

## ARTICLE

# Applying Occam's razor in modeling cognition: A Bayesian approach

IN JAE MYUNG and MARK A. PITT  
Ohio State University, Columbus, Ohio

In mathematical modeling of cognition, it is important to have well-justified criteria for choosing among differing explanations (i.e., models) of observed data. This paper introduces a Bayesian model selection approach that formalizes Occam's razor, choosing the simplest model that describes the data well. The choice of a model is carried out by taking into account not only the traditional model selection criteria (i.e., a model's fit to the data and the number of parameters) but also the extension of the parameter space, and, most importantly, the functional form of the model (i.e., the way in which the parameters are combined in the model's equation). An advantage of the approach is that it can be applied to the comparison of non-nested models as well as nested ones. Application examples are presented and implications of the results for evaluating models of cognition are discussed.

A goal of research in psychology, as in other behavioral sciences, is to infer the underlying process that generated observed data. The use of sophisticated mathematical models to describe these processes has grown considerably, especially in cognitive psychology (see, e.g., J. R. Anderson & Sheu, 1995; N. H. Anderson, 1981; Ashby & Townsend, 1986; Busemeyer & Townsend, 1993; Gillund & Shiffrin, 1984; Green & Swets, 1966; Hintzman, 1986; Kruschke, 1992; Massaro & Friedman, 1990; Medin & Schaffer, 1978; Murdock, 1982; Nosofsky, 1986; Oden & Massaro, 1978; Reed, 1972; van Zandt & Ratcliff, 1995). Yet, the development of equally sophisticated and well-justified methods for evaluating the adequacy of the models themselves has lagged behind. Jacobs and Grainger (1994) recently summarized a number of criteria for choosing among models: (1) generality (does the model generalize well across different experimental settings?); (2) explanatory adequacy (are the assumptions of the model plausible and compatible with established findings?); (3) descriptive adequacy (does the model fit the pattern of data well?); and (4) complexity (is the formulation

[e.g., the number of parameters] of the model simple?). Among these, descriptive adequacy and complexity have been used most frequently, probably because they are easier to quantify than the other two. Together they embody the principle of Occam's razor, which states that "entities should not be multiplied beyond necessity" (William of Occam, ca. 1290-1349). The goal of model selection is to choose the *simplest* (i.e., least complex) model that describes the data *well* (i.e., descriptive adequacy).

In this paper, we introduce a new Bayesian method of formalizing Occam's razor in model selection. It goes beyond current selection methods by taking into account dimensions of complexity that are not captured by its predecessors. We begin with a tutorial on model selection methods that are currently in use. Next, fundamentals of the Bayesian model selection approach are described, and its desirable properties are discussed. The utility of the Bayesian approach is then demonstrated using concrete examples with simulated data. Finally, the merits and shortcomings of the Bayesian approach are discussed and contrasted with traditional approaches.

## MODEL SELECTION CRITERIA

### Descriptive Adequacy

The goal of mathematical modeling in cognitive psychology is straightforward: Given observed data, identify the underlying processes that generated the data. Because a model is defined as a set of assumptions about underlying processes, the goal of the researcher is to determine the viability of the model. There are, however, at least two obstacles to such an endeavor. First, given the nearly infinite number of distinct models that can be defined by combining different assumptions, the true model might

A portion of this work was presented at the 27th annual meeting of the Society for Mathematical Psychology held at the University of California, Irvine, in August 1995. Many people provided very useful feedback on earlier versions of this paper. They include Greg Ashby, Michael Browne, Jerry Busemeyer, Dan Friedman, Lester Krueger, Duncan Luce, Robert MacCallum, Dominic Massaro, Richard Schweickert, James Townsend, Michael Wenger, and Patricia van Zandt. Greg Ashby and Lester Krueger were especially helpful in sharpening our thinking on model complexity. This research was supported in part by Ohio Supercomputer Center Grant PAS887-1. Correspondence should be addressed to I. J. Myung, Department of Psychology, Ohio State University, 1885 Neil Avenue Mall, Columbus, OH 43210-1222 (e-mail: myung.1@osu.edu).

not be one of a particular set of models that is being tested. Second, random noise in data can obscure model identification. Consequently, a realistic goal of mathematical modeling is to choose the model that represents the closest approximation to the "true" model.

How can the closest approximation be identified when the true model has yet to be discovered? Consider the situation in which the true model is included in the set of models being tested and further, data are noise free. In this ideal situation, the true model must fit the data perfectly (e.g., as measured using a metric such as sum of squared errors). Note that this is a *necessary*—but not sufficient—condition, for there could be more than one model that fits the data perfectly. By extending this logic to less ideal situations (e.g., noisy data), the following model selection rule is obtained: Choose the model that provides the best fit to the data. Accordingly, the foremost criterion of model selection, *descriptive adequacy*, is born. Examples of descriptive adequacy measures that are in use include the percent variance accounted for by the model (i.e., coefficient of determination), the sum of squared errors (SSE) between observed and predicted outcomes, and the maximum likelihood, in which the probability of obtaining the observed data is maximized with respect to the model's possible parameter values (see Bickel & Doksum, 1977).

For a model to be considered true, it must satisfy the minimal condition of sufficiency in fitting data well. A failure to do so invalidates the model. An illustrative example is shown in Figure 1. Each of the four solid lines in the figure represents a model's best fit to the same data set (solid dots) using the least squares estimation method. Model 1 is a two-parameter linear model. As can be seen, it fits the data poorly, with only 79.5% of the variance accounted for. Systematic deviations from the line are evident at the endpoints and in the middle of the range.

Clearly, Model 1 fails the test of descriptive adequacy and can be dropped from further consideration. In contrast, Model 2, a three-parameter exponential model, fits the data fairly well, accounting for 96% of the variance with no systematic deviation from the fit. Between the first two models, Model 2 would be chosen as the preferred description of the data. Models 3 and 4 are discussed in the next section.

### Model Complexity

It is important to note that descriptive adequacy is a heuristic. Selection of the best fitting model may be useful in identifying the true model or closest approximation, but the rule's accuracy is not guaranteed. This is because model fit can be improved by increasing model complexity. *Complexity* refers to the flexibility inherent in a model that enables it to fit diverse patterns of data.<sup>1</sup> It can be understood by contrasting the data-fitting capabilities of simple and complex models. A simple model is one that assumes that a specific pattern will be found in the data. If this pattern occurs, the model will fit the data well. Simple models make clear and falsifiable predictions precisely because a specific pattern is assumed to be present. In terms of actually testing the model, what this means is that the model's fit will be good over a sizable range of parameter values. A complex model, on the other hand, is more flexible than a simple model, providing good fits to a wide range of data patterns. To do so, however, the complex model's parameters must be finely tuned. This is because as the model's parameters change, even slightly, the postulated data pattern also changes.

There are at least three dimensions of a model that contribute to its complexity, thereby significantly affecting model fit: the number of parameters, the model's functional form, and the extension of the parameter space. In

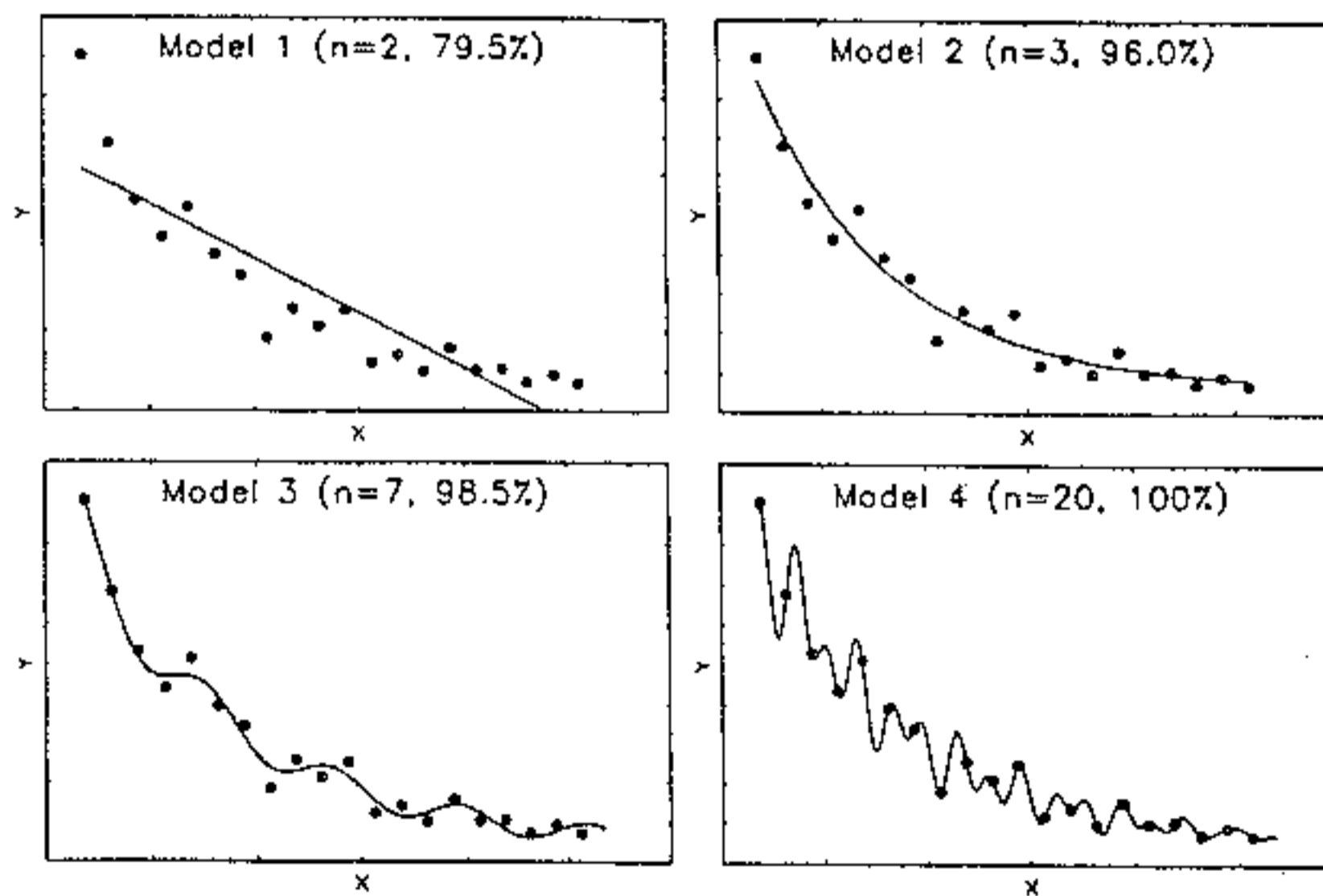


Figure 1. The effect of the number of parameters in a model on the model's ability to fit data. A single data set (dots) was fitted to four models (lines) differing in the number of parameters ( $n$ ). Percentage of variance accounted for by each model is shown in parentheses.

the following subsections, we discuss the implications of each of these for model selection.

**Number of parameters.** In general, a model with many parameters fits data better than a model with few parameters, even if the latter generated the data (Collyer, 1985). The effect of excessive parameters on model fit is illustrated in the bottom panels in Figure 1. Model 3 was created from Model 2 by introducing an additional cyclic component with four new parameters. Its fit is improved over Model 2, but only by a meager 2.5%. The extra parameters seem unnecessary, capturing what appear to be a few idiosyncracies in the data. Model 4 is a more dramatic example of a model with excessive parameters. It fits the data perfectly and even tells us more than the data do. Model 4 is an interpolation model with 20 parameters, the same number as data points. Such a model can always be found and by design, it will fit the data perfectly. A problem with Models 3 and 4 is that they generalize poorly to other data because they precisely fit only one data set. Thus it seems unlikely that these models accurately reflect the mental processes responsible for generating the data.

The preceding examples should make it clear that there is a tradeoff between fit and generalizability as the number of parameters increases. The best fit should be preferred when it is not achieved at the expense of additional parameters. (For a discussion of the related issue of the theoretical justification of the extra parameters, see Jacobs & Grainger, 1994, and Cudeck & Henly, 1991.) Three model selection methods were proposed that adjust for variation in the number of parameters among models. They do so by penalizing more heavily models with many parameters as opposed to those with few parameters.

Akaike (1973, 1983) introduced the Akaike information criterion (AIC), which is defined as

$$AIC_i = -2\ln(ML_i) + 2n_i \quad (1)$$

In this equation,  $ML_i$  is the maximum likelihood for Model  $i$  and  $n_i$  is the number of free parameters in the model. The criterion prescribes that the model that minimizes the AIC should be chosen. If the choice is between two models with equal maximum likelihood values but different numbers of parameters, AIC favors the model with fewer parameters. The AIC has been employed in time series analysis (e.g., Cryer, 1986), psychometric data analysis (e.g., Bozdogan, 1987), and mathematical modeling of categorization (e.g., Maddox & Ashby, 1993; Takane & Shibayama, 1992).

Schwarz (1978) introduced another criterion called the Bayesian information criterion (BIC), defined as

$$BIC_i = -2\ln(ML_i) + 2n_i \ln s \quad (2)$$

where  $\ln(s)$  is the natural log of sample size ( $s$ : number of observations per stimulus) of the data. The model that minimizes BIC should be chosen. Note its similarity to AIC. Only the penalty term for excessive parameters is slightly different ( $2n_i$  vs.  $n_i \ln(s)$ ). A comparison of both criteria shows that AIC favors more complex models than BIC for large sample sizes when  $\ln(s) > 2$  (i.e.,  $s > 8$ ).

Steiger and colleagues (Steiger, 1990; Steiger & Lind, 1980; see also Browne & Cudeck, 1992, Cudeck & Henly, 1991) introduced yet another criterion, called the root mean square error of approximation (RMSEA), defined as

$$RMSEA = \sqrt{\frac{F_i}{N - n_i}} \quad (3)$$

In this equation, the function  $F_i$  is a measure of the lack of fit for model  $i$  (e.g., using SSE) and  $N$  is the data size.<sup>2</sup> RMSEA penalizes complex models by subtracting the number of parameters ( $n_i$ ) from the divisor,  $N$ . The RMSEA closely approximates the root mean squared deviation (RMSD), which is a sample statistic often used in mathematical modeling (see, e.g., Massaro & Friedman, 1990).

**Functional form.** The number of parameters is the only dimension of model complexity that AIC, BIC, and RMSEA consider. An often unrecognized dimension of model complexity that can significantly affect model fit by simply capturing irrelevant patterns of data is *functional form*, which can be defined as the way in which parameters are combined in the model equation.

Figure 2 illustrates the effect that functional form can have on model fit. There is a universe ( $U$ ) of possible data, subsets of which are consistent with different models. The set  $M_a$  denotes a data region that Model  $M_a$  fits "well" in an appropriately defined sense (e.g., maximum likelihood). Model  $M_b$ , which has more parameters than  $M_a$  but the same functional form, accounts for a wider range of data. Now consider the third model,  $M_c$ , which has the same number of parameters as Model  $M_a$  but assumes a functional form different from that of  $M_a$ . Not only is the  $M_c$  region larger than that of  $M_a$  (i.e., more data are accounted for by  $M_c$ ) but also, most of the data that  $M_a$  fits can also be fitted by  $M_c$ . Importantly, the reverse is not true. Thus the functional form of Model  $M_c$  makes it a more flexible model than  $M_a$ .

Another example of the importance of functional form in model behavior is a comparison of two psychophysical models. Townsend (1975) pointed out that Stevens's law ( $\Psi(x) = k \cdot x^a$ ) is more complex, and thus less falsifiable, than Fechner's law ( $\Psi(x) = k \cdot \ln(x + \beta)$ ), even though both have the same number of parameters. This is

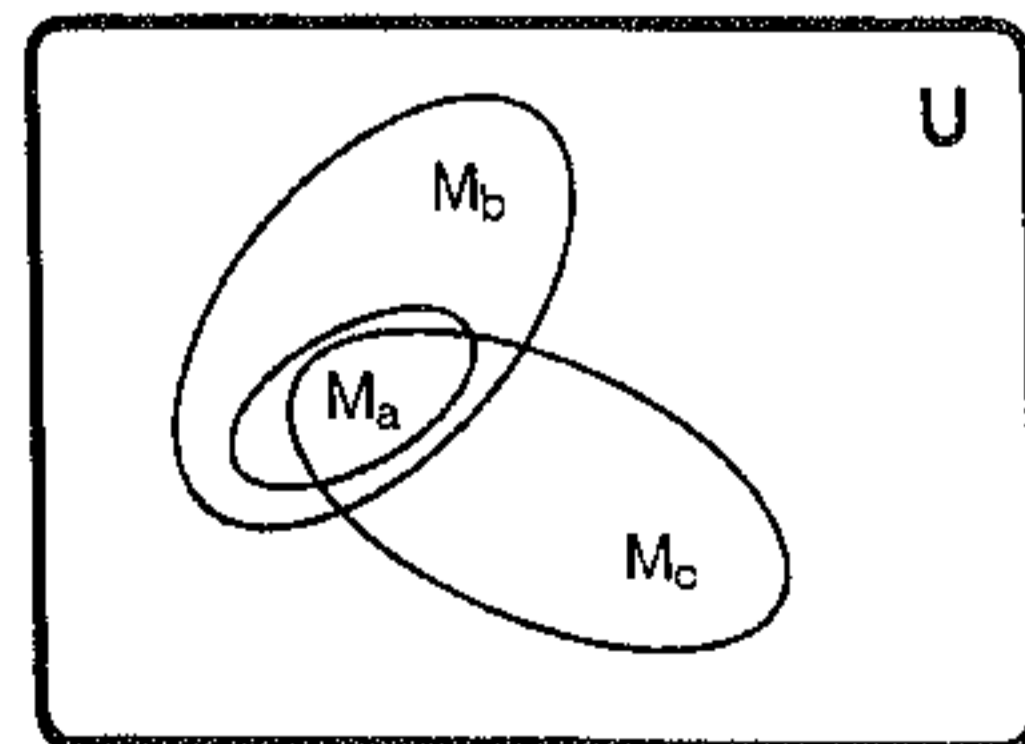


Figure 2. Data-fitting capabilities of various models.

because psychological and physical dimensions are assumed to be related by a power function in Stevens's law, making it capable of fitting data that have negative, positive, and zero curvature. Fechner's law assumes a logarithmic relationship, which can fit data patterns with a negative curvature only.

More recently, the superfluous effects of functional form in model fitting were shown in an insightful simulation study by Cutting, Bruno, Brady, and Moore (1992). They compared two *non-nested* models of perception, the linear integration model (LIM; Anderson, 1981) and the fuzzy logic model of perception (FLMP; Oden & Massaro, 1978).<sup>3</sup> The two models have the same number of parameters but different functional forms, a linear additive function in LIM and a nonlinear multiplicative function in FLMP. In the first simulation, data patterns were generated that spanned the range of the parameter space. Both models were then fitted to all data patterns. Results showed that FLMP was more flexible than LIM, providing a superior fit (i.e., smaller RMSD) for 80.3% of the data patterns with exponential functions and 95.7% with logistic functions. If the number of parameters were the only factor affecting model complexity, the percentage of superior fits should have been about 50%.

A more convincing demonstration of FLMP's superior functional form can be found in the results from another simulation, in which the models were fitted to a set of random numbers (dependent variable) generated from a uniform distribution ranging from 0 to 1. Again, FLMP performed at better than chance level (60.8%). Further exploring this issue, Cutting et al. (1992) also fitted the models to noisy data whose means were generated by LIM. Again, FLMP fitted the data better than LIM itself about 60% of the time, and did so even more decisively when the added noise was relatively large. Given that LIM produced the data, the percentage of superior fits by FLMP should have been at least below the chance level of 50%, though ideally it should have been zero.

Massaro and Cohen (1993) argued that Cutting et al.'s (1992) simulations distorted the true data-fitting abilities of FLMP. Massaro and Cohen stated that in the simulations in which the FLMP and LIM were fitted to a range of exponential and logistic functions, the functions were ones that favored FLMP, and that a set of plausible functions could be created that could just as easily cause the LIM to provide a superior fit. This criticism would have been more convincing if supporting evidence had been provided. Massaro and Cohen also claimed that the simulations in which the models were fitted to random data provided an unfair test of the models, because both models fit the data so poorly (very large RMSDs) that FLMP's superior fits are not meaningful. This reasoning ignores the consistency with which FLMP provided better fits than LIM and places an arbitrary limit on the interpretation of goodness-of-fit measures. Regardless of the quality of the fit, the simulations clearly show FLMP's superior ability at fitting data. Finally, Massaro and Cohen argued that the simulations using noisy data were invalid because Cutting et al. had used an artificially narrow range of data

values (between 0.3 and 0.7, whereas observed values often span from 0 to 1). Our own simulations in the present paper go a long way toward dispelling this criticism. Data values ranged between 0.1 and 0.9, and we consistently found that FLMP provided superior fits. Furthermore, we also found that LIM *never* provided superior fits to those provided by FLMP. If the two models could mimic each other equally well (i.e., if they were equally flexible), then LIM should have fit FLMP data as well as FLMP fit LIM data. Without assuming that FLMP is more flexible than LIM, it is difficult to explain the simulation results, in particular the asymmetric pattern of model mimicry.

The point of this discussion is not to argue against the psychological validity of FLMP or LIM, but rather to make clear the importance of functional form in model selection. Because both models have the same number of parameters, FLMP's advantage must be due to the functional form of the model equation. The nonlinear function of FLMP seems more flexible than the linear function of LIM in fitting noisy data, thus improving model fit.

**Extension of parameter space.** Another dimension of model complexity that can also influence model fit is the extension of the parameter space. To illustrate the point, consider two one-parameter models that share the same functional form of  $y = 1/(1 + e^{-\theta x})$ , but assume different ranges of the parameter  $\theta$ . In Model 1,  $\theta$  ranges from  $-R$  to  $R$ , where  $R$  is a constant. In Model 2,  $\theta$ 's range is cut in half, spanning from 0 to  $R$ . Model 1 allows the parameter to be either positive or negative, whereas Model 2 allows it to be only positive. The parameter range of Model 1 is twice that of Model 2. As shown in Figure 3, this difference in parameter space means that Model 1 can fit decreasing ( $\theta < 0$ ) and increasing ( $\theta > 0$ ) patterns of data. Model 2 can fit only the latter. Solely by virtue of the larger parameter space, Model 1 is guaranteed to fit data showing a decreasing pattern better than Model 2.

To summarize, mathematical models can vary along a number of dimensions. It is necessary to incorporate all of them into the model selection procedure to maximize its accuracy. Selection of a model based solely on descriptive adequacy can lead to incorrect conclusions unless appropriate adjustments are made for the three dimensions of model complexity (number of parameters, functional form, and extension of parameter space), which can significantly and independently affect model fit. Standard model selection criteria such as AIC and BIC consider only the first of these.<sup>4</sup> Such criteria are appropriate to use when comparing nested models, which tends to be done in psychometric research (e.g., analysis of variance models, regression models, latent variable models). However, their use in comparing non-nested models, which is virtually always the case in models of cognition, is not justified and can lead to erroneous conclusions, as we demonstrate in the section Application Example of Bayesian Model Selection.

The Bayesian method described next provides a way to overcome the limitation of the standard model selection

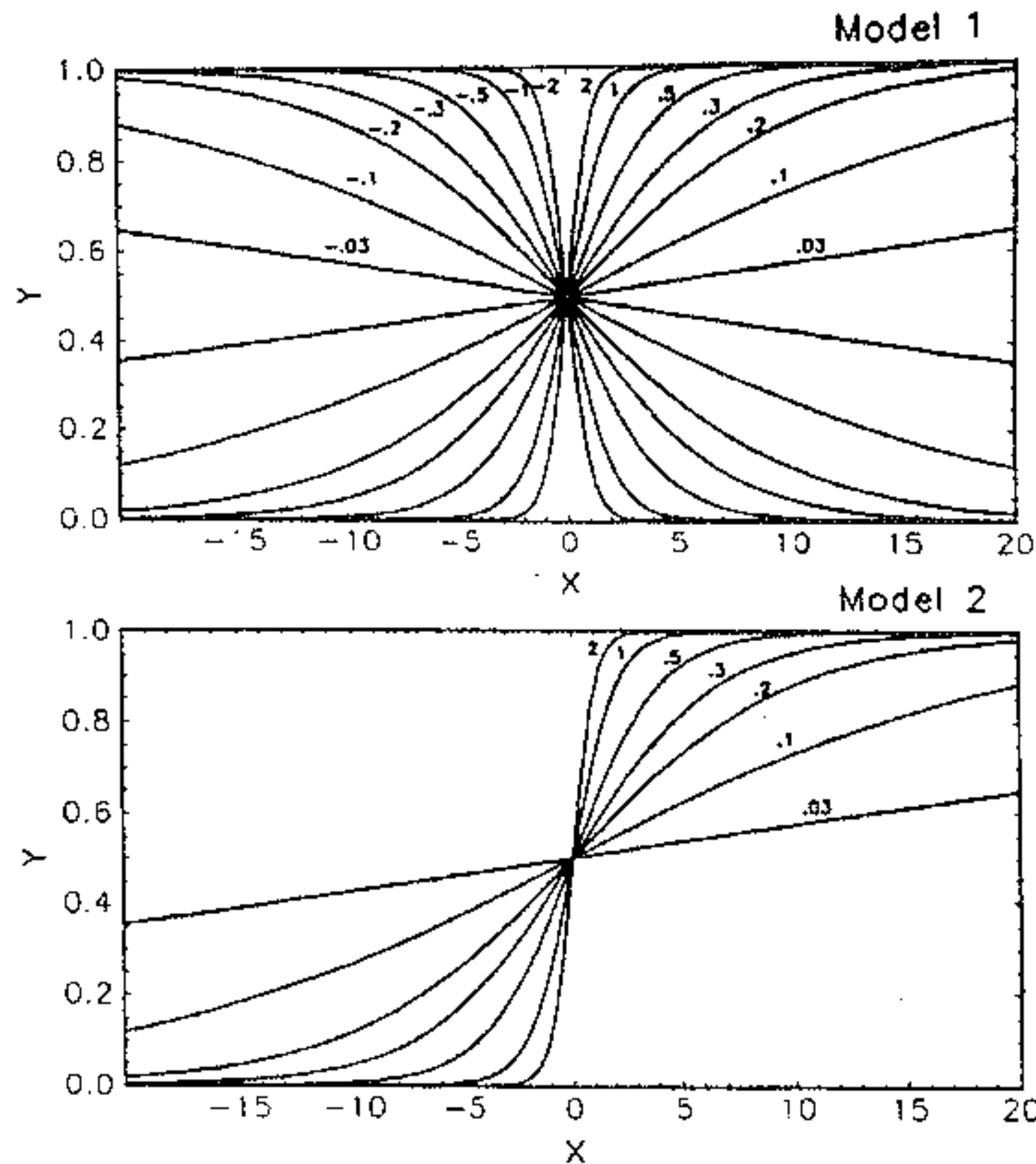


Figure 3. An example of how the extension of the parameter space can affect model fit. A range of possible data patterns that each of two different models can fit is shown. The two models share the same functional form of  $y = 1/(1 + e^{-\theta x})$  but assume different ranges of the parameter  $\theta$ :  $-R < \theta < R < \infty$  for Model 1;  $0 < \theta < R < \infty$  for Model 2. Values of  $R$  are listed next to the corresponding function.

criteria. The method combines *all three* dimensions of model complexity, as well as descriptive adequacy, into a single measure, thus making it particularly suitable for testing models of cognition.<sup>5</sup>

## BAYESIAN MODEL SELECTION

We begin this section by describing the basics of the Bayesian model selection approach, followed by a consideration of the computational issues surrounding its implementation.

### Fundamentals of Bayesian Model Selection

Recently there has been a surge of interest in Bayesian methods of model selection among statisticians, mathematicians, physicists, chemists, and neural network researchers (e.g., Berger & Perrichi, 1996; Bretthorst, 1989; Carlin & Chib, 1995; Gelfand & Dey, 1994; Gregory & Lored, 1992; W. H. Jeffreys & Berger, 1992; Le & Raftery, 1996; MacKay, 1992; Raftery, 1993, 1994; Rissanen, 1986, 1990; Smith & Roberts, 1993). In this section, we provide a comprehensive overview of Bayesian model selection methods, emphasizing their utility in eval-

uating models of cognition and discussing their implications for mathematical modeling of psychological data (see Kass & Raftery, 1995, for a more technically rigorous presentation of the material).

Application of the Bayesian approach requires a statistical formulation of a model that specifies the probability density function (i.e., distribution) of the data. For example, a model of semantic priming may assume that response times in the lexical decision task follow a normal probability density with a certain mean and standard deviation. Moreover, the mean may increase or decrease depending on theoretical predictions, whereas the standard deviation is an unknown but fixed quantity.

Model selection is carried out by computing the posterior probability that each of the models is correct (i.e., true) given a particular data set. This process is formally stated as follows:

For each model  $M_i$  and the data  $D$ , compute the posterior probability,  $P(M_i|D)$ , from the prior probability,  $P(M_i)$ , and the evidence for  $M_i$ ,  $P(D|M_i)$ , using the Bayes rule, and then choose the model that maximizes the posterior probability.

For simplicity, we will consider only the two-model case of models  $M_1$  and  $M_2$ , although extension to multimodel cases is straightforward. From the Bayes rule, we obtain

$$\frac{P(M_1|D)}{P(M_2|D)} = \frac{P(M_1)}{P(M_2)} \cdot \frac{P(D|M_1)}{P(D|M_2)} \quad (4)$$

The first term,  $P(M_1|D)/P(M_2|D)$ , is called the posterior odds ratio, and represents the ratio of the model likelihoods given the data. The second term,  $P(M_1)/P(M_2)$ , is called the prior odds ratio; it represents the ratio of the model likelihoods before evaluating the data. Model priors are determined independently of the data. The third term,  $P(D|M_1)/P(D|M_2)$  is called the Bayes factor or evidence ratio. The numerator and denominator each represent the probabilities of the data given each model. Model selection is achieved by computing the prior odds ratio and the Bayes factor. In other words, the equation above can be rewritten as

posterior odds ratio = (prior odds ratio)  $\times$  (Bayes factor).

**Prior odds  $P(M_i)$ .**  $P(M_i)$  represents the probability that the model  $M_i$  is a true description of the events under study *before* data are collected. For real-world problems, this probability may not exist, or not be known even if it does. We assume equal model priors; that is,  $P(M_1)/P(M_2) = 1$ . This assumption is often made by researchers in Bayesian model selection (Berger & Perrichi, 1996; Carlin & Chib, 1995; Gregory & Loredó, 1992; Raftery, 1993), and is an unavoidable simplification of a problem for which there is not as yet a satisfactory solution. An implication of the assumption is that model selection is based solely on the Bayes factor. The decision to ignore model priors does not undermine use of the Bayesian approach. Rather, it represents a calculated approximation to the true Bayesian model selection criterion.<sup>6</sup>

**Bayes factor  $P(D|M_1)/P(D|M_2)$ .** The Bayes factor is defined as a ratio of two *marginal likelihoods*, which can best be understood by noting their relationship to likelihood functions. A likelihood function is the probability of the given data computed for a particular value of the parameter of a model,  $P(D|\theta, M_i)$ , and thus is a function of the parameter  $\theta$ . In essence, the likelihood function is a goodness-of-fit measure of a model for a given parameter value. If the value of the parameter changes, the fit might also change. The marginal likelihood of a model,  $P(D|M_i)$ , is the probability of the data as a whole, independent of parameter values. It is an average of likelihoods under a prior distribution of the parameter. The marginal likelihood is expressed in the following integral form:

$$P(D|M_i) = \int P(D|\theta, M_i)P(\theta|M_i)d\theta \quad (i = 1, 2). \quad (5)$$

where  $\theta$  is a parameter vector under Model  $i$ ,  $P(D|\theta, M_i)$  is the likelihood function, and  $P(\theta|M_i)$  is the prior density of  $\theta$  for Model  $i$ .

Under the particular assumption of a uniform prior density function of the parameter [i.e.,  $P(\theta|M_i) = \text{constant}$ ], the marginal likelihood is proportional to the area

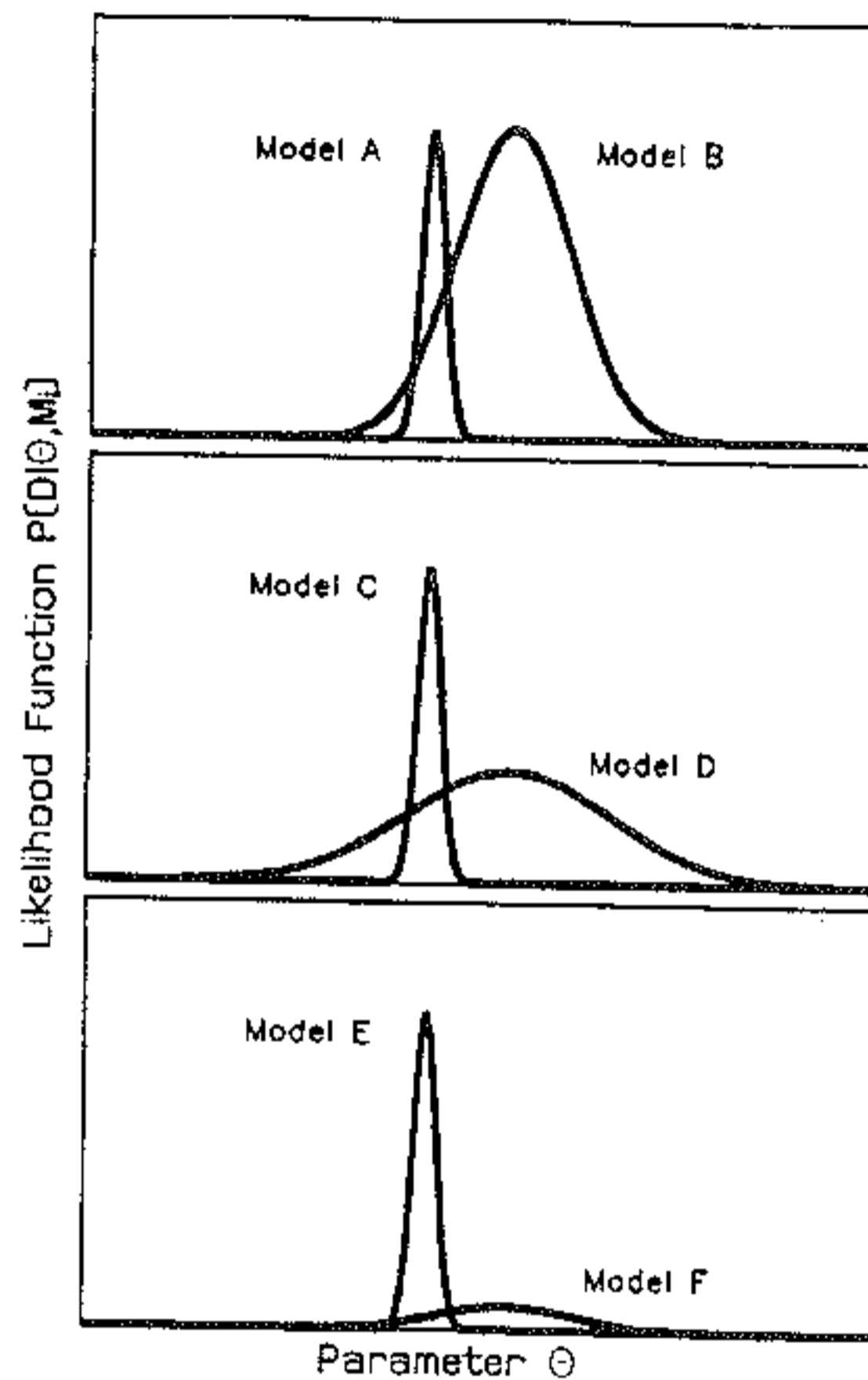


Figure 4. Maximum likelihood functions for three pairs of models as a function of the parameter  $\theta$ . Models A, C, and E are more complex (more peaked functions) than are models B, D, and F. See text for details.

under the curve representing the likelihood function. To illustrate, likelihood functions for two models, A and B, are shown in the top panel of Figure 4. Note that the abscissa of the figure represents parameter value, not data value. Although the maximum likelihoods (i.e., the highest values) are the same for both models, the marginal likelihood of Model B is larger because the area under its function is larger. Note that the functions should be thought of as actual representations of the models. Their peakedness is determined by model complexity, to which we now turn.

**Model Complexity: A Bayesian View**

A model is selected in the Bayesian approach by maximizing the marginal likelihood. Under some simplifying assumptions, the marginal likelihood can be expressed in words as

$$\text{marginal likelihood} = \frac{\text{goodness of fit}}{\text{model complexity}} \quad (6)$$

The numerator, goodness of fit, is the maximum likelihood of the data given the model,  $P(D|\theta_o, M_i)$ , where  $\theta_o$  is the parameter value that maximizes the likelihood

function. The denominator term represents a model complexity measure that embodies all three dimensions of complexity: number of parameters, functional form, and extension of parameter space. A simplified mathematical function depicting their relationship is shown below. (A technical discussion of the mathematics is presented in Appendix A.)

$$\text{model complexity} = g(n_i, F_i(\theta_o), R_i) \quad (7a)$$

$$g(n_i, F_i(\theta_o), R_i) = \beta F_i(\theta_o) (R_{i1} \cdots R_{in_i}) \quad (7b)$$

In the equation above,  $\beta$  is a positive scaling factor,  $n_i$  is the number of parameters of Model  $i$ , and  $R_i = (R_{i1}, \dots, R_{in_i})$  is a vector representing the range of the parameter vector  $\theta$  of Model  $i$  (e.g.,  $R_{ik}$  equals 2 if the parameter  $\theta_k$  is defined on  $-1 \leq \theta_k \leq +1$ ).  $F_i(\theta_o)$  is a "functional-form" factor whose value depends on how the parameters are combined mathematically in the model. Because this factor is evaluated at  $\theta_o$ , which is determined by the data, the value of the functional-form factor may depend on the data.<sup>7</sup> Roughly speaking, the peakedness of the likelihood function near  $\theta_o$  is positively correlated with the value of  $F_i(\theta_o)$ . For example, in Figure 3,  $F_i(\theta_o)$  is higher for Model A than for Model B. As  $F_i(\theta_o)$ ,  $R_{ik}$ , and  $n_i$  increase, so does model complexity. According to Equation 6, a more complex model yields a smaller marginal likelihood; consequently, it is less likely to be selected as the best fitting model. Model B would be chosen over Model A using the Bayesian method if other things were equal (i.e., same  $R_{ik}$  and  $n_i$  for both models).

The middle and lower panels of Figure 4 are examples of other model relationships that illustrate further how model complexity affects model selection. For the pairs of models, assume that they have the same number of parameters and the same extension of the parameter space, and further, that the prior density of the parameter follows a uniform distribution. In the lower panel, both selection methods (Bayesian and standard) would likely reach the same conclusion when evaluating Models E and F. Not only does Model F provide a much poorer fit than Model E, but also its marginal likelihood is smaller. Whether the maximum likelihood or the marginal likelihood were used to select the most appropriate model, Model E would be the clear winner.

In the middle panel, Model C would be chosen over Model D using the standard method because C's maximum likelihood is larger than D's. Bayesian model selection is less straightforward in this case. Even though D's fit is not as good, it would be chosen because its marginal likelihood is larger than C's. Maximization of the marginal likelihood favors a model that not only fits data well, but also does so over a wide range of parameters. Model D does so to a far greater extent than C. Comparisons such as this, in which goodness of fit and model complexity differ greatly between models, should benefit from application of the Bayesian approach.

In this last example, the Bayesian method might seem to yield counterintuitive results because the model that provides the best fit to the data is not chosen as the pre-

ferred model. This is because model complexity is just as important as model fit in influencing model selection. In the Bayesian approach, the ability of the model to fit the data reasonably well over a range of parameter values is what is important because it indicates that the model captures the structure in the data with a minimum of complexity (e.g., minimal reliance on functional form, number of parameters, etc.). Models in which parameters must be finely tuned to achieve a good fit (Models A and C) are valued less, precisely because the fit is achieved by relying on the complexity of the model. The contribution of complexity to model fit should not be underestimated. By finely tuning parameters, a complex model can fit a range of data patterns, not because it is the "true" model, but because it is the most flexible, easily adjusting to variations in the data. Taken to an extreme, one could imagine a highly complex model fitting almost any data set very well.

The preceding examples should make it clear that the difference in the pattern of data to be maximized by the likelihood function versus the marginal likelihood can have nontrivial implications for model selection and is at the heart of the Bayesian approach. Maximization of the marginal likelihood is accomplished by pitting maximization of the likelihood function against complexity minimization. A more complex model will be favored only if its goodness of fit is large enough to justify the additional complexity of the model. It is in this sense that Bayesian model selection is a quantitative implementation of Occam's razor.

### Computation of the Bayes Factor

Although the simplified expression in Equation 6 is useful to illustrate basic ideas of Bayesian model selection, it is derived under some restrictive assumptions about the shape of the likelihood function. Goodness of fit and model complexity are not, in general, separable from each other, but instead, the two work together implicitly within a single measure, which is the marginal likelihood. Computation of the Bayes factor  $P(D|M_1)/P(D|M_2)$  then requires evaluating the multiple integral in Equation 5. Readers who are not interested in the mathematical details of computing the Bayes factor should skip to the section, Application Example of Bayesian Model Selection.

**Priors  $P(\theta|M_i)$ .** To evaluate the integral, the prior density of the parameter,  $P(\theta|M_i)$ , must be specified. Two straightforward methods for obtaining the prior are described below. Others are discussed in Appendix B. Whatever method is used, it is very important to perform some form of sensitivity analysis to ensure that the resulting Bayes factor is reasonably stable over a range of values. For a more comprehensive treatment of this crucial topic, see Berger (1985) and Kass and Raftery (1995).

A simple-minded method for choosing priors is to use a noninformative prior, which by definition assumes no information about the parameter  $\theta$ . For example, the uniform density function  $\pi(\theta) = 1$  on  $0 \leq \theta \leq 1$  is a noninformative prior. Although a noninformative prior would be an obvious choice when absolutely no information is

available, often the noninformative prior is an *improper* prior having infinite mass (i.e.,  $\int \pi(\theta) d\theta = \infty$ ). One example is the uniform density  $\pi(\theta) = c$  on  $0 \leq \theta < \infty$  where  $c > 0$ . Another reasonable choice of a noninformative prior for  $0 \leq \theta < \infty$  that has the desirable scale-invariance property [i.e.,  $\pi(\theta) = \pi(\theta/\alpha)/\alpha$  for any  $\alpha > 0$ ] is  $\pi(\theta) = 1/\theta$ ; this is also an improper prior. Although the justification and interpretation of noninformative improper priors is disputed (see Berger, 1985, pp. 89–90), they are routinely used in Bayesian inference (e.g., Gelfand & Dey, 1994; H. Jeffreys, 1961).

Informative priors can be used when information about  $\theta$  is available. One way to do this is to obtain an informative prior from an experimental setup. To illustrate, consider a signal detection experiment consisting of  $n$  conditions in which the signal-to-noise ratio is systematically varied from lowest in Condition 1 to highest in Condition  $n$ . Assume that the number of correct responses out of a total of  $s$  trials in Condition  $k$  ( $k = 1, \dots, n$ ) is binomially distributed with  $Bin(s, \theta_k)$  where  $\theta_k$  ( $0 \leq \theta_k \leq 1$ ) is the binomial probability parameter. A reasonable way to obtain the prior  $\pi(\theta_1, \dots, \theta_n)$  is to assume that  $\theta_k$ s are the order statistics (i.e.,  $\theta_i \leq \theta_j$  for any  $i < j$ ), resulting from an independent random sample from the uniform distribution on  $[0, 1]$ . For example, given a particular random sample of three ( $n = 3$ ) observations, (.53, .29, .87), from the uniform distribution, the desired order statistics are obtained as  $\theta_1 = .29$ ,  $\theta_2 = .53$ , and  $\theta_3 = .87$ . For other similar cases, by being creative and reasonable, one can find priors that best capture the information available in the experimental design as well as in the data.

**Numerical methods for computing the Bayes factor.** A closed-form solution to the integral in Equation 5 is the preferred method for computing the Bayes factor. Most often, however, the integral must be evaluated numerically. A simple Monte Carlo integration method (Thisted, 1988, p. 302) can be used for integrals involving fewer than 10 parameters. More efficient methods such as Markov chain Monte Carlo (MCMC) methods should be used for integrals involving large numbers of parameters. For comprehensive treatments of MCMC methods and numerical integration methods in general, readers are directed to Smith (1991), Smith and Roberts (1993), and Thisted (1988, sec. 5.6).

The logic of the MCMC method is to obtain a sufficiently large sample of random observations generated from a probability density function in a *prespecified* manner, and then to use the sample to evaluate a desired integral. More specifically, the integral in Equation 5 is approximated by the following time average:

$$\frac{1}{T} \sum_{i=1}^T P(D|\theta^i, M_i) \quad (8)$$

In this equation,  $\{\theta^1, \theta^2, \dots, \theta^T\}$  is a random sample of  $T$  vectors drawn from the prior density function  $\pi(\theta)$  of Model  $i$ . Various versions of MCMC, such as the Gibbs sampler (Gelfand & Smith, 1990; Geman & Geman, 1984; Wakefield, Smith, Racine-Poon, & Gelfand, 1994) and the

Metropolis-Hastings algorithm (Hastings, 1970), have been proposed. The Gibbs sampler will be described here, primarily because it is easy to apply once the set of fully conditional distributions of the parameters (the probability distributions of each parameter conditional on all remaining parameters) is identified, at least up to proportionality.

The iteration procedure for generating a random sample of any size by the Gibbs sampler is summarized in the following steps. A hypothetical example of the Gibbs sampler is shown in square brackets, [...], for a three-parameter case with the parameters ( $\theta_1, \theta_2, \theta_3$ ) defined between 0 and 1:

*Initialization:*

Determine  $T$ ;

Pick an arbitrary starting vector  $\theta^0 = (\theta_1^0, \dots, \theta_n^0)$

[ $\theta^0 = (.52, .19, .34)$ ];

Define and set a temporary vector  $\theta^{\text{temp}} = \theta^0$

[ $\theta^{\text{temp}} = (.52, .19, .34)$ ];

Let  $\theta^{\text{temp}}\{-k\}$  denote a set of the current values in  $\theta^{\text{temp}}$  without the  $k$ th element

[e.g., for  $k = 2$ ,  $\theta^{\text{temp}}\{-2\} = (\theta_1^{\text{temp}} = .52, \theta_3^{\text{temp}} = .34)$ .]

Let  $\text{SUM} = 0$ ;

Let  $t = 1$ .

*Step 1:* Take a random sample of  $\theta_1^t$  from the conditional distribution  $\pi(\theta_1 | \theta^{\text{temp}}\{-1\})$

[ $\theta_1^t = .91$ ];

Replace  $\theta_1^{t-1}$  of  $\theta^{\text{temp}}$  with the new  $\theta_1^t$  [ $\theta^{\text{temp}} = (.91, .19, .34)$ ].

*Step 2:* Take a random sample of  $\theta_2^t$  from the conditional distribution  $\pi(\theta_2 | \theta^{\text{temp}}\{-2\})$

[ $\theta_2^t = .48$ ];

Replace  $\theta_2^{t-1}$  of  $\theta^{\text{temp}}$  with the new  $\theta_2^t$  [ $\theta^{\text{temp}} = (.91, .48, .34)$ ].

...

*Step  $n$ :* Take a random sample of  $\theta_n^t$  from the conditional distribution  $\pi(\theta_n | \theta^{\text{temp}}\{-n\})$

[ $\theta_n^t = .20$ ];

Replace  $\theta_n^{t-1}$  of  $\theta^{\text{temp}}$  with the new  $\theta_n^t$  [ $\theta^{\text{temp}} = (.91, .48, .20)$ ].

*Evaluation:*

Set  $\theta^t = \theta^{\text{temp}}$  and evaluate  $P(D|\theta^t, M_i)$  at  $\theta^t$

[ $\theta^t = (.91, .48, .20)$ ];

$\text{SUM} = \text{SUM} + P(D|\theta^t, M_i)$ ;

If  $t < T$  then let  $t = t + 1$  and go to Step 1 or else go to End.

*End:*

$\text{SUM} = \text{SUM} / T$ .

The value of the variable  $\text{SUM}$  represents a numerical solution to the integral in Equation 5 for the specified number of iterations ( $T$ ). The number of iterations nec-

essary to obtain an accurate approximation of the integral is not easy to determine beforehand. Raftery and Lewis (1991) have provided useful guidelines for estimating  $T$ , which depends on the type of problem, the model equation, and the specific priors used. Probably the most practical method is to examine the converging pattern of the sum by plotting it against the number of iterations to identify where the curve becomes stationary. A few runs usually give a good estimate for a given problem.

### APPLICATION EXAMPLE OF BAYESIAN MODEL SELECTION

In this section, an implementation of the Bayesian model selection method is illustrated through numerical examples with simulated data.

#### Models of Information Integration

A long-standing question that continues to receive considerable attention in cognitive psychology is, How is information from different sources (e.g., sensory and contextual) integrated during perception? The answer to this question has broad implications for theory development because information integration is a common denominator among models: Models must specify when and how independent sources of information are combined during processing. Entire classes of models will have to be modified if they are shown to have the wrong functional architecture.

Interest in information integration spawned a variety of mathematical models of integration (see N. H. Anderson, 1981; Massaro & Friedman, 1990). Their number and relative simplicity provided an appropriate context in which to test the Bayesian selection method. In a typical experiment on information integration, two or more stimulus dimensions are factorially manipulated and presented to one or two modalities (e.g., auditory, visual). Participants identify stimuli along one dimension. Of interest is the influence of the orthogonal dimension on perception. For example, context effects in speech perception are investigated by measuring categorization of perceptually ambiguous phonemes embedded in words and nonwords.

For simplicity, only data from a two-factor, two-response-category experiment will be considered. Extension of the approach to more general cases is straightforward. Suppose that the first factor has  $n_1$  levels and the second has  $n_2$  levels, so that  $S_{ij}$  ( $i = 1, \dots, n_1; j = 1, \dots, n_2$ ) denotes the stimulus constructed with Level  $i$  of the first factor and Level  $j$  of the second. The stimuli are presented to participants in a random order and equally often over  $s$  independent trials. Participants categorize the stimuli as one of two possibilities, A or B. Let a random variable  $X_{ij}$  denote the number of Category A responses participants made when Stimulus  $S_{ij}$  was presented. Then,  $X_{ij}$  will follow a binomial probability distribution with parameters  $p_{ij}$  (probability of Category A response) and  $s$  (number of independent observations).

Three non-nested models of information integration were chosen for evaluation on the basis of their similarities in implementation. They were the fuzzy logic model of perception (FLMP) by Oden and Massaro (1978), the linear integration model (LIM) by N. H. Anderson (1981), and the theory of signal detection (TSD) by Green and Swets (1966). All three assume that the probability of classifying a stimulus as a member of Category A is a function of the extent to which the two feature dimensions of the stimulus ( $i$  and  $j$ ) support the category response (Massaro & Friedman, 1990). Specifically, the response probability,  $p_{ij}$ , is assumed to be a function of two independent parameters,  $\theta_i$  and  $\lambda_j$ , each of which represents the degree of support for a Category A response given the specific  $i$  and  $j$  feature dimensions of a stimulus. The three models, however, differ in how the two parameters are combined to produce  $p_{ij}$ .

According to the FLMP,  $p_{ij}$  takes the following non-linear form (see Massaro & Friedman, 1990):

$$p_{ij, \text{FLMP}} = \frac{\theta_i \lambda_j}{\theta_i \lambda_j + (1 - \theta_i)(1 - \lambda_j)} \quad (9a)$$

The LIM assumes a linear combination rule for  $p_{ij}$  as follows:

$$p_{ij, \text{LIM}} = \frac{\theta_i + \lambda_j}{2} \quad (9b)$$

For the TSD, the response probability is given by

$$p_{ij, \text{TSD}} = \Phi \left[ s_{ij} \sqrt{|\Phi^{-1}(\theta_i)|^2 + v_{ij} \Phi^{-1}(\lambda_j)^2} \right] \quad (9c)$$

where  $\Phi(\cdot)$  and  $\Phi^{-1}(\cdot)$  are the cumulative and inverse cumulative normal functions, respectively. The  $s_{ij}$  is the sign ( $\pm 1$ ) of  $(\Phi^{-1}(\theta_i) + \Phi^{-1}(\lambda_j))$  and  $v_{ij}$  is the sign ( $\pm 1$ ) of  $(\Phi^{-1}(\theta_i) \cdot \Phi^{-1}(\lambda_j))$ .

Note that all three models possess the same number of parameters (two) as well as the same extension of the parameter space (i.e.,  $0 < \theta_i, \lambda_j < 1$ ), but differ in their functional form. Model selection using a standard method, such as AIC, would be decided solely on goodness of fit. Numerical examples in the following section demonstrate the improvement in model selection that is achieved with the Bayesian method.

#### Specification of Priors

To apply the Bayesian method, a probability density distribution for the two parameters of the model must be specified. We assumed that the distributions were the same for all three models, and further, that the parameters belonging to the first experimental variable were independent of those belonging to the second. In determining the characteristics of the distributions, a sensible choice is informative priors, because *prior* information about the parameters is available from the experimental setup that the present examples are intended to simulate. The method in which levels of the independent variable are created frequently ensures that ordinal information about the parameters is available, even before data are collected. For example, in speech perception experi-

ments, the levels of one of the variables (e.g., steps along a /ba/-/da/ phonetic continuum) are varied incrementally so that the probability of a Category A response at level  $i$  is greater than or equal to that at level  $i'$  if  $i < i'$ , and vice versa. The levels of the other variable can be manipulated in a similar fashion. Details on how this ordinal information was incorporated in the prior probabilities are provided in Appendix C.

### Simulated Data and Model Fitting

Data simulating three response patterns in a  $2 \times 8$  factorial design (i.e.,  $n_1 = 2$  and  $n_2 = 8$ ) were created. Three distinct parameter sets, chosen to represent a range of response patterns in categorization experiments, were used to generate the simulated data.<sup>8</sup> Each set contained 10 parameter values ( $2+8$ ). For each set, 16 ( $N = 16$ ) binomial response probabilities ( $p_{ij}$ ) were then computed using one of the three model equations in Equation 9. Each of the nine panels in Figure 5 shows a plot of the 16 probabilities for a given model and parameter set. Note that these probabilities represent ideal, error-free performance.

To simulate actual performance by human participants, samples were created using each parameter set in Figure 5 according to the binomial probability distribution as fol-

lows. For each of the 16 response probabilities ( $p_{ij}$ ), a series of 20 ( $s = 20$ ) independent binary outcomes (0 or 1) were generated in such a way that the probability of 1 was  $p_{ij}$ . Next, the number of 1s in the series was summed and divided by  $s$  to obtain an observed proportion for the particular combination of  $i$  and  $j$ . Finally, the above procedure was repeated for the remaining 15 values to obtain 16 observed proportions, which together constituted a single sample. One hundred samples were created for each of the nine panels in Figure 5.

Each of the three models was fitted to each simulated sample separately. This was done using the standard method and the Bayesian method to compare their abilities to recover (i.e., identify) the model that generated the original data. For the standard method, a nonlinear optimization routine (Marquardt compromise method; Marquardt, 1963) was used to find the least squares estimates that minimized the sum of squared errors between the simulated and predicted data. Because the number of parameters was the same across models, this procedure was equivalent to other methods, such as AIC. For the Bayesian method, the integral in Equation 5 was evaluated numerically using an extended series of simple Monte Carlo simulations.<sup>9</sup>

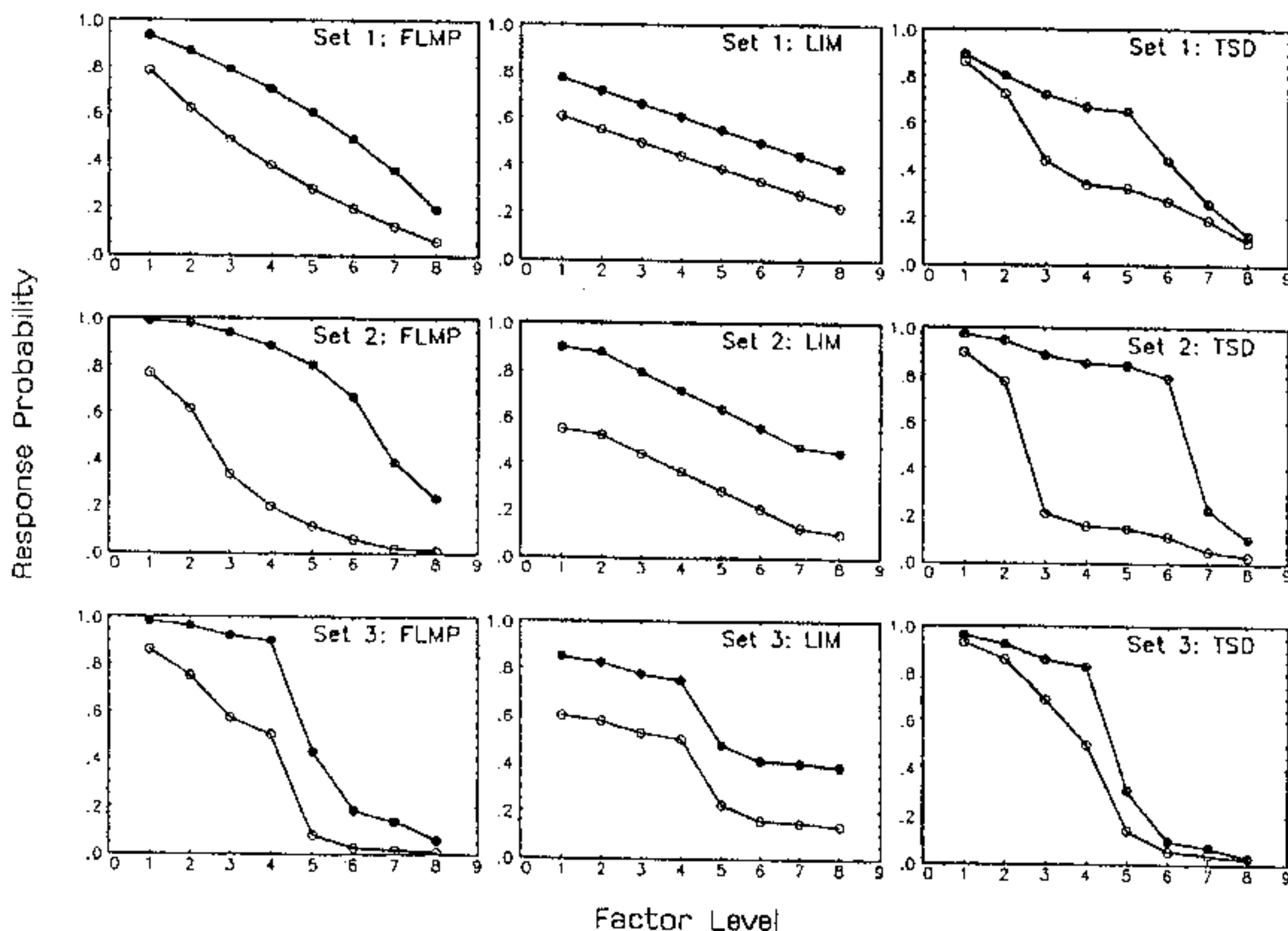


Figure 5. Binary response probabilities used to create the simulated data. The three panels in each row represent three different response patterns generated from the fuzzy logic model of perception (FLMP), the linear integration model (LIM), and the theory of signal detection (TSD) using Equations 9a-9c with a single set of 10 parameter values. Three sets of parameter values, corresponding to the three rows in the figure, were used.

