



UvA-DARE (Digital Academic Repository)

Bayesian explorations in mathematical psychology

Matzke, D.

Publication date

2014

Document Version

Final published version

[Link to publication](#)

Citation for published version (APA):

Matzke, D. (2014). *Bayesian explorations in mathematical psychology*. [Thesis, fully internal, Universiteit van Amsterdam].

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, P.O. Box 19185, 1000 GD Amsterdam, The Netherlands. You will be contacted as soon as possible.

Model Comparison and the Principle of Parsimony

This chapter is a modified version of:
Joachim Vandekerckhove, Dora Matzke, and Eric-Jan Wagenmakers (in press).
Model comparison and the principle of parsimony.
In J. Busemeyer, J. Townsend, Z. J. Wang, & A. Eidels (Eds.), *Oxford handbook of computational and mathematical psychology*. Oxford, UK: Oxford University Press.

6.1 Introduction

At its core, the study of psychology is concerned with the discovery of plausible explanations for human behavior. For instance, one may observe that “practice makes perfect”: as people become more familiar with a task, they tend to execute it more quickly and with fewer errors. More interesting is the observation that practice tends to improve performance such that most of the benefit is accrued early on, a pattern of diminishing returns that is well described by a power law (Logan, 1988; but see Heathcote et al., 2000). This pattern occurs across so many different tasks (e.g., cigar rolling, maze solving, fact retrieval, and a variety of standard psychological tasks) that it is known as the “power law of practice”. Consider, for instance, the lexical decision task, a task in which participants have to decide quickly whether a letter string is an existing word (e.g., *sunscreen*) or not (e.g., *tolphin*). When repeatedly presented with the same stimuli, participants show a power law decrease in their mean response latencies; in fact, they show a power law decrease in the entire response time distribution, that is, both the fast responses and the slow responses speed up with practice according to a power law (Logan, 1992).

The observation that practice makes perfect is trivial, but the finding that practice-induced improvement follows a general law is not. Nevertheless, the power law of practice only provides a descriptive summary of the data and does not explain the reasons why practice should result in a power law improvement in performance. In order to go beyond direct observation and statistical summary, it is necessary to bridge the divide between observed performance on the one hand and the pertinent psychological processes on the other. Such bridges are built from a coherent set of assumptions about the underlying cognitive processes—a theory. Ideally, substantive psychological theories are formalized as quantitative models (Busemeyer & Diederich, 2010; Lewandowsky & Farrell, 2010). For example, the power law of practice has been explained by instance theory (Logan, 1992, 2002). Instance theory stipulates that earlier experiences are stored in memory as

individual traces or instances; upon presentation of a stimulus, these instances race to be retrieved, and the winner of the race initiates a response. Mathematical analysis shows that, as instances are added to memory, the finishing time of the winning instance decreases as a power function. Hence, instance theory provides a simple and general explanation of the power law of practice.

For all its elegance and generality, instance theory has not been the last word on the power law of practice. The main reason is that single phenomena often afford different competing explanations. For example, the effects of practice can also be accounted for by Rickard's component power laws model (Rickard, 1997), Anderson's ACT-R model (Anderson et al., 2004), Cohen et al.'s PDP model (J. D. Cohen, Dunbar, & McClelland, 1990), Ratcliff's diffusion model (Dutilh, Vandekerckhove, Tuerlinckx, & Wagenmakers, 2009; Ratcliff, 1978), or Brown and Heathcote's linear ballistic accumulator model (Brown & Heathcote, 2005, 2008; Heathcote & Hayes, 2012). When various models provide competing accounts of the same data set, it can be difficult to choose between them. The process of choosing between models is called model comparison, model selection, or hypothesis testing, and it is the focus of this chapter.

A careful model comparison procedure includes both qualitative and quantitative elements. Important qualitative elements include the plausibility, parsimony, and coherence of the underlying assumptions, the consistency with known behavioral phenomena, the ability to explain rather than describe data, and the extent to which model predictions can be falsified through experiments. Here we ignore these important aspects and focus solely on the quantitative elements. The single most important quantitative element of model comparison relates to the ubiquitous tradeoff between parsimony and goodness-of-fit (Pitt & Myung, 2002). The motivating insight is that the appeal of an excellent fit to the data (i.e., high descriptive adequacy) needs to be tempered to the extent that the fit was achieved with a highly complex and powerful model (i.e., low parsimony).

The topic of quantitative model comparison is as important as it is challenging; fortunately, the topic has received—and continues to receive—considerable attention in the field of statistics, and the results of those efforts have been made accessible to psychologists through a series of recent special issues, books, and articles (e.g., Grünwald, 2007; Myung, Forster, & Browne, 2000; Pitt & Myung, 2002; Wagenmakers & Waldorp, 2006). Here we discuss several procedures for model comparison, with an emphasis on minimum description length and the Bayes factor. Both procedures entail principled and general solutions to the tradeoff between parsimony and goodness-of-fit.

The outline of this chapter is as follows. The first section describes the principle of parsimony and the unavoidable tradeoff with goodness-of-fit. The second section summarizes the research of Wagenaar and Boer (1987) who carried out an experiment to compare three competing multinomial processing tree models (MPTs; Batchelder & Riefer, 1980); this model comparison exercise is used as a running example throughout the chapter. The third section outlines different methods for model comparison and applies them to Wagenaar and Boer's MPT models. We focus on two popular information criteria, the AIC and the BIC, on the Fisher information approximation of the minimum description length principle, and on Bayes factors as obtained from importance sampling. The fourth section contains conclusions and take-home messages.

6.2 The Principle of Parsimony

Throughout history, prominent philosophers and scientists have stressed the importance of parsimony. For instance, in the *Almagest*—a famous 2nd-century book on astronomy—Ptolemy writes: “We consider it a good principle to explain the phenomena by the simplest hypotheses that can be established, provided this does not contradict the data in an important way.” Ptolemy's principle

Occam's razor (sometimes *Ockham's*) is named after the English philosopher and Franciscan friar Father William of Occam (c.1288-c.1348), who wrote "Numquam ponenda est pluralitas sine necessitate" (plurality must never be posited without necessity), and "Frustra fit per plura quod potest fieri per pauciora" (it is futile to do with more what can be done with less). Occam's metaphorical razor symbolizes the principle of parsimony: by cutting away needless complexity, the razor leaves only theories, models, and hypotheses that are as simple as possible without being false. Throughout the centuries, many other scholars have espoused the principle of parsimony; the list predating Occam includes Aristotle, Ptolemy, and Thomas Aquinas ("it is superfluous to suppose that what can be accounted for by a few principles has been produced by many"), and the list following Occam includes Isaac Newton ("We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances. Therefore, to the same natural effects we must, so far as possible, assign the same causes."), Bertrand Russell, Albert Einstein ("Everything should be made as simple as possible, but no simpler"), and many others.

In the field of statistical reasoning and inference, Occam's razor forms the foundation for the principle of minimum description length (Grünwald, 2000, 2007). In addition, Occam's razor is automatically accommodated through Bayes factor model comparisons (e.g., Jeffreys, 1961; Jefferys & Berger, 1992; MacKay, 2003). Both minimum description length and Bayes factors feature prominently in this chapter as principled methods to quantify the tradeoff between parsimony and goodness-of-fit.

Box 6.1 Occam's razor.

of parsimony is widely known as Occam's razor (see Box 6.1); the principle is intuitive as it puts a premium on elegance. In addition, most people feel naturally attracted to models and explanations that are easy to understand and communicate. Moreover, the principle also gives ground to reject propositions that are without empirical support, including extrasensory perception, alien abductions, or mysticism. In an apocryphal interaction, Napoleon Bonaparte asked Pierre-Simon Laplace why the latter's book on the universe did not mention its creator, only to receive the curt reply "I had no need of that hypothesis".

However, the principle of parsimony finds its main motivation in the benefits that it bestows those who use models for prediction. To see this, note that empirical data are composed of a structural, replicable part and an idiosyncratic, non-replicable part. The former is known as the signal, and the latter is known as the noise (Silver, 2012). Models that capture all of the signal and none of the noise provide the best possible predictions to unseen data from the same source. Overly simplistic models, however, fail to capture part of the signal; these models underfit the data and provide poor predictions. Overly complex models, on the other hand, mistake some of the noise for actual signal; these models overfit the data and again provide poor predictions. Thus, parsimony is essential because it helps discriminate the signal from the noise, allowing better prediction and generalization to new data.

Goodness-of-Fit

"From the earliest days of statistics, statisticians have begun their analysis by proposing a distribution for their observations and then, perhaps with somewhat less enthusiasm, have checked on whether this distribution is true. Thus over the years a vast number of test procedures have appeared, and the study of these procedures has come to be known as goodness-of-fit" (D'Agostino & Stephens, 1986, p. v).

The *goodness-of-fit* of a model is a quantity that expresses how well the model is able to account for a given set of observations. It addresses the following question: Under the assumption that a certain model is a true characterization of the population from which we have obtained a sample, and given the best fitting parameter estimates for that model, how well does our sample of data agree with that model?

Various ways of quantifying goodness-of-fit exist. One common expression involves a Euclidean distance metric between the data and the model's best prediction (the least squared error or LSE metric is the most well-known of these). Another measure involves the likelihood function, which expresses the likelihood of observing the data under the model, and is maximized by the best fitting parameter estimates (Myung, 2000).

Parsimony

Goodness-of-fit must be balanced against model complexity in order to avoid overfitting—that is, to avoid building models that well explain the data at hand, but fail in out-of-sample predictions. The principle of parsimony forces researchers to abandon complex models that are tweaked to the observed data in favor of simpler models that can generalize to new data sets.

A common example is that of polynomial regression. Figure 6.1 gives a typical example. The observed data are the circles in both the left and right panels. Crosses indicate unobserved, out-of-sample data points to which the model should generalize. In the left panel, a quadratic function is fit to the 8 observed data points, whereas the right panel shows a 7th order polynomial function fitted to the same data. Since a polynomial of degree 7 can be made to contain any 8 points in the plane, the observed data are perfectly captured by the best fitting polynomial. However, it is clear that this function generalizes poorly to the unobserved samples, and it shows undesirable behavior for larger values of x .

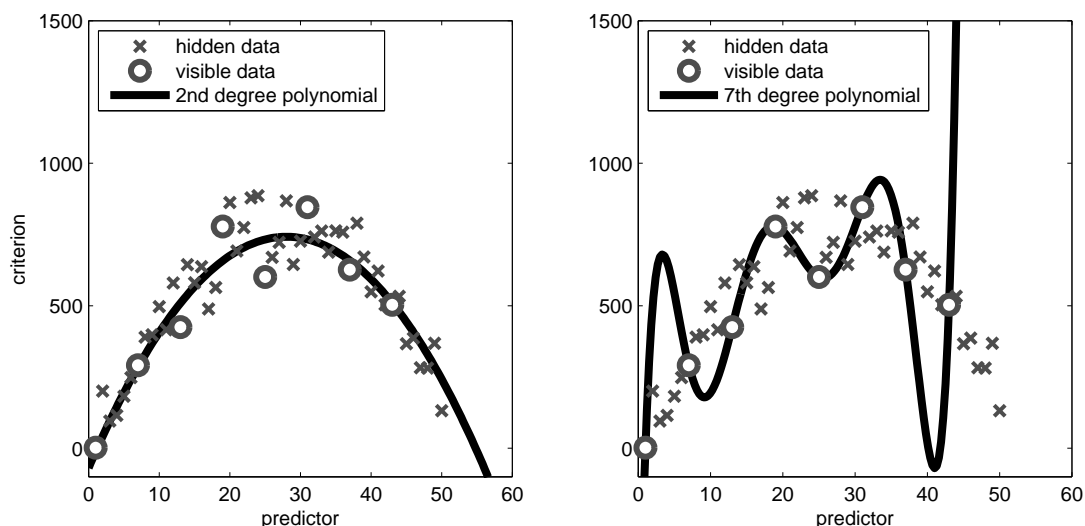


Figure 6.1 A polynomial regression of degree d is characterized by $\hat{y} = \sum_{i=0}^d a_i x^i$. This model has $d + 1$ free parameters a_i ; hence, in the right panel, a polynomial of degree 7 perfectly accounts for the 8 visible data points. This 7th order polynomial, however, accounts poorly for the out-of-sample data points.

In sum, an adequate model comparison method needs to discount goodness-of-fit with model complexity. But how exactly can this be accomplished? As we will describe shortly, several model comparison methods are currently in vogue; all resulting from principled ideas on how to obtain *measures of generalizability*¹, meaning that these methods attempt to quantify the extent to which a model predicts unseen data from the same source (cf. Figure 6.1). Before outlining the details of various model comparison methods, we now introduce a data set that serves as a working example throughout the remainder of the chapter.

6.3 Example: Competing Models of Interference in Memory

For an example model comparison scenario, we revisit a study by Wagenaar and Boer (1987) on the effect of misleading information on the recollection of an earlier event. The effect of misleading postevent information was first studied systematically by E. F. Loftus, Miller, and Burns (1978); for a review of relevant literature, see Wagenaar and Boer (1987) and references therein.

Wagenaar and Boer (1987) proposed three competing theoretical accounts of the effect of misleading postevent information. To evaluate the three accounts, Wagenaar and Boer set up an experiment and introduced three quantitative models that translate each of the theoretical accounts into a set of parametric assumptions that together give rise to a probability density over the data, given the parameters.

Abstract Accounts

Wagenaar and Boer (1987) outlined three competing theoretical accounts of the effect of misleading postevent information on memory. Loftus' *destructive updating* model (DUM) posits that the conflicting information replaces and destroys the original memory. A *coexistence* model (CXM) asserts that an inhibition mechanism suppresses the original memory, which nonetheless remains viable though temporarily inaccessible. Finally, a *no-conflict* model (NCM) simply states that misleading postevent information is ignored, except when the original information was not encoded or already forgotten.

Experimental Design

The experiment by Wagenaar and Boer (1987) proceeded as follows. In Phase I, a total of 562 participants were shown a sequence of events in the form of a pictorial story involving a pedestrian-car collision. One picture in the story would show a car at an intersection, and a traffic light that was either red, yellow, or green. In Phase II, participants were asked a set of test questions with (potentially) conflicting information: Participants might be asked whether they remembered a pedestrian crossing the road when the car approached the “traffic light” (in the consistent group), the “stop sign” (in the inconsistent group) or the “intersection” (the neutral group). Then, in Phase III, participants were given a recognition test about elements of the story using picture pairs. Each pair would contain one picture from Phase I and one slightly altered version of the original picture. Participants were then asked to identify which of the pair had featured in the original story. A picture pair is shown in Figure 6.2, where the intersection is depicted with either a traffic light or a stop sign. Finally, in Phase IV, participants were informed that the correct choice in Phase III was the picture with the traffic light, and were then asked to recall the color of the traffic light.

¹This terminology is due to Pitt and Myung (2002), who point out that measures often referred to as “model fit indices” are in fact more than mere measures of fit to the data—they combine fit to the data with parsimony and hence measure generalizability. We adopt their more accurate terminology here.

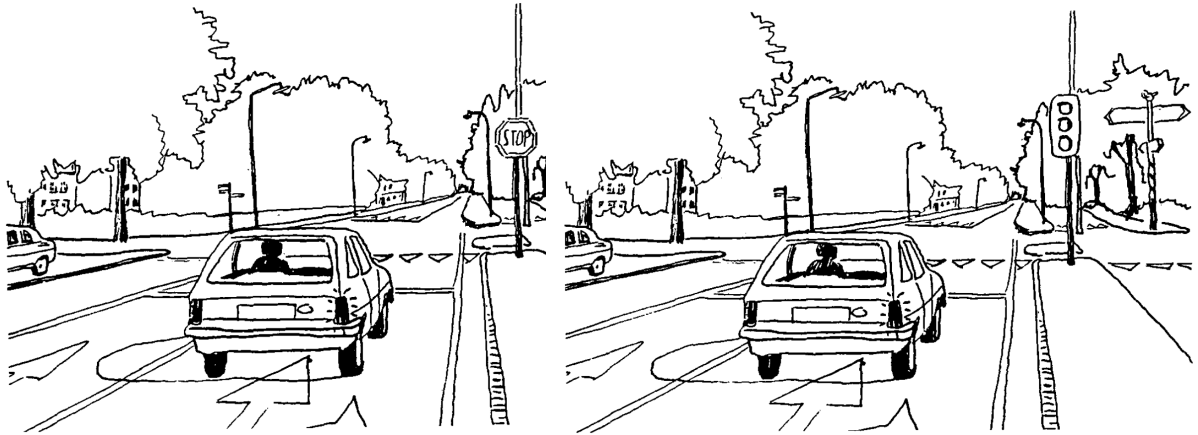


Figure 6.2 A pair of pictures from the third phase (i.e., the recognition test) of (Wagenaar & Boer, 1987, reprinted with permission), containing the critical episode at the intersection.

By design, this experiment should yield different response patterns depending on whether the conflicting postevent information destroys the original information (destructive updating model), only suppresses it temporarily (coexistence model), or does not affect the original information unless it is unavailable (no-conflict model).

Concrete Models

Wagenaar and Boer (1987) developed a series of MPT models (see Box 6.2) to quantify the predictions of the three competing theoretical accounts. Figure 6.3 depicts the no-conflict MPT model in the inconsistent condition. The figure is essentially a decision tree that is navigated from left to right. In Phase I of the collision narrative, the traffic light is encoded with probability p , and if so, the color is encoded with probability c . In Phase II, the stop sign is encoded with probability q . In Phase III, the answer may be known, or may be guessed correctly with probability $1/2$, and in Phase IV the answer may be known or may be guessed correctly with probability $1/3$. The probability of each path is given by the product of all the encountered probabilities, and the total probability of a response pattern is the summed probability of all branches that lead to it. For example, the total probability of getting both questions wrong is $(1-p) \times q \times 2/3 + (1-p) \times (1-q) \times 1/2 \times 2/3$. We would then, under the no-conflict model, expect that proportion of participants to fall in the response pattern with two errors.

The destructive updating model (Figure 2 in Wagenaar & Boer, 1987) extends the three-parameter no-conflict model by adding a fourth parameter d : the probability of destroying the traffic light information, which may occur whenever the stop sign was encoded. The coexistence model (Figure 3 in Wagenaar & Boer, 1987), on the other hand, posits an extra probability s that the traffic light is suppressed (but not destroyed) when the stop sign is encoded. A critical difference between the latter two is that a destruction step will lead to chance accuracy in Phase IV if every piece of information was encoded, whereas a suppression step will not affect the underlying memory and lead to accurate responding. Note here that if $s = 0$, the coexistence model reduces to the no-conflict model, as does the destructive updating model with $d = 0$. The models only make different predictions in the inconsistent condition, so that for the consistent and neutral conditions the trees are identical.

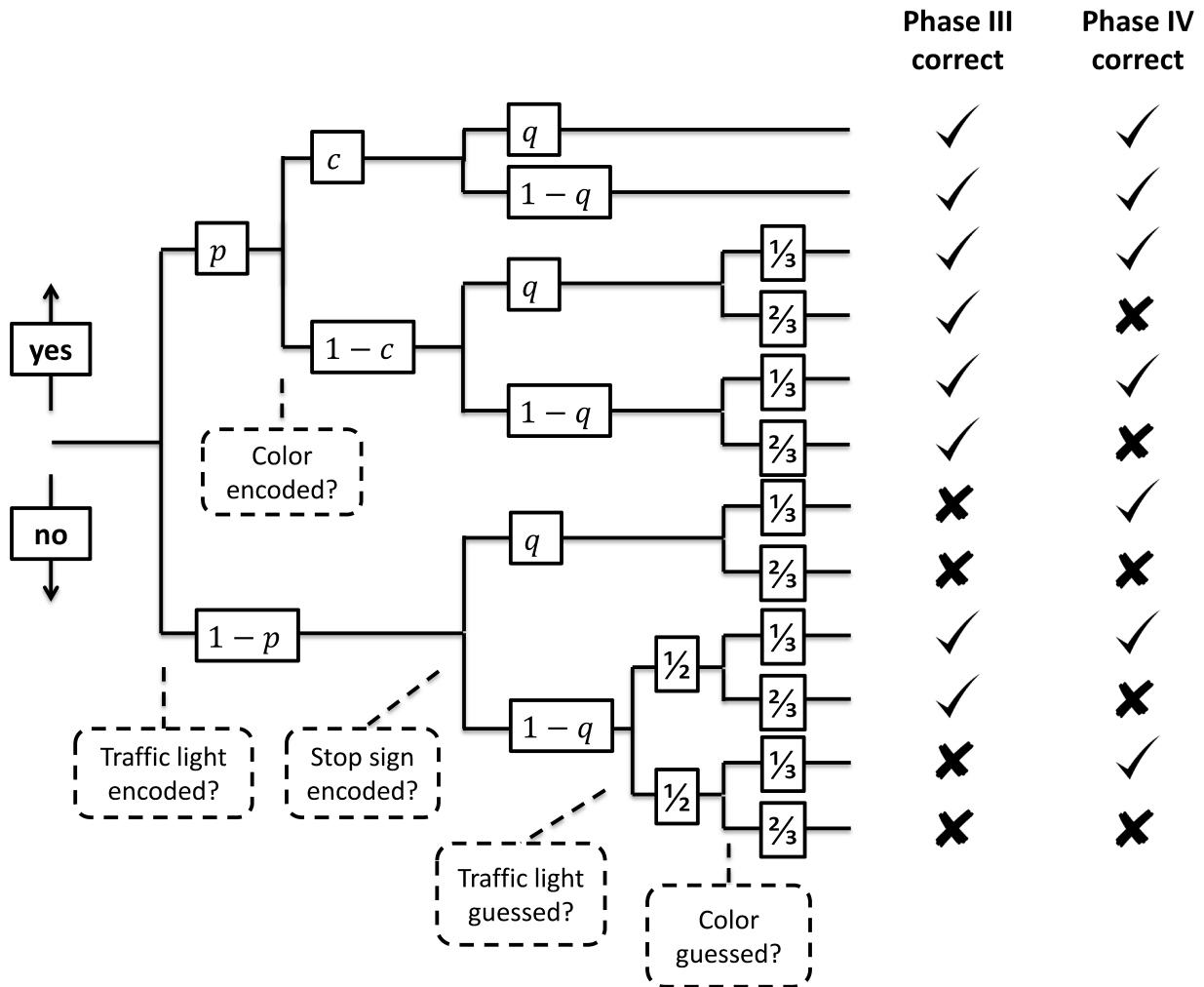


Figure 6.3 Multinomial processing tree representation of the inconsistent condition according to the no-conflict model (adapted from Wagenaar & Boer, 1987).

Previous Conclusions

After fitting the three competing MPT models, Wagenaar and Boer (1987) obtained the parameter point estimates in Table 6.1. Using a χ^2 model fit index, they concluded that “a distinction among the three model families appeared to be impossible in actual practice” (p. 304), after noting that the no-conflict model provides “an almost perfect fit” to the data. They propose, then, “to accept the most parsimonious model, which is the no-conflict model.” In the remainder of this chapter, we re-examine this conclusion using various model comparison methods.

6.4 Three Methods for Model Comparison

Many model comparison methods have been developed, all of them attempts to address the ubiquitous tradeoff between parsimony and goodness-of-fit. Here we focus on three main classes of interrelated methods: (1) AIC and BIC, the most popular information criteria; (2) minimum de-

Table 6.1 Parameter Point Estimates From Wagenaar and Boer (1987).

	p	c	q	d	s
No-conflict model (NCM)	0.50	0.57	0.50	n/a	n/a
Destructive updating model (DUM)	0.50	0.57	0.50	0.00	n/a
Coexistence model (CXM)	0.55	0.55	0.43	n/a	0.20

Multinomial processing tree models (Batchelder & Riefer, 1980; Chechile, 1973; Chechile & Meyer, 1976; Riefer & Batchelder, 1988) are psychological process models for categorical data. MPT models are used in two ways: as a psychometric tool to measure unobserved cognitive processes, and as a convenient formalization of competing psychological theories. Over time, MPTs have been applied to a wide range of psychological tasks and processes. For instance, MPT models are available for recognition, recall, source monitoring, perception, priming, reasoning, consensus analysis, the process dissociation procedure, implicit attitude measurement, and many other phenomena. For more information about MPTs, we recommend the review articles by Batchelder and Riefer (1999), Batchelder and Riefer (2007, pp. 24–32), and Erdfelder et al. (2009). The latter review article also discusses different software packages that can be used to fit MPT models. Necessarily missing from that list is the recently developed R package `MPTinR` (Singmann & Kellen, 2013) with which we have good experiences. As will become apparent throughout this chapter, however, our preferred method for fitting MPT models is Bayesian (Chechile & Meyer, 1976; Klauer, 2010; M. D. Lee & Wagenmakers, 2013; Matzke, Dolan, Batchelder, & Wagenmakers, in press; Rouder et al., 2008; J. B. Smith & Batchelder, 2010).

Box 6.2 Popularity of multinomial processing tree models.

scription length; (3) Bayes factors. Below we provide a brief description of each method and then apply it to the model comparison problem that confronted Wagenaar and Boer (1987).

Information Criteria

Information criteria are among the most popular methods for model comparison. Their popularity is explained by the simple and transparent manner in which they quantify the tradeoff between parsimony and goodness-of-fit. Consider for instance the oldest information criterion, AIC (“an information criterion”), proposed by Akaike (1973, 1974a):

$$\text{AIC} = -2 \ln p(y | \hat{\theta}) + 2k. \quad (6.1)$$

The first term $\ln p(y | \hat{\theta})$ is the log maximum likelihood that quantifies goodness-of-fit, where y is the data set and $\hat{\theta}$ the maximum-likelihood parameter estimate; the second term $2k$ is a penalty for model complexity, measured by the number of adjustable model parameters k . The AIC estimates the expected information loss incurred when a probability distribution f (associated with the true data-generating process) is approximated by a probability distribution g (associated with the model under evaluation). Hence, the model with the lowest AIC is the model with the smallest expected information loss between reality f and model g , where the discrepancy is quantified by the Kullback-Leibler divergence $I(f, g)$ (for full details, see Burnham & Anderson, 2002). The AIC is unfortunately not *consistent*: as the number of observations grows infinitely large, AIC is

not guaranteed to choose the true data generating model. Many researchers believe that the AIC tends to select complex models that overfit the data (O’Hagan & Forster, 2004; for a discussion see Vrieze, 2012).

Another information criterion, the BIC (“Bayesian information criterion”) was proposed by G. Schwarz (1978):

$$\text{BIC} = -2 \ln p(y | \hat{\theta}) + k \ln n. \quad (6.2)$$

Here, the penalty term is $k \ln n$, where n is the number of observations. Hence, the BIC penalty for complexity increases with sample size, outweighing that of AIC as soon as $n \geq 8$. The BIC was derived as an approximation of a Bayesian hypothesis test using default parameter priors (the “unit information prior”; see below for more information on Bayesian hypothesis testing, and see Raftery, 1995, for more information on the BIC). The BIC is consistent: as the number of observations grows infinitely large, BIC is guaranteed to choose the true data generating model. Nevertheless, some researchers believe that in practical applications the BIC tends to select simple models that underfit the data (Burnham & Anderson, 2002).

Now consider a set of candidate models, $\mathcal{M}_i, i = 1, \dots, m$, each with a specific IC (AIC or BIC) value. The model with the smallest IC value should be preferred, but the extent of this preference is not immediately apparent. For better interpretation we can calculate IC model weights (Akaike, 1974b; Burnham & Anderson, 2002; Wagenmakers & Farrell, 2004); First, we compute, for each model i , the difference in IC with respect to the IC of the best candidate model:

$$\Delta_i = \text{IC}_i - \min \text{IC}. \quad (6.3)$$

This step is taken to increase numerical stability, but it also serves to emphasize the point that only differences in IC values are relevant. Next, we obtain the model weights by transforming back to the likelihood scale and normalizing:

$$w_i = \frac{\exp(-\Delta_i/2)}{\sum_{m=1}^M \exp(-\Delta_m/2)}. \quad (6.4)$$

The resulting AIC and BIC weights are called Akaike weights and Schwarz weights, respectively. These weights not only convey the relative preference among a set of candidate models, but also provide a method to combine predictions across multiple models using model averaging (Hoeting, Madigan, Raftery, & Volinsky, 1999). Both AIC and BIC rely on an assessment of model complexity that is relatively crude, as it is determined entirely by the number of free parameters but not by their functional form.

Application to Multinomial Processing Tree Models

In order to apply AIC and BIC to the three competing MPTs proposed by Wagenaar and Boer (1987), we first need to compute the maximum log likelihood. Note that the MPT model parameters determine the predicted probabilities for the different response outcome categories (cf. Figure 6.3 and Box 6.2); these predicted probabilities are deterministic parameters from a multinomial probability density function. Hence, the maximum log likelihood parameter estimates for an MPT model produce multinomial parameters that maximize the probability of the observed data (i.e., the occurrence of the various outcome categories).

Several software packages exist that can help find the maximum log likelihood parameter estimates for MPTs (e.g. Singmann & Kellen, 2013). With these estimates in hand, we can compute the information criteria described in the previous section. Table 6.2 shows the maximum log likelihood as well as AIC, BIC, and their associated weights (wAIC and wBIC; from Equation 6.4).

Table 6.2 AIC and BIC for the Wagenaar and Boer (1987) MPT Models.

	Log likelihood	k	AIC	wAIC	BIC	wBIC
No-conflict model (NCM)	-24.41	3	54.82	0.41	67.82	0.86
Destructive updating model (DUM)	-24.41	4	56.82	0.15	74.15	0.04
Coexistence model (CXM)	-23.35	4	54.70	0.44	72.03	0.10

Note. k is the number of free parameters.

Interpreting wAIC and wBIC as measures of relative preference, we see that the results in Table 6.2 are mostly inconclusive. According to wAIC, the no-conflict model and coexistence model are virtually indistinguishable, though both are preferable to the destructive updating model. According to wBIC, however, the no-conflict model should be preferred over both the destructive updating model and the coexistence model. The extent of this preference is noticeable but not decisive.

Minimum Description Length

The minimum description length principle is based on the idea that statistical inference centers around capturing regularity in data; regularity, in turn, can be exploited to compress the data. Hence, the goal is to find the model that compresses the data the most (Grünwald, 2007). This is related to the concept of Kolmogorov complexity—for a sequence of numbers, Kolmogorov complexity is the length of the shortest program that prints that sequence and then halts (Grünwald, 2007). Although Kolmogorov complexity cannot be calculated, a suite of concrete methods are available based on the idea of model selection through data compression. These methods, most of them developed by Jorma Rissanen, fall under the general heading of minimum description length (MDL; Rissanen, 1978, 1987, 1996, 2001). In psychology, the MDL principle has been applied and promoted primarily by Grünwald (2000), Grünwald (2007), Grünwald, Myung, and Pitt (2005), as well as Myung, Navarro, and Pitt (2006), Pitt and Myung (2002), and Pitt, Myung, and Zhang (2002).

Here we mention three versions of the MDL principle. First, there is the so-called *crude two-part code* (Grünwald, 2007); here, one sums the description of the model (in bits) and the description of the data encoded with the help of that model (in bits). The penalty for complex models is that they take many bits to describe, increasing the summed code length. Unfortunately, it can be difficult to define the number of bits required to describe a model.

Second, there is the Fisher information approximation (FIA; Pitt et al., 2002; Rissanen, 1996):

$$\text{FIA} = -\ln p(y | \hat{\theta}) + \frac{k}{2} \ln \left(\frac{n}{2\pi} \right) + \ln \int_{\Theta} \sqrt{\det [I(\theta)]} d\theta, \quad (6.5)$$

where $I(\theta)$ is the Fisher information matrix of sample size 1. Note that FIA is similar to AIC and BIC in that it includes a first term that represents goodness-of-fit, and additional terms that represent a penalty for complexity. The second term resembles that of BIC, and the third term reflects a more sophisticated penalty that represents the number of distinguishable probability distributions that a model can generate (Pitt et al., 2002). Hence, FIA differs from AIC and BIC in that it also accounts for functional form complexity, not just complexity due to the number of free parameters. Note that FIA weights (or Rissanen weights) can be obtained by multiplying FIA by 2 and then applying Equations 6.3 and 6.4.

Table 6.3 Minimum Description Length Values for the Wagenaar and Boer (1987) MPT Models.

	Complexity	FIA	wFIA
No-conflict model (NCM)	6.44	30.86	0.44
Destructive updating model (DUM)	7.39	31.80	0.17
Coexistence model (CXM)	7.61	30.96	0.39

The third version of the MDL principle discussed here is normalized maximum likelihood (NML; Myung et al., 2006; Rissanen, 2001):

$$\text{NML} = \frac{p(y | \hat{\theta}_y)}{\int_x p(x | \hat{\theta}_x)}. \quad (6.6)$$

This equation shows that NML tempers the enthusiasm about a good fit to the observed data y (i.e., the numerator) to the extent that the model could also have provided a good fit to random data x (i.e., the denominator). NML is simple to state but can be difficult to compute; for instance, the denominator may be infinite and this requires additional measures to be taken (for details, see Grünwald, 2007).

Application to Multinomial Processing Tree Models

Using the parameter estimates from Table 6.1 and the code provided by Wu, Myung, and Batchelder (2010), we can compute the FIA for the three competing MPT models considered by Wagenaar and Boer (1987).² Table 6.3 displays, for each model, the FIA along with its associated complexity measure (the other one of its two constituent components, the maximum log likelihood, can be found in Table 6.2). The conclusions from the MDL analysis mirror those from the AIC measure, expressing a slight disfavor for the destructive updating model, and approximately equal preference for the no-conflict model versus the coexistence model.

Bayes Factors

In Bayesian model comparison, the posterior odds for models \mathcal{M}_1 and \mathcal{M}_2 are obtained by updating the prior odds with the diagnostic information from the data:

$$\frac{p(\mathcal{M}_1 | y)}{p(\mathcal{M}_2 | y)} = \frac{p(\mathcal{M}_1)}{p(\mathcal{M}_2)} \times \frac{m(y | \mathcal{M}_1)}{m(y | \mathcal{M}_2)}. \quad (6.7)$$

Equation 6.7 shows that the change from prior odds $p(\mathcal{M}_1)/p(\mathcal{M}_2)$ to posterior odds $p(\mathcal{M}_1 | y)/p(\mathcal{M}_2 | y)$ is given by the ratio of marginal likelihoods $m(y | \mathcal{M}_1)/m(y | \mathcal{M}_2)$, a quantity known as the *Bayes factor* (Jeffreys, 1961; Kass & Raftery, 1995). The log of the Bayes factor is often interpreted as the weight of evidence provided by the data (Good, 1985; for details, see Berger & Pericchi, 1996; Bernardo & Smith, 1994; Gill, 2002; O’Hagan, 1995).

Thus, when the Bayes factor $BF_{12} = m(y | \mathcal{M}_1)/m(y | \mathcal{M}_2)$ equals 5, the observed data y are 5 times more likely to occur under \mathcal{M}_1 than under \mathcal{M}_2 ; when BF_{12} equals 0.1, the observed data are 10 times more likely under \mathcal{M}_2 than under \mathcal{M}_1 . Even though the Bayes factor has an

²Analysis using the `MPTinR` package from Singmann and Kellen (2013) gave virtually identical results.

unambiguous and continuous scale, it is sometimes useful to summarize the Bayes factor in terms of discrete categories of evidential strength. Jeffreys (1961, Appendix B) proposed the classification scheme shown in Table 6.4. We replaced the labels “not worth more than a bare mention” with “anecdotal”, “decisive” with “extreme”, and “substantial” with “moderate”. These labels facilitate scientific communication but should be considered only as an approximate descriptive articulation of different standards of evidence.

Table 6.4 Evidence Categories for the Bayes Factor BF_{12} (Based on Jeffreys, 1961).

Bayes factor BF_{12}		Interpretation
$>$	100	Extreme evidence for \mathcal{M}_1
30	— 100	Very strong evidence for \mathcal{M}_1
10	— 30	Strong evidence for \mathcal{M}_1
3	— 10	Moderate evidence for \mathcal{M}_1
1	— 3	Anecdotal evidence for \mathcal{M}_1
	1	No evidence
1/3	— 1	Anecdotal evidence for \mathcal{M}_2
1/10	— 1/3	Moderate evidence for \mathcal{M}_2
1/30	— 1/10	Strong evidence for \mathcal{M}_2
1/100	— 1/30	Very strong evidence for \mathcal{M}_2
$<$	1/100	Extreme evidence for \mathcal{M}_2

Bayes factors negotiate the tradeoff between parsimony and goodness-of-fit and implement an automatic Occam’s razor (Jefferys & Berger, 1992; MacKay, 2003; Myung & Pitt, 1997). To see this, consider that the marginal likelihood $m(y)$ can be expressed as $\int_{\Theta} p(y | \theta)p(\theta) d\theta$: an average across the entire parameter space, with the prior providing the averaging weights. It follows that complex models with high-dimensional parameter spaces are not necessarily desirable—large regions of the parameter space may yield a very poor fit to the data, dragging down the average. The marginal likelihood will be highest for parsimonious models that use only those parts of the parameter space that are required to provide an adequate account of the data (M. D. Lee & Wagenmakers, 2013). By using marginal likelihood, the Bayes factor punishes models that hedge their bets and make vague predictions. Models can hedge their bets in different ways: by including extra parameters, by assigning very wide prior distributions to the model parameters, or by using parameters that have a complicated functional form. By computing a weighted average likelihood across the entire parameter space, the marginal likelihood (and, consequently, the Bayes factor) automatically takes all these aspects into account.

Bayes factors represent “the standard Bayesian solution to the hypothesis testing and model selection problems” (Lewis & Raftery, 1997, p. 648) and “the primary tool used in Bayesian inference for hypothesis testing and model selection” (Berger, 2006, p. 378), but their application is not without challenges (Box 6.3). Below we show how these challenges can be overcome for the general class of MPT models. Next we compare the results of our Bayes factor analysis with those of the other model comparison methods using Jeffreys weights (i.e., normalized marginal likelihoods).

Application to Multinomial Processing Tree Models

In order to compute the Bayes factor, we seek to determine each model’s marginal likelihood $m(y | \mathcal{M}_{(\cdot)})$. In the following, we omit the conditioning on a particular model. As indicated above,

Bayes factors (Jeffreys, 1961; Kass & Raftery, 1995) come with two main challenges, one practical and one conceptual. The practical challenge arises because Bayes factors are defined as the ratio of two marginal likelihoods, each of which requires integration across the entire parameter space. This integration process can be cumbersome and hence the Bayes factor can be difficult to obtain. Fortunately, there are many approximate and exact methods to facilitate the computation of the Bayes factor (e.g., Ardia, Baştürk, Hoogerheide, & van Dijk, 2012; Chen, Shao, & Ibrahim, 2002; Gamerman & Lopes, 2006); in this chapter we focus on BIC (a crude approximation), the Savage-Dickey density ratio (applies only to nested models) and importance sampling. The conceptual challenge that Bayes factors bring is that the prior on the model parameters has a pronounced and lasting influence on the result. This should not come as a surprise: the Bayes factor punishes models for needless complexity, and the complexity of a model is determined in part by the prior distributions that are assigned to the parameters. The difficulty arises because researchers are often not very confident about the prior distributions that they specify. To overcome this challenge, one can either spend more time and effort on the specification of realistic priors, or else one can choose default priors that fulfill general desiderata (e.g., Jeffreys, 1961; Liang et al., 2008). Finally, the robustness of the conclusions can be verified by conducting a sensitivity analysis in which one examines the effect of changing the prior specification (e.g., Wagenmakers, Wetzels, Borsboom, & van der Maas, 2011).

Box 6.3 Two challenges for Bayes factors.

the marginal likelihood $m(y)$ is given by integrating the likelihood over the prior:

$$m(y) = \int p(y | \theta) p(\theta) d\theta. \quad (6.8)$$

The most straightforward manner to obtain $m(y)$ is to draw samples from the prior $p(\theta)$ and average the corresponding values for $p(y | \theta)$:

$$m(y) \approx \frac{1}{N} \sum_{i=1}^N p(y | \theta_i), \quad \theta_i \sim p(\theta). \quad (6.9)$$

For MPT models, this brute force integration approach may often be adequate. An MPT model usually has few parameters, and each is conveniently bounded from 0 to 1. However, brute force integration is inefficient, particularly when the posterior is highly peaked relative to the prior: in this case, draws from $p(\theta)$ tend to result in low likelihoods and only few chance draws may have high likelihood. This problem can be overcome by a numerical technique known as *importance sampling* (Hammersley & Handscomb, 1964).

In importance sampling, efficiency is increased by drawing samples from an importance density $g(\theta)$ instead of from the prior $p(\theta)$. Consider an importance density $g(\theta)$. Then,

$$\begin{aligned} m(y) &= \int p(y | \theta) p(\theta) \frac{g(\theta)}{g(\theta)} d\theta \\ &= \int \frac{p(y | \theta) p(\theta)}{g(\theta)} g(\theta) d\theta \\ &\approx \frac{1}{N} \sum_{i=1}^N \frac{p(y | \theta_i) p(\theta_i)}{g(\theta_i)}, \quad \theta_i \sim g(\theta). \end{aligned} \quad (6.10)$$

Note that if $g(\theta) = p(\theta)$, the importance sampler reduces to the brute force integration shown in Equation 6.9. Also note that if $g(\theta) = p(\theta | y)$, a single draw suffices to determine $p(y)$ exactly.

Importance sampling was invented by Stan Ulam and John von Neumann. Here we use it to estimate the marginal likelihood by repeatedly drawing samples and averaging—the samples are, however, not drawn from the prior (as per Equation 6.9, the brute force method), but instead they are drawn from some convenient density $g(\theta)$ (as per Equation 6.10; Andrieu, De Freitas, Doucet, & Jordan, 2003; Hammersley & Handscomb, 1964). The parameters in MPT models are constrained to the unit interval, and therefore the family of Beta distributions is a natural candidate for $g(\theta)$. The middle panel of Figure 6.4 shows an importance density (dashed line) for MPT parameter c in the no-conflict model for the data from Wagenaar and Boer (1987). This importance density is a Beta distribution that was fit to the posterior distribution for c using the method of moments. The importance density provides a good description of the posterior (the dashed line tracks the posterior almost perfectly) and therefore is more efficient than the brute force method illustrated in the left panel of Figure 6.4, which uses the prior as the importance density. Unfortunately, Beta distributions do not always fit MPT parameters so well; specifically, the Beta importance density may sometimes have tails that are thinner than the posterior, and this increases the variability of the marginal likelihood estimate. To increase robustness and ensure that the importance density has relatively fat tails, we can use a Beta mixture, shown in the right panel of Figure 6.4. The Beta mixture consists of a uniform prior component (i.e., the Beta(1, 1) prior as in the left panel) and a Beta posterior component (i.e., a Beta distribution fit to the posterior, as in the middle panel). In this example, the mixture weight for the uniform component is $w = 0.2$. Small mixture weights retain the efficiency of the Beta posterior approach but avoid the extra variability due to thin tails. It is possible to increase efficiency further by specifying a multivariate importance density, but the present univariate approach is intuitive, easy to implement, and appears to work well in practice. The accuracy of the estimate can be confirmed by increasing the number of draws from the importance density, and by varying the w parameter.

Box 6.4 Importance sampling for MPT models using the Beta mixture method.

In sum, when the importance density equals the prior we have brute force integration, and when it equals the posterior we have a zero-variance estimator. However, knowledge of the posterior implies knowledge of its normalizing constant (i.e., the marginal likelihood), and this is exactly the quantity we wish to determine. In practice then, we want to use an importance density that is similar to the posterior, is easy to evaluate, and is easy to draw samples from. In addition, we want to use an importance density with tails that are not thinner than those of the posterior; thin tails cause the estimate to have high variance. These desiderata are met by the *Beta mixture* importance density described in Box 6.4: a mixture between a Beta(1, 1) density and a Beta density that provides a close fit to the posterior distribution. Here we use a series of univariate Beta mixtures, one for each separate parameter, but acknowledge that a multivariate importance density is potentially even more efficient as it accommodates correlations between the parameters.

In our application to MPT models, we assume that all model parameters have uniform Beta(1, 1) priors. For most MPT models, this assumption is fairly uncontroversial. The uniform priors can be thought of as a default choice; in the presence of strong prior knowledge one can substitute more informative priors. The uniform priors yield a default Bayes factor that can be a reference point for an analysis with more informative priors, if such an analysis is desired.

Before turning to the results of the Bayes factor model comparison, we first inspect the posterior distributions. The posterior distributions were approximated using Markov chain Monte Carlo sampling implemented in JAGS (Plummer, 2003) and WinBUGS (Lunn et al., 2012).³ All code

³The second author used WinBUGS, the first and third authors used JAGS.

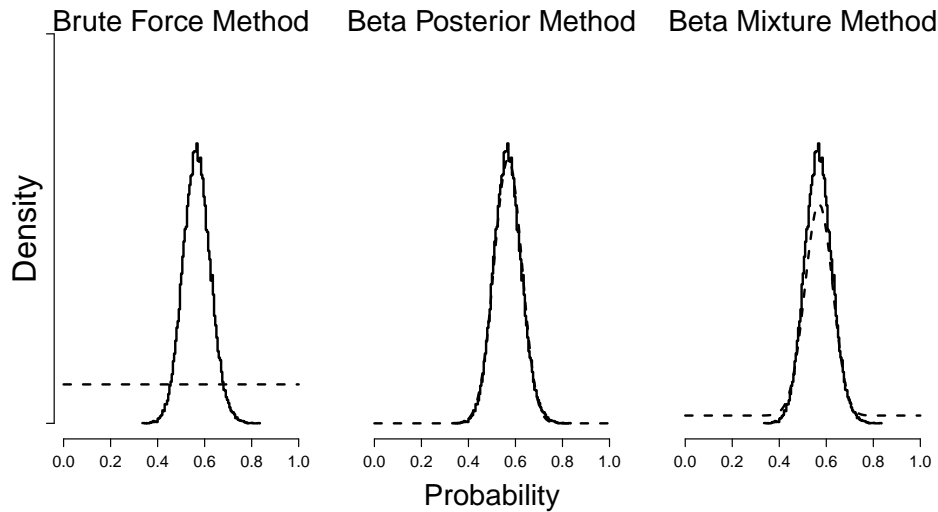


Figure 6.4 Three different importance sampling densities (dashed lines) for the posterior distribution (solid lines) of the c parameter in the no-conflict model as applied to the data from Wagenaar and Boer (1987). Left panel: a uniform Beta importance density (i.e., the brute force method); middle panel: a Beta posterior importance density (i.e., a Beta distribution that provides the best fit to the posterior); right panel: a Beta mixture importance density (i.e., a mixture of the uniform Beta density and the Beta posterior density, with a mixture weight $w = 0.2$ on the uniform component).

is available at <http://www.ejwagenmakers.com/papers.html>. Convergence was confirmed by visual inspection and the \hat{R} statistic (Gelman & Rubin, 1992). The top panel of Figure 6.5 shows the posterior distributions for the no-conflict model. Although there is slightly more certainty about parameter p than there is about parameters q and c , the posterior distributions for all three parameters are relatively wide considering that they are based on data from as many as 562 participants.

The middle panel of Figure 6.5 shows the posterior distributions for the destructive-updating model. It is important to realize that when $d = 0$ (i.e., no destruction of the earlier memory), the destructive-updating model reduces to the no-conflict model. Compared to the no-conflict model, parameters p , q , and c show relatively little change. The posterior distribution for d is very wide, indicating considerable uncertainty about its true value. A frequentist point-estimate yields $\hat{d} = 0$ (Wagenaar & Boer, 1987; see also Table 6.1), but this does not convey the fact that this estimate is highly uncertain.

The lower panel of Figure 6.5 shows the posterior distributions for the coexistence model. When $s = 0$ (i.e., no suppression of the earlier memory), the coexistence model reduces to the no-conflict model. Compared to the no-conflict model and the destructive-updating model, parameters p , q , and c again show relatively little change. The posterior distribution for s is very wide, indicating considerable uncertainty about its true value.

The fact that the no-conflict model is nested under both the destructive-updating model and the no-conflict model allows us to inspect the extra parameters d and s and conclude that we have not learned very much about their true values. This suggests that, despite having tested 562 participants, the data do not firmly support one model over the other. We will now see how Bayes

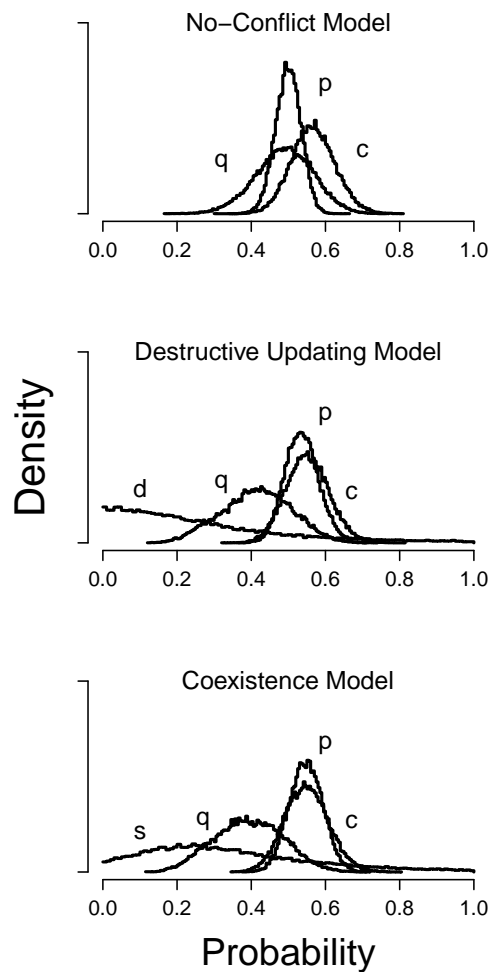


Figure 6.5 Posterior distributions for the parameters of the no-conflict MPT model, the destructive updating MPT model, and the coexistence MPT model, as applied to the data from Wagenaar and Boer (1987).

factors can make this intuitive judgment more precise.

We applied the Beta mixture importance sampling method to estimate marginal likelihoods for the three models under consideration. The results were confirmed by varying the mixture weight w , by independent implementations by the authors, and by comparison to the Savage-Dickey density ratio test presented later. Table 6.5 shows the results.

From the marginal likelihoods and the Jeffreys weights, we can derive the Bayes factors for the pair-wise comparisons; the Bayes factor is 2.77 in favor of the no-conflict model over the destructive updating model, the Bayes factor is 1.39 in favor of the coexistence model over the no-conflict model, and the Bayes factor is 3.86 in favor of the coexistence model over the destructive updating model. The first two of these Bayes factors are anecdotal or “not worth more than a bare mention” (Jeffreys, 1961), and the third one just makes the criterion for “moderate” evidence, although any enthusiasm about this level of evidence should be tempered by the realization that Jeffreys himself described a Bayes factor as high as 5.33 as “odds that would interest a gambler, but would be hardly worth more than a passing mention in a scientific paper” (Jeffreys, 1961, pp.

Table 6.5 Bayesian Evidence, Jeffreys Weights, and Pairwise Bayes Factors Computed From the Jeffreys Weights or Through the Savage-Dickey Density Ratio for the Wagenaar and Boer (1987) MPT Models.

	Bayesian evidence	Jeffreys weight	Bayes factor (Savage-Dickey)		
			over NCM	over DUM	over CXM
No-conflict model (NCM)	-30.55	0.36	1	2.77 (2.81)	0.72 (0.80)
Destructive updating model (DUM)	-31.57	0.13	0.36 (0.36)	1	0.26 (0.28*)
Coexistence model (CXM)	-30.22	0.51	1.39 (1.25)	3.86 (3.51*)	1

Note. * Derived through transitivity: $2.81 \times 1/0.80 = 3.51$.

256-257). In other words, the Bayes factors are consistent with the intuitive visual assessment of the posterior distributions: the data do not allow us to draw strong conclusions.

We should stress that Bayes factors apply to a comparison of any two models, regardless of whether or not they are structurally related or *nested*, so that one model is a special, simplified version of a larger, encompassing model. As is true for the information criteria and minimum description length methods, Bayes factors can be used to compare structurally very different models, such as for example REM (Shiffrin & Steyvers, 1997) versus ACT-R (Anderson et al., 2004), or the diffusion model (Ratcliff, 1978) versus the linear ballistic accumulator model (Brown & Heathcote, 2008). In other words, Bayes factors can be applied to nested and non-nested models alike. For the models under consideration, however, there exists a nested structure that allows one to obtain the Bayes factor through a computational shortcut.

Specifically, consider first the comparison between the no-conflict model \mathcal{M}_{NCM} and the destructive updating model \mathcal{M}_{DUM} . As shown above, we can obtain the Bayes factor for \mathcal{M}_{NCM} versus \mathcal{M}_{DUM} by computing the marginal likelihoods using importance sampling. However, because the models are nested we can also obtain the Bayes factor by considering only \mathcal{M}_{DUM} , and dividing the posterior ordinate at $d = 0$ by the prior ordinate at $d = 0$. This surprising result was first published by Dickey and Lientz (1970), who attributed it to Leonard J. “Jimmie” Savage. The result is now generally known as the *Savage-Dickey density ratio* (e.g., Dickey, 1971; for extensions and generalizations, see Chen, 2005; Verdinelli & Wasserman, 1995; Wetzels, Grasman, & Wagenmakers, 2010; for an introduction for psychologists, see Wagenmakers et al., 2010; a short mathematical proof is presented in O’Hagan & Forster, 2004, pp. 174-177). Thus, we can exploit the fact that \mathcal{M}_{NCM} is nested in \mathcal{M}_{DUM} and use the Savage-Dickey density ratio to obtain the Bayes factor:

$$BF_{\text{NCM,DUM}} = \frac{m(y | \mathcal{M}_{\text{NCM}})}{m(y | \mathcal{M}_{\text{DUM}})} = \frac{p(d = 0 | y, \mathcal{M}_{\text{DUM}})}{p(d = 0 | \mathcal{M}_{\text{DUM}})}. \quad (6.11)$$

The Savage-Dickey density ratio test is visualized in Figure 6.6; the posterior ordinate at $d = 0$ is higher than the prior ordinate at $d = 0$, indicating that the data have increased the plausibility that d equals 0. This means that the data support \mathcal{M}_{NCM} over \mathcal{M}_{DUM} . The prior ordinate equals 1, and hence $BF_{\text{NCM,DUM}}$ simply equals the posterior ordinate at $d = 0$. A nonparametric density estimator (Stone, Hansen, Kooperberg, & Truong, 1997) that respects the bound at 0 yields an estimate of 2.81. This estimate is close to 2.77, the estimate from the importance sampling approach.

The Savage-Dickey density ratio test can be applied similarly to the comparison between the no-conflict model \mathcal{M}_{NCM} versus the coexistence model \mathcal{M}_{CXM} , where the critical test is at $s = 0$. Here the posterior ordinate is estimated to be 0.80, and hence the Bayes factor for \mathcal{M}_{CXM} over

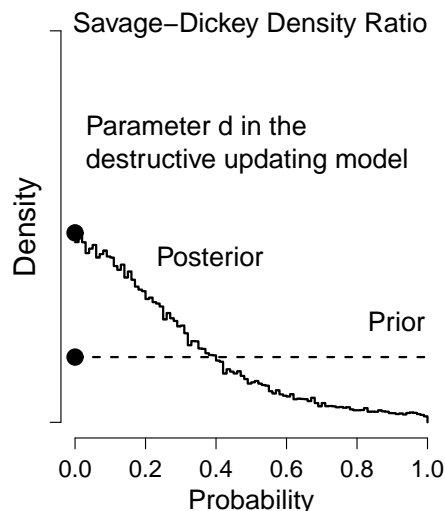


Figure 6.6 Illustration of the Savage-Dickey density ratio test. The dashed and solid lines show the prior and the posterior distribution for parameter d in the destructive updating model. The black dots indicate the height of the prior and the posterior distributions at $d = 0$.

\mathcal{M}_{NCM} equals $1/0.80 = 1.25$, close to the Bayes factor obtained through importance sampling, $BF_{\text{CXM,NCM}} = 1.39$.

With these two Bayes factors in hand, we can immediately derive the Bayes factor for the comparison between the destructive updating model \mathcal{M}_{DUM} versus the coexistence model \mathcal{M}_{CXM} through transitivity, that is, $BF_{\text{CXM,DUM}} = BF_{\text{NCM,DUM}} \times BF_{\text{CXM,NCM}}$. Alternatively, we can also obtain $BF_{\text{CXM,DUM}}$ by directly comparing the posterior density for $d = 0$ against that for $s = 0$:

$$\begin{aligned}
 BF_{\text{CXM,DUM}} &= BF_{\text{NCM,DUM}} \times BF_{\text{CXM,NCM}} \\
 &= \frac{p(d = 0 \mid y, \mathcal{M}_{\text{DUM}})}{p(d = 0 \mid \mathcal{M}_{\text{DUM}})} \times \frac{p(s = 0 \mid \mathcal{M}_{\text{CXM}})}{p(s = 0 \mid y, \mathcal{M}_{\text{CXM}})} \\
 &= \frac{p(d = 0 \mid y, \mathcal{M}_{\text{DUM}})}{p(s = 0 \mid y, \mathcal{M}_{\text{CXM}})},
 \end{aligned} \tag{6.12}$$

where the second step is allowed because we have assigned uniform priors to both d and s , so that $p(d = 0 \mid \mathcal{M}_{\text{DUM}}) = p(s = 0 \mid \mathcal{M}_{\text{CXM}})$. Hence, the Savage-Dickey estimate for the Bayes factor between the two non-nested models \mathcal{M}_{DUM} and \mathcal{M}_{CXM} equals the ratio of the posterior ordinates at $d = 0$ and $s = 0$, resulting in the estimate $BF_{\text{CXM,DUM}} = 3.51$, close to the importance sampling result of 3.86.

Comparison of Model Comparisons

We have now implemented and performed a variety of model comparison methods for the three competing MPT models introduced by Wagenaar and Boer (1987): we computed and interpreted the Akaike information criteria (AIC), Bayesian information criteria (BIC), the Fisher information approximation of the minimum description length principle (FIA), and two computational implementations of the Bayes factor (BF).

The general tenor across most of the model comparison exercises has been that the data do not convincingly support one particular model. However, the destructive updating model is consistently ranked the worst of the set. Looking at the parameter estimates, it is not difficult to see why this is so: the d parameter of the destructive updating model (i.e., the probability of destroying memory through updating) is estimated at 0, thereby reducing the destructive updating model to the no-conflict model, and yielding an identical fit to the data (as can be seen in the likelihood column of Table 6.2). Since the no-conflict model counts as a special case of the destructive updating model, it is by necessity less complex from a model selection point of view—the d parameter is an unnecessary entity, the inclusion of which is not warranted by the data. This is reflected in the inferior performance of the destructive updating model according to all measures of generalizability.

The difference between the no-conflict model and the coexistence model is less clear-cut. Following AIC, the two models are virtually indistinguishable—compared to the coexistence model, the no-conflict model sacrifices one unit of log-likelihood for two units of complexity (one parameter). As a result, both models perform equally well under the AIC measure. Under the BIC measure, however, the penalty for the number of free parameters is more substantial, and here the no-conflict model trades a unit of log likelihood for $\log(N) = 6.33$ units of complexity, outdistancing both the destructive updating model and the coexistence model. The BIC is the exception in clearly preferring the no-conflict model over the coexistence model. The MDL, like the AIC, would have us hedge on the discriminability of the no-conflict model and the coexistence model.

The BF, finally, allows us to express evidence for the models using standard probability theory. Between any two models, the BF tells us how much the balance of evidence has shifted due to the data. Using two methods of computing the BF, we determined that the odds of the coexistence model over the destructive updating model almost quadrupled ($BF_{\text{CXM,DUM}} \approx 3.86$), but the odds of the coexistence model over the no-conflict model barely shifted at all ($BF_{\text{CXM,NCM}} \approx 1.39$). Finally, we can use the same principles of probability to compute Jeffreys weights, which express, for each model under consideration, the probability that it is true, assuming prior indifference. Furthermore, we can easily recompute the probabilities in case we wish to express a prior preference between the candidate models (for example, we might use the prior to express a preference for sparsity, as was originally proposed by Jeffreys, 1961).

6.5 Concluding Comments

Model comparison methods need to implement the principle of parsimony: goodness-of-fit has to be discounted to the extent that it was accomplished by a model that is overly complex. Many methods of model comparison exist (Myung et al., 2000; Wagenmakers & Waldorp, 2006), and our selective review focused on methods that are popular, easy-to-compute approximations (i.e., AIC and BIC) and methods that are difficult-to-compute “ideal” solutions (i.e., minimum description length and Bayes factors). We applied these model comparison methods to the scenario of three competing MPT models introduced by Wagenaar and Boer (1987). Despite collecting data from 562 participants, the model comparison methods indicate that the data are somewhat ambiguous; at any rate, the data do not support the destructive updating model. This echoes the conclusions drawn by Wagenaar and Boer (1987).

It is important to note that the model comparison methods discussed in this chapter can be applied regardless of whether the models are nested. This is not just a practical nicety; it also means that the methods are based on principles that transcend the details of a specific model implementation. In our opinion, a method of inference that is necessarily limited to the comparison of nested models is incomplete at best and misleading at worst. It is also important to realize that

model comparison methods are *relative* indices of model adequacy; when, say, the Bayes factor expresses an extreme preference for model A over model B, this does not mean that model A fits the data at all well. Because it is a mistake to base inference on a model that fails to describe the data, a complete inference methodology features both relative and absolute indices of model adequacy. For the MPT models under consideration here, Wagenaar and Boer (1987) reported that the no-conflict model provided “an almost perfect fit” to the data.⁴

The example MPT scenario considered here was relatively straightforward. More complicated MPT models contain order-restrictions, feature individual differences embedded in a hierarchical framework (Klauer, 2010; Matzke et al., in press), or contain a mixture-model representation with different latent classes of participants (for application to other models, see Frühwirth-Schnatter, 2006; Scheibehenne, Rieskamp, & Wagenmakers, 2013). In theory, it is relatively easy to derive Bayes factors for these more complicated models. In practice, however, Bayes factors for complicated models may require the use of numerical techniques more involved than importance sampling. Nevertheless, for standard MPT models, the Beta mixture importance sampler appears to be a convenient and reliable tool to obtain Bayes factors. We hope that this methodology will facilitate the principled comparison of MPT models in future applications.

⁴We confirmed the high quality of fit in a Bayesian framework using posterior predictives (Gelman & Hill, 2007), results not reported here.