideas.

Two good recent examples are models of feature induction and stimulus similarity (Kemp, Bernstein, & Tenenbaum, 2005; Kemp, Perfors, & Tenenbaum, 2004) and of sequential decision-making behavior (Lee, 2006). Both of these are hierarchical Bayesian models and rely entirely on the Bayesian approach to statistical inference to relate model parameters to data. Both also apply Bayesian methods to model the mind where those ideas are available, using, for example, Bayes's theorem as an account of how information updates mental representations and how model averaging combines different mental hypotheses. But both models need to introduce non-Bayesian components to address the full range of phenomena they aim to explain. Kemp et al. (2005) used generative node-replacement graph grammars and diffusion processes over graphical structure to generate stimulus representations and model their relationships to one another. Lee (2006) relies on a simple finite state account for generating thresholds to guide decision making. None of these theoretical mechanisms could be regarded as following directly and uniquely from Bayesian principles, although all are compatible. Accordingly, both models use Bayesian inference as theoretical accounts of the mind *and* as a method for analyzing data.

Conclusion

The underlying philosophy of Bayesian inference is to represent uncertainty about quantities of interest, using probability distributions, and then to use relevant information to update these representations as it comes to hand, all within the coherent framework for inference offered by probability theory. These capabilities fit naturally with

the goals of modeling in empirical sciences such as psychology. Psychological variables and processes are given formal expression by parameters and models, and the rationale for collecting experimental data is to refine our understanding of those variables and processes.

In this article, the aim has been to give worked examples applying Bayesian methods to models, showing how Bayesian analysis provides the tools to make inferences about hard but important research questions. The use of graphical models and posterior sampling has been emphasized as an easy and powerful method to undertake Bayesian analyses. The adoption of Bayesian methods for analysis promises to improve the way models are related to data, maximizing what we can learn from our ongoing efforts to develop theoretical models and gather empirical information.

AUTHOR NOTE

I thank Geoffrey Iverson, Simon Dennis, Charles Kemp, E.-J. Wagenmakers, Wolf Vanpaemel, Jared Smith, and an anonymous reviewer, and also Evan Heit and Teresa Treat for generously sharing their raw data. Correspondence concerning this article should be addressed to M. D. Lee, Department of Cognitive Sciences, University of California, 3151 Social Sciences Plaza, Irvine, CA 92697-5100 (e-mail: mdlee@uci.edu).

REFERENCES

- ANDERSON, J. R. (1991). The adaptive nature of human categorization. *Psychological Review*, **98**, 409-429.
- ASHBY, F. G., MADDOX, W. T., & LEE, W. W. (1994). On the dangers of averaging across subjects when using multidimensional scaling or the similarity-choice model. *Psychological Science*, **5**, 144-151.
- CHATER, N., & OAKSFORD, M. (2000). The rational analysis of mind and behavior. *Synthese*, **122**, 93-131.
- CHEN, M.-H., SHAO, Q.-M., & IBRAHIM, J. G. (2000). *Monte Carlo methods in Bayesian computation*. New York: Springer.
- DAYAN, P., & KAKADE, S. (2001). Explaining away in weight space. In T. K. Leen, T. G. Dietterich, & V. Tresp (Eds.), *Advances in neural information processing systems 13* (pp. 451-457). Cambridge, MA: MIT Press.
- ESCOBAR, M. D., & WEST, M. (1995). Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association*, **90**, 577-588.
- FERGUSON, T. S. (1973). A Bayesian analysis of some nonparametric problems. *Annals of Statistics*, **1**, 209-230.
- GARNER, W. R. (1974). *The processing of information and structure*. Potomac, MD: Erlbaum.
- GATI, I., & TVERSKY, A. (1982). Representations of qualitative and quantitative dimensions. *Journal of Experimental Psychology: Human Perception & Performance*, **8**, 325-340.
- GELMAN, A. (2006). Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, **1**, 515-534.
- GHOSH, J. K., & RAMAMOORTHY, R. V. (2003). *Bayesian nonparametrics*. New York: Springer.
- GIGERENZER, G., & GOLDSTEIN, D. G. (1996). Reasoning the fast and frugal way: Models of bounded rationality. *Psychological Review*, **103**, 650-669.
- GILKS, W. R., RICHARDSON, S., & SPIEGELHALTER, D. J. (Eds.) (1996). *Markov chain Monte Carlo in practice*. Boca Raton, FL: Chapman & Hall/CRC.
- GREEN, D. M., & SWETS, J. A. (1966). *Signal detection theory and psychophysics*. New York: Wiley.
- GRIFFITHS, T. L., KEMP, C., & TENENBAUM, J. B. (in press). Bayesian models of cognition. In R. Sun (Ed.), *Cambridge handbook of computational cognitive modeling*. Cambridge: Cambridge University Press.
- GRIFFITHS, T. L., & STEYVERS, M. (2002). A probabilistic approach

- to semantic representation. In W. G. Gray & C. D. Schunn (Eds.), *Proceedings of the 24th Annual Conference of the Cognitive Science Society* (pp. 381-386). Mahwah, NJ: Erlbaum.
- GRIFFITHS, T. L., & STEYVERS, M. (2004). Finding scientific topics. *Proceedings of the National Academy of Sciences*, **101**, 5228-5235.
- GRIFFITHS, T. L., & TENENBAUM, J. B. (2005). Structure and strength in causal induction. *Cognitive Psychology*, **51**, 354-384.
- GRÜNWARD, P. D. (1998). *The minimum description length principle and reasoning under uncertainty*. Amsterdam: University of Amsterdam, Institute for Logic, Language, and Computation.
- GRÜNWARD, P. D. (1999). Viewing all models as "probabilistic." In S. Ben-David & P. Long (Eds.), *Proceedings of the 12th Annual Conference on Computational Learning Theory (COLT '99)* (pp. 171-182). Santa Cruz, CA: ACM Press.
- HEIT, E. (2000). Properties of inductive reasoning. *Psychonomic Bulletin & Review*, **7**, 569-592.
- HEIT, E., & ROTELLO, C. (2005). Are there two kinds of reasoning? In B. G. Bara, L. W. Barsalou, & M. Bucciarelli (Eds.), *Proceedings of the 27th Annual Conference of the Cognitive Science Society* (pp. 923-928). Mahwah, NJ: Erlbaum.
- HELM, C. E. (1959). *A multidimensional ratio scaling analysis of color relations*. Princeton, NJ: Princeton University and Educational Testing Service.
- JORDAN, M. I. (2004). Graphical models. *Statistical Science*, **19**, 140-155.
- KASS, R. E., & RAFTERY, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, **90**, 773-795.
- KEMP, C., BERNSTEIN, A., & TENENBAUM, J. B. (2005). A generative theory of similarity. In B. G. Bara, L. W. Barsalou, & M. Bucciarelli (Eds.), *Proceedings of the 27th Annual Conference of the Cognitive Science Society* (pp. 1132-1137). Mahwah, NJ: Erlbaum.
- KEMP, C., PERFORIS, A., & TENENBAUM, J. B. (2004). Learning domain structures. In K. Forbus, D. Gentner, & T. Regier (Eds.), *Proceedings of the 26th Annual Conference of the Cognitive Science Society* (pp. 720-725). Mahwah, NJ: Erlbaum.
- KRUSCHKE, J. K. (1992). ALCOVE: An exemplar-based connectionist model of category learning. *Psychological Review*, **99**, 22-44.
- KRUSCHKE, J. K. (1993). Human category learning: Implications for backpropagation models. *Connection Science*, **5**, 3-36.
- KRUSKAL, J. B. (1964). Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, **29**, 1-27.
- LEE, M. D. (2001). Determining the dimensionality of multidimensional scaling representations for cognitive modeling. *Journal of Mathematical Psychology*, **45**, 149-166.
- LEE, M. D. (2006). A hierarchical Bayesian model of human decision-making on an optimal stopping problem. *Cognitive Science*, **30**, 555-580.
- LEE, M. D., & CUMMINS, T. D. R. (2004). Evidence accumulation in decision making: Unifying the "take the best" and the "rational" models. *Psychonomic Bulletin & Review*, **11**, 343-352.
- LEE, M. D., FUSS, I. G., & NAVARRO, D. J. (2007). A Bayesian approach to diffusion models of decision-making and response time. In B. Schölkopf, J. C. Platt, & T. Hofmann (Eds.), *Advances in neural information processing systems 19* (pp. 809-816). Cambridge, MA: MIT Press.
- LEE, M. D., & POPE, K. J. (2003). Avoiding the dangers of averaging across subjects when using multidimensional scaling. *Journal of Mathematical Psychology*, **47**, 32-46.
- LEE, M. D., & WAGENMAKERS, E.-J. (2005). Bayesian statistical inference in psychology: Comment on Trafimow (2003). *Psychological Review*, **112**, 662-668.
- LEE, M. D., & WEBB, M. R. (2005). Modeling individual differences in cognition. *Psychonomic Bulletin & Review*, **12**, 605-621.
- LOVE, B. C., MEDIN, D. L., & GURECKIS, T. (2004). SUSTAIN: A network model of category learning. *Psychological Review*, **111**, 309-332.
- MACKAY, D. J. C. (2003). *Information theory, inference, and learning algorithms*. Cambridge: Cambridge University Press.
- MACMILLAN, N. A., & CREELMAN, C. D. (2005). *Detection theory: A user's guide* (2nd ed.). Mahwah, NJ: Erlbaum.
- MYUNG, I. J., FORSTER, M., & BROWNE, M. W. (2000). A special issue on model selection. *Journal of Mathematical Psychology*, **44**, 1-2.
- MYUNG, I. J., & PITT, M. A. (1997). Applying Occam's razor in modeling cognition: A Bayesian approach. *Psychonomic Bulletin & Review*, **4**, 79-95.
- NAVARRO, D. J., GRIFFITHS, T. L., STEYVERS, M., & LEE, M. D. (2006). Modeling individual differences using Dirichlet processes. *Journal of Mathematical Psychology*, **50**, 101-122.
- NAVARRO, D. J., & LEE, M. D. (2004). Common and distinctive features in stimulus similarity: A modified version of the contrast model. *Psychonomic Bulletin & Review*, **11**, 961-974.
- NEAL, R. M. (2000). Markov chain sampling methods for Dirichlet processes. *Journal of Computational & Graphical Statistics*, **9**, 619-629.
- NOSOFSKY, R. M. (1984). Choice, similarity, and the context theory of classification. *Journal of Experimental Psychology: Learning, Memory, & Cognition*, **10**, 104-114.
- NOSOFSKY, R. M. (1986). Attention, similarity, and the identification-categorization relationship. *Journal of Experimental Psychology: General*, **115**, 39-57.
- PARSONS, L. M., & OSHERSON, D. (2001). New evidence for distinct right and left brain systems for deductive and probabilistic reasoning. *Cerebral Cortex*, **11**, 954-965.
- PITT, M. A., MYUNG, I. J., & ZHANG, S. (2002). Toward a method of selecting among computational models of cognition. *Psychological Review*, **109**, 472-491.
- RAMSAY, J. O. (1982). Some statistical approaches to multidimensional scaling data. *Journal of the Royal Statistical Society A*, **145**, 285-312.
- RESCORLA, R. A., & WAGNER, A. R. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.), *Classical conditioning II: Current research and theory* (pp. 64-99). New York: Appleton-Century-Crofts.
- RIDGEWAY, G., & MADIGAN, D. (2003). A sequential Monte Carlo method for Bayesian analysis of massive datasets. *Data Mining & Knowledge Discovery*, **7**, 301-319.
- RIPS, L. J. (2001). Two kinds of reasoning. *Psychological Science*, **12**, 129-134.
- ROUDER, J. N., & LU, J. (2005). An introduction to Bayesian hierarchical models with an application in the theory of signal detection. *Psychonomic Bulletin & Review*, **12**, 573-604.
- ROUDER, J. N., LU, J., SPECKMAN, P., SUN, D., & JIANG, Y. (2005). A hierarchical model for estimating response time distributions. *Psychonomic Bulletin & Review*, **12**, 195-223.
- SANBORN, A. N., GRIFFITHS, T. L., & NAVARRO, D. J. (2006). A more rational model of categorization. In R. Sun & N. Miyake (Eds.), *Proceedings of the 28th Annual Conference of the Cognitive Science Society* (pp. 726-731). Mahwah, NJ: Erlbaum.
- SHEPARD, R. N. (1957). Stimulus and response generalization: A stochastic model relating generalization to distance in psychological space. *Psychometrika*, **22**, 325-345.
- SHEPARD, R. N. (1962). The analysis of proximities: Multidimensional scaling with an unknown distance function. I. *Psychometrika*, **27**, 125-140.
- SHEPARD, R. N. (1980). Multidimensional scaling, tree-fitting, and clustering. *Science*, **210**, 390-398.
- SHEPARD, R. N. (1987). Toward a universal law of generalization for psychological science. *Science*, **237**, 1317-1323.
- SHEPARD, R. N. (1991). Integrality versus separability of stimulus dimensions: From an early convergence of evidence to a proposed theoretical basis. In J. R. Pomerantz & G. L. Lockhead (Eds.), *The perception of structure: Essays in honor of Wendell R. Garner* (pp. 53-71). Washington, DC: American Psychological Association.
- SHEPARD, R. N. (1994). Perceptual-cognitive universals as reflections of the world. *Psychonomic Bulletin & Review*, **1**, 2-28.
- SLOMAN, S. A. (1998). Categorical inference is not a tree: The myth of inheritance hierarchies. *Cognitive Psychology*, **35**, 1-33.
- SPIEGELHALTER, D. J., THOMAS, A., & BEST, N. G. (2004). *WinBUGS Version 1.4 user manual*. Cambridge: Medical Research Council Biostatistics Unit.
- SPIEGELHALTER, D. J., THOMAS, A., BEST, N. G., & GILKS, W. R. (1996). *BUGS Examples Volume 1, Version 0.5*. Cambridge: Medical Research Council Biostatistics Unit.
- TENENBAUM, J. B., & GRIFFITHS, T. L. (2001). Generalization, similarity, and Bayesian inference. *Behavioral & Brain Sciences*, **24**, 629-640.

- TREAT, T. A., MACKAY, D. B., & NOSOFKY, R. M. (1999, July). *Probabilistic scaling: Basic research and clinical applications*. Paper presented at the 32nd Annual Meeting of the Society for Mathematical Psychology, Santa Cruz, CA.
- VAN HAMME, L. J., & WASSERMAN, E. A. (1994). Cue competition in causality judgments: The role of nonrepresentation of compound stimulus elements. *Learning & Motivation*, **25**, 127-151.
- YUILLE, A. L., & KERSTEN, D. (2006). Vision as Bayesian inference: Analysis by synthesis? *Trends in Cognitive Sciences*, **10**, 301-308.

NOTES

1. We are aware of the argument that it is not clear that the exponential decay relationship applies as well to direct judgment of (dis)similarity as it does to the conditional probabilities obtained from identification confusion or generalization experiments. In particular, nonmetric analyses of direct judgment data often show a relationship that is more nearly linear.

2. These interpretations of $0 < r \leq 2$ also map naturally onto Shepard's (1987) framework, corresponding to the full possible range of -1 to $+1$ correlations for the *consequential regions* that underpin his theoretical results. In contrast, other than for the supremum metric, it is difficult to give psychological meaning to metrics with $r > 2$, and so our attention will be restricted to $0 < r \leq 2$.

3. Gelman (2006) points out that this prior distribution, although very widely used, is problematic for small variances. For all of the applications reported here, the inferred variances are large enough to avoid this difficulty.

4. In an alternative analysis, r was given a prior that allowed values greater than 2.0. The posterior for r for the color data under these assumptions retained a mode below 2.0 but had some significant density extending to 2.0 and beyond.

5. In this case, it would have been highly desirable to have observed mixing within the chains, rather than relying on random initial assignments for a large number of chains. In other words, what posterior sampling should ideally be able to produce is a single chain that includes all of the patterns of latent assignment in proportion to their posterior density. What has been done is to combine many different chains to approximate this output. Nevertheless, the patterns of latent assignments observed are intuitively sensible, and we believe that they can be used to make inferences. But it is obviously necessary to treat measures such as the posterior densities of each latent assignment that could be derived from this analysis with caution.

(Manuscript received November 14, 2006;
revision accepted for publication May 17, 2007.)