

1

A Tutorial on Fisher Information

2

Alexander Ly, Josine Verhagen, Raoul Grasman, and Eric-Jan

3

Wagenmakers

4

University of Amsterdam

5

Abstract

6

The concept of Fisher information plays a crucial role in many statistical applications that are of key importance to mathematical psychologists. Here we explain Fisher information by illustrating its use across three different statistical paradigms: first, in the frequentist paradigm, Fisher information is used to determine the sample size with which we design an experiment; second, in the Bayesian paradigm, Fisher information is used to define a default parameter prior; finally, in the minimum description length paradigm, Fisher information is used to measure model complexity. Each application is illustrated with simple examples.

Keywords: Statistical inference; Jeffreys' prior; Model complexity; Model selection; Minimum description length

7

8

9

10

Mathematical psychologists develop and apply statistical models in order to describe human behavior and understand latent psychological processes. Examples include Stevens' law of psychophysics that describes the relation between the objective physical intensity of a stimulus and its subjectively experienced intensity (Stevens, 1957); Ratcliff's diffusion model

This work was supported by the starting grant "Bayes or Bust" awarded by the European Research Council. Correspondence concerning this article may be addressed to Alexander Ly, University of Amsterdam, Department of Psychology, Weesperplein 4, 1018 XA Amsterdam, the Netherlands. Email address: alexander.ly.nl@gmail.com.

1 of decision making that quantifies the different processes that drive behavior in speeded
2 response time tasks (Ratcliff, 1978); multinomial processing tree models that decompose
3 performance in memory tasks into the contribution of separate latent mechanisms (Batchelder
4 & Riefer, 1980; Chechile, 1973), and so on and so forth.

5 When applying their models to data, mathematical psychologists may operate from
6 within different statistical paradigms and focus on different substantive questions. For in-
7 stance, working within the classical or frequentist paradigm a researcher may wish to decide
8 upon the number of trials to be presented to a participant in order to estimate the partic-
9 ipant's latent abilities. Working within the Bayesian paradigm a researchers may wish to
10 know how to determine a suitable default prior on the model parameters. Working within the
11 minimum description length (MDL) paradigm a researcher may wish to compare rival models
12 and quantify their complexity. Despite the diversity of these paradigms and purposes, they
13 are connected through the concept of Fisher information.

14 Fisher information plays a pivotal role throughout statistical modeling, but an accessi-
15 ble introduction for mathematical psychologists is lacking. The goal of this tutorial is to fill
16 this gap and illustrate the use of Fisher information in the three statistical paradigms men-
17 tioned above: frequentist, Bayesian, and MDL. This work builds directly upon the *Journal*
18 *of Mathematical Psychology* tutorial article by Myung (2003) on maximum likelihood estima-
19 tion. The intended target group for this tutorial are graduate students and researchers with
20 an affinity for cognitive modeling and mathematical statistics.

21 The outline of this article is as follows: To keep this tutorial self-contained, we start
22 by establishing the standard notations for probabilistic modeling and introduce two running
23 examples featuring the Bernoulli distribution and the normal distribution. The next section
24 provides the definition of Fisher information and shows how it can be calculated in general.
25 The next three sections exemplify the use of Fisher information for different purposes. Sec-
26 tion 1 shows how Fisher information can be used to determine the desired sample size for
27 an experiment; Section 2 shows how Fisher information can be used to define a default prior
28 on model parameters, and Section 3 shows how Fisher information can be used to measure
29 model complexity.

1

Basic Concepts

2

3

4

5

6

7

For concreteness, consider an experiment comprised of n trials of equal difficulty. The goal of the experiment is to determine a participant's ability to identify particular species of animals. On each trial, the participant is presented with a picture of an animal on a computer screen and is asked to identify the species. On each trial, the participant can either respond correctly or in error. We assume that if the participant does not recognize the species, the probability of guessing the correct answer is close to zero.

8

9

10

11

12

13

14

15

16

17

18

Notation Let n denote the planned number of trials in an experiment and let $\vec{X} = (X_1, \dots, X_n)$ denote a *future* trial sequence that will be presented to a participant. More generally, we call \vec{X} a random variable as we are uncertain about its outcomes. In a typical experimental set up, we design each trial X_i to be a replication of a prototypical trial, which we denote by X . The prototypical trial in our running example can take on two outcomes x , and we write $x = 1$ when the participant responded correctly, and $x = 0$ otherwise. Similarly, we denote $\vec{x} = (x_1, x_2, \dots, x_n)$ for an outcome of the n -trial sequence \vec{X} . For example, after a participant completed the experiment with $n = 10$ trials, one possible outcome is $\vec{x} = (1, 0, 0, 0, 0, 0, 1, 0, 0)$, implying that the participant responded correctly only on the first and the eighth trial. In order to model the chance with which this outcome occurs we construct a probabilistic model: the Bernoulli model.

19

20

21

22

Example 1 (The Bernoulli model). *To model the behavior of the participant, we first construct a model for all possible outcomes of the prototypical trial X , assuming all such outcomes are caused by the participant's general "species recognition ability", which we denote by θ . This latent ability θ relates to the possible outcomes as follows:*

$$f(x | \theta) = P(X = x) = \theta^x(1 - \theta)^{1-x}, \text{ where } x = 0 \text{ or } x = 1. \quad (1)$$

23

We call this relation a probability density function (pdf).¹ The pdf allows us to deduce the

¹Formally, Eq. (1) defines a probability mass function, as X takes on discrete values; for brevity we do not distinguish between the two here.

1 chances with which an outcome occurs when the latent ability is known. Hence, if $\theta = 1$ the
 2 participant always correctly recognizes the animal species, as the chance of a correct response
 3 is 1. Conversely, if $\theta = 0$ the participant never correctly recognizes the animal species, as
 4 the chance of a correct response is 0. In reality, a participant's latent ability is likely to fall
 5 somewhere in between these two extremes, $\theta \in (0, 1)$. The relation Eq. (1) describes how a
 6 participant's true latent ability θ "causes" the behavioral outcomes \vec{x} . This particular relation
 7 is known as the Bernoulli model, denoted $X \sim \text{Ber}(\theta)$. As the full trial sequence \vec{X} consists
 8 of n independent and identically distributed (i.i.d.) replications of X we can extrapolate a pdf
 9 for the vector-valued outcomes \vec{x} by taking n products of Eq. (1):

$$f(\vec{x} | \theta) = \underbrace{\theta^{x_1} (1 - \theta)^{1-x_1}}_{X_1 \sim \text{Ber}(\theta)} \cdot \dots \cdot \underbrace{\theta^{x_n} (1 - \theta)^{1-x_n}}_{X_n \sim \text{Ber}(\theta)} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}. \quad (2)$$

10 Note that with $n = 10$ we have $2^{10} = 1024$ possible outcomes \vec{x} for the experiment \vec{X} .

11 Typically, we do not report the raw outcome series \vec{x} but only the number of successful
 12 responses y . We can view y as an outcome of the summary statistic $Y = \sum_{i=1}^n X_i$. A statistic
 13 is by definition a function of possible outcomes and it therefore inherits the random mechanism
 14 of \vec{x} , in this case a participant's latent ability θ . The pdf that Y inherits from Eq. (2) is called
 15 a binomial distribution $Y \sim \text{Bin}(\theta, n = 10)$. For example, when $\theta = 0.4$ the chances of the
 16 outcomes are then given by

$$P(Y = y) = f(y | n = 10, \theta = 0.4) = \frac{10!}{y!(10 - y)!} (0.4)^y (0.6)^{10-y} \text{ for } y = 0, 1, \dots, 10. \quad (3)$$

17 Note that Y can only take on 11 values compared to the $2^{10} = 1024$ outcomes of \vec{X} , a
 18 compression rate of almost a hundred, due to the fact that Y ignores the order within \vec{x} ,
 19 compare Eq. (2) with Eq. (4). Fig. 1 shows the chances corresponding to the outcomes of Y
 20 with $\theta = 0.4$.

21 Fig. 1 emphasizes the nature of probabilistic modeling as no possible outcome is entirely

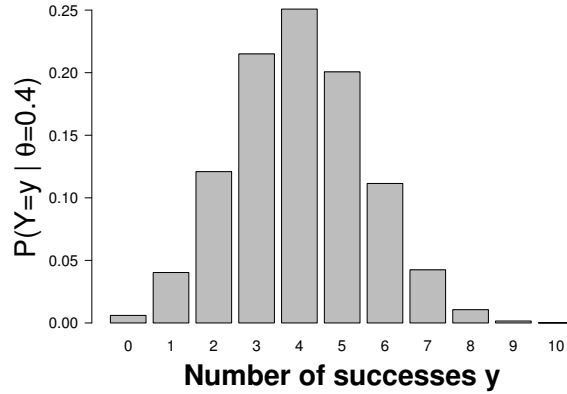


Figure 1. The chances of observing a particular number of successes in a binomial experiment with $n = 10$ and $\theta = 0.4$. The area under the bars sum to one. Note that, when a participant's latent ability is truly $\theta = 0.4$ an outcome of $Y = 0$ is still possible with chance $P(Y = 0 | \theta = 0.4) = 0.006$.

- 1 *excluded. In particular, even when a participant's latent ability is truly $\theta = 0.4$ there is*
 2 *a slim chance, $P(Y = 0 | \theta = 0.4) = 0.006$, that she will respond incorrectly on all trials.*
 3 *More generally, any number of trials n and latent ability θ between zero and one would yield*
 4 *outcomes of Y with the following probabilities:*

$$f(y | n, \theta) = \frac{n!}{y!(n-y)!} \theta^y (1-\theta)^{n-y} \text{ for } y = 0, 1, \dots, n, \text{ for all } \theta \in (0, 1). \quad (4)$$

- 5 *Equivalently, we say that the number of successes is modeled according to a binomial model,*
 6 *which we denote by $Y \sim \text{Bin}(\theta, n)$. Moreover, this interpretation of Eq. (4) is in line with the*
 7 *idea that any outcome y is generated from the latent ability θ ; schematically, $\theta \xrightarrow{\text{Bin}(\theta, n)} y$. \diamond*

- 8 **Likelihood** In practice, we only observe data and do not know θ . To infer something about
 9 a participant's latent ability θ , we invert the data generating process by exchanging the roles
 10 of y and θ within a pdf. For example, when a participant responded correctly on seven out
 11 of ten trials, we plug this information into Eq. (4) and consider it as a function of possible
 12 latent abilities θ :

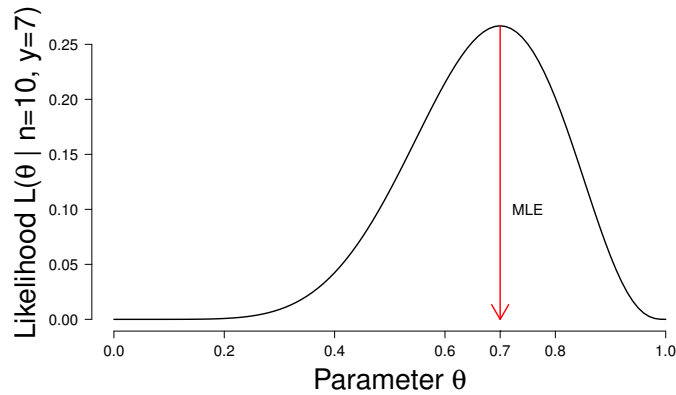


Figure 2. The likelihood function based on observing seven $y = 7$ out of ten $n = 10$ correct responses in a binomial experiment. The MLE equals $y/n = 0.7$. Note that the likelihood function is a continuous function of the parameter θ , whereas the observed data are discrete.

$$L(\theta | n = 10, y = 7) = \frac{10!}{7!3!} \theta^7 (1 - \theta)^3. \quad (5)$$

1 This function is shown in Fig. 2. To distinguish Eq. (5) from a pdf, we call it a likelihood
 2 function. Fisher used the likelihood function to develop his general purpose method of maxi-
 3 mum likelihood estimation (Fisher, 1912; Fisher, 1922; Fisher, 1925; Myung, 2003). The idea
 4 is to use the modeled relation $f(\vec{x} | \theta)$ to “reverse engineer” the most likely value of θ that
 5 could have been responsible for the observed data. It can be easily shown that the maxi-
 6 mum likelihood estimator (i.e., MLE) for the Bernoulli model is given by the sample mean
 7 $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$.

8 **Example 2** (The n -trial Normal Model). *Pre-experimental:* As a second example we
 9 consider the normal model, which forms the building block for common models in psychology
 10 (e.g., t -tests, ANOVA, and linear regression). In this model, the trial sequence \vec{X} consists of
 11 n replications of a prototypical random variable X with a pdf given by

$$f(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x - \mu^2}{2\sigma^2}\right), \quad (6)$$

1 where the parameter vector consists of the population mean and the population variance $\vec{\theta} =$
 2 (μ, σ^2) . To extrapolate the pdf for the n -trial experiment \vec{X} we exploit the i.i.d. assumption
 3 by taking n products of Eq. (6), which yields:

$$f(\vec{x} | \vec{\theta}) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right). \quad (7)$$

4 We then say that the possible outcomes follow a normal distribution, $X \sim \mathcal{N}(\mu, \sigma^2)$.²

5 **Inference:** To invert the generating model we simply plug the observed data \vec{x} the
 6 relation into Eq. (7) and the resulting likelihood function then consists of a two-dimensional
 7 surface spanned by both μ and σ^2 . Maximizing the likelihood is equivalent to maximizing the
 8 natural logarithm of the likelihood function, i.e., $\log f(\vec{x} | \vec{\theta})$, a function shown in Fig. 3 for
 9 fictitious data.

10 For the normal model, the MLE consists of a pair $\hat{\mu}, \hat{\sigma}^2$: the sample mean $\hat{\mu} = \bar{X} = \frac{1}{n} \sum X_i$
 11 and $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$. See the online appendix for a derivation. \diamond

12 In Section 1 we will show how Fisher information is related to the performance of
 13 the aforementioned maximum likelihood estimators. First, however, we introduce the general
 14 definition of Fisher information and show how it can be computed for the two models discussed
 15 above.

16 Fisher Information

17 **Definition of Fisher Information** The (unit) Fisher information is a measure for the
 18 amount of information that is expected within the prototypical trial X about the parameter
 19 of interest θ . It is defined as the variance of the so-called score function, i.e., the derivative
 20 of the log-likelihood function with respect to the parameter,

²This is a slight abuse of notation; however, we believe the meaning is well understood, as we assume that the random variable $\vec{X} = X_1, \dots, X_n$ consists of n i.i.d. copies of a prototypical trial X .

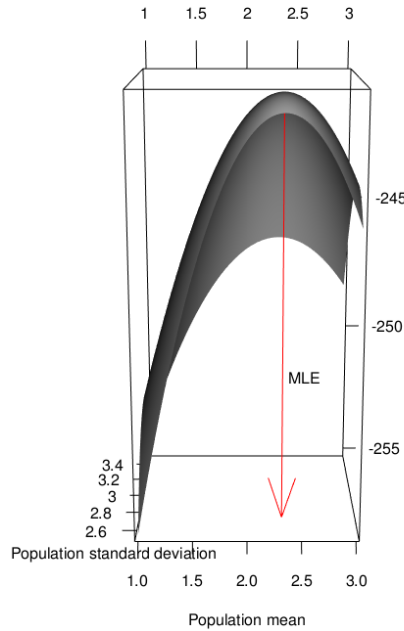


Figure 3. Log-likelihood function for 100 simulated normally distributed observations \vec{x} . The log-likelihood is a surface spanned by the two parameters μ and σ^2 . On the x -axis is the population mean, and on the y -axis is the population standard deviation. The MLE $(\hat{\mu}, \hat{\sigma}) = (2.32, 2.68)$ corresponds to the projection of the highest point on the log-likelihood surface onto the parameter space.

$$I(\theta) = \text{Var}\left(\underbrace{\frac{d}{d\theta} \log f(X|\theta)}_{\text{Score function}}\right) \text{ or as } I(\theta) = -E\left(\frac{d^2}{d\theta^2} \log f(X|\theta)\right), \quad (8)$$

1 under additional (mild) regularity conditions. To calculate $I(\theta)$ we keep θ fixed and take the
 2 expectation with respect to all possible outcomes x :

$$I(\theta) = -E\left(\frac{d^2}{d\theta^2} \log f(X|\theta)\right) = -\int_{\mathcal{X}} \left(\frac{d^2}{d\theta^2} \log f(x|\theta)\right) f(x|\theta) dx, \quad (9)$$

3 where \mathcal{X} denotes the outcome space of X . When X is discrete, we replace the integral by a
 4 summation, see Eq. (10) below.

5 **Example 1** (Fisher Information in the Bernoulli Model, Continued). *To calculate the infor-*
 6 *mation of a Bernoulli distribution, we take Eq. (1) and plug this into Eq. (9) which yields:*

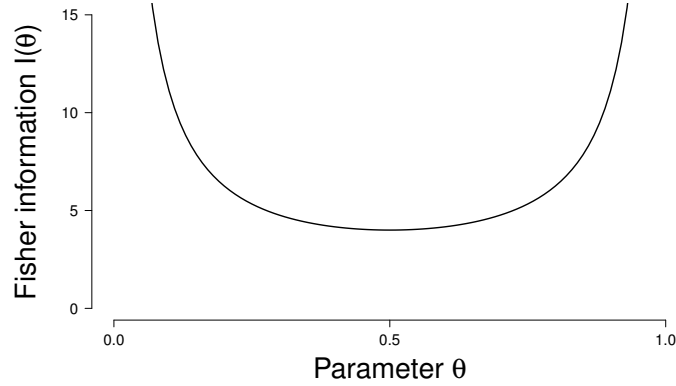


Figure 4. The Fisher information as a function of θ within the Bernoulli model. As θ reaches zero or one the expected information goes to infinity. When $\theta = 1$ the outcomes will always be 1, therefore clearly conveying this information within the data.

$$I(\theta) = - \sum_{x=0}^1 \frac{d^2}{d\theta^2} \log P(X=x)P(X=x) = - \left(-\frac{1}{\theta^2}\theta - \frac{1}{(1-\theta)^2}(1-\theta) \right) = \frac{1}{\theta(1-\theta)}. \quad (10)$$

1 Hence, when a participant's true latent ability is θ^* the expected information about her ability
 2 is then given by $I(\theta^*) = \frac{1}{\theta^*(1-\theta^*)}$ for a prototypical trial, for instance, $I(0.5) = 4$ when
 3 $\theta^* = 0.5$. Thus, the Fisher information contains two aspects of the model: (1) the sensitivity
 4 of the relationship $f(x|\theta)$ with respect to the parameter θ expressed by the score function at
 5 the true value θ^* , and (2) how this sensitivity at θ^* varies over (all possible) outcomes x a
 6 model can generate according to $f(x|\theta^*)$. This dependence on the true parameter value θ^* is
 7 shown in Fig. 4.

8 As the full trial sequence \vec{X} consists of n i.i.d. replications of the prototypical X we
 9 formally refer to Eq. (8) as the unit Fisher information. To extrapolate the information about
 10 θ within n trials we simply multiply the unit information by n , $I_n(\theta) = nI(\theta)$. For the 10-
 11 trial Bernoulli model, this means $I_n(\theta) = \frac{10}{\theta(1-\theta)}$. This connection can be formally shown by
 12 exploiting the i.i.d. assumption and by computing the integral Eq. (8) with the pdf Eq. (2)
 13 instead.

14 Furthermore, we can also use Eq. (8) to calculate the Fisher information about θ con-

1 tained in the summary statistic Y for which we also get $I_Y(\theta) = \frac{10}{\theta(1-\theta)} = I_n(\theta)$. This means
 2 that no information is lost when we infer θ from the summary statistic Y rather than from
 3 the full trial sequence \vec{X} . ◇

4 **Unit Fisher Information Within an Observed Sample** The expected (unit) Fisher
 5 information as defined in Eq. (8) is an weighted average over all the possible outcomes of X .
 6 We might, however, be interested in the observed (unit) Fisher information within an observed
 7 sample \vec{x} instead. We then replace the expectation in Eq. (9) by its empirical version

$$I_{\text{Obs}}(\theta) = -\frac{1}{n} \sum_{i=1}^n \frac{d^2}{d\theta^2} \log f(x_i | \theta). \tag{11}$$

8 For example, if we observed 7 successes out of 10 Bernoulli trials we have:

$$I_{\text{Obs}}(\theta) = \frac{1}{10} \sum_{i=1}^{10} \frac{x_i}{\theta^2} + \frac{n - x_i}{(1 - \theta)^2} = \frac{1}{10} \left(\frac{7}{\theta^2} + \frac{3}{(1 - \theta)^2} \right). \tag{12}$$

9 If the data are truly governed by $\theta^* = 0.7$, we would indeed expect to see $10\theta^* = 7$ successes
 10 and $10(1 - \theta^*) = 3$ failures and the observed and expected Fisher information then coincide:

$$I_{\text{Obs}}(\theta^*) = \frac{1}{10} \left(\frac{10\theta^*}{(\theta^*)^2} + \frac{10(1 - \theta^*)}{(1 - \theta^*)^2} \right) = \frac{10}{10} \left(\frac{\theta^*}{(\theta^*)^2} + \frac{(1 - \theta^*)}{(1 - \theta^*)^2} \right) = \frac{1}{\theta^*(1 - \theta^*)} = I(\theta^*). \tag{13}$$

11 On the other hand, if $\theta^* = 0.15$ the probability of seeing 7 out of 10 heads is about .01 and we
 12 then see that there is a big discrepancy between expected and observed Fisher information,
 13 $I(0.15) \approx 8$ versus $I_{\text{Obs}}(0.15) = 31.5$ respectively. This might imply that the hypothesis
 14 $\theta^* = 0.15$ was wrong to begin with.

15 In more realistic cases, we do not know θ^* and to calculate the observed Fisher infor-
 16 mation we replace θ^* by the MLE, which yields $I_{\text{Obs}}(\hat{\theta})$. For the examples we discuss this
 17 coincides with plugging in the MLE into the expected Fisher information, i.e., $I(\hat{\theta}) = I_{\text{Obs}}(\hat{\theta})$.³

³This equality holds for distributions that belong to the so-called exponential family, see Huzurbazar (1949).

1 **Example 2** (Fisher Information in the Normal Model, Continued). *When there are multiple*
 2 *parameters of interest, say $\vec{\theta} = (\theta_1, \dots, \theta_d)$, the Fisher information turns into a $d \times d$ matrix.*
 3 *The i, j -th element of this matrix is then given by*

$$I(\vec{\theta})_{i,j} = \text{Cov} \left(\frac{\partial}{\partial \theta_i} \log f(X | \vec{\theta}), \frac{\partial}{\partial \theta_j} \log f(X | \vec{\theta}) \right) \text{ or } I(\vec{\theta})_{i,j} = -E \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(X | \vec{\theta}) \right), \quad (14)$$

4 *under additional (mild) regularity conditions. For example, when X is normally distributed*
 5 *with both $\vec{\theta} = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$ unknown, the unit Fisher information matrix $I(\vec{\theta})$ turns into a 2×2 matrix,*
 6 *consisting of expectations of partial derivatives:*

$$I(\vec{\theta}) = -E \begin{pmatrix} \frac{\partial^2}{\partial \mu \partial \mu} \log f(x | \mu, \sigma^2) & \frac{\partial^2}{\partial \mu \partial \sigma^2} \log f(x | \mu, \sigma^2) \\ \frac{\partial^2}{\partial \sigma^2 \partial \mu} \log f(x | \mu, \sigma^2) & \frac{\partial^2}{\partial \sigma^2 \partial \sigma^2} \log f(x | \mu, \sigma^2) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}, \quad (15)$$

7 *see the online appendix for a derivation. The off-diagonal elements are in general not zero,*
 8 *unless the pdfs are symmetric around its location μ . When the data are modeled as n i.i.d.*
 9 *normal distributions, the expected Fisher information within \vec{X} is then given by $I_n(\vec{\theta}) = nI(\vec{\theta})$.*
 10 *As in the Bernoulli model, the observed (unit) Fisher information can be calculated by plugging*
 11 *the MLE $\hat{\sigma}$ into its expected version Eq. (15). \diamond*

12 In this section we showed how to calculate the Fisher information, a measure for the
 13 expected amount of information within the trial sequence \vec{X} about the parameter of interest
 14 θ . We now highlight the role of Fisher information in three different statistical paradigms: in
 15 the frequentist paradigm, Fisher information is used to determine the sample size with which
 16 we design an experiment; in the Bayesian paradigm, Fisher information is used to determine
 17 a default prior; and in the minimum description length paradigm, Fisher information is used
 18 to measure model complexity.

1 1. Using Fisher Information to Determine Sample Size

2 **Introduction** As mentioned above, a statistical model for a psychological process can be
 3 thought of as a particular mapping f between an individual's latent properties θ and that
 4 individual's possible behavior, i.e., the outcomes \vec{x} , schematically: $\theta \xrightarrow{f(\vec{x}|\theta)} \vec{x}$. Consequently,
 5 when we know an individual's latent ability θ , we then know the chances with which all
 6 possible outcomes \vec{x} occur (see the discussion below Eq. (4)).

7 In practice, however, we only observe a single outcome \vec{x} of size n and have to infer θ
 8 instead. We write θ^* for the true value of the parameter θ that uses the relationship $f(\vec{x}|\theta)$
 9 to generate data \vec{x} . The goal of classical or frequentist point-estimation is to provide an
 10 educated guess of the true value θ^* by applying an estimator (i.e., a function), $T = t(\vec{X})$, to
 11 an anticipated outcome \vec{x} .

12 In this section we discuss how Fisher information can be used to determine the (asymptotically)
 13 least number of trials n that needs to be presented to a participant such that an esti-
 14 mator yields estimates at a certain level of accuracy. As n is decided upon pre-experimentally,
 15 we are uncertain of the outcomes of T . The randomness within T is inherited from the data
 16 that are assumed to be governed by $f(\vec{x}|\theta)$ and the participant's true latent ability θ^* . Fur-
 17 thermore, we call an observed outcome of T an estimate, which we denote as $t(\vec{x})$.

18 For instance, in Example 1 we infer a participant's latent ability θ using the mean
 19 estimator $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and write \bar{x} for the sample mean, an outcome of \bar{X} . As Example 1
 20 concerns an experiment we have to choose the number of trials n we plan to present to a
 21 participant.

22 Below we first show how n can be determined by the following scheme: (i) choose an
 23 estimator on which to base our inference, (ii) check whether this estimator is consistent, (iii)
 24 derive its sampling distribution, and (iv) derive n based on the sampling distribution for a
 25 certain target level of accuracy. We then apply this scheme to the normal model and show
 26 how it can be simplified by using the MLE in conjunction with the Fisher information.

27 **Measuring the Performance of an Estimator: Consistency, Sampling Distribu-**
 28 **tions, and Sample Size Determination Step (i): Choose an estimator** When the

1 trial sequence \vec{X} is assumed to be distributed according to the Bernoulli model, $X \sim \text{Ber}(\theta)$,
 2 as in the animal discrimination task, we typically infer θ using the mean estimator $\bar{X} =$
 3 $\frac{1}{n} \sum_{i=1}^n X_i$. **Step (ii): Assess consistency** The general belief is that the mean estima-
 4 tor is consistent, meaning that it is expected to extract more information about θ from an
 5 anticipated sample \vec{x} as the number of trials n increases. This qualitative behavior can be de-
 6 duced by studying the (asymptotic) mean and variance of the estimator, as the next example
 7 illustrates.

8 **Example 1** (Consistent Estimation within the Bernoulli Model, Continued). Use θ^* to denote
 9 a participant's unknown true latent ability, which we plan to estimate with the mean estimator
 10 \bar{X} . To study the expected behavior of a generic estimator T on the anticipated data, we average
 11 its possible outcomes across the sample space, resulting in the mean of the estimator: $E(T)$.
 12 To do so for $T = \bar{X}$, we use the linearity of the expectation and the i.i.d. assumption to
 13 express $E(\bar{X})$ in terms of the mean of the prototypical trial X :

$$E(\bar{X}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) \stackrel{i.i.d.}{=} \frac{n}{n} E(X) \quad (16)$$

$$= \sum_{x=0}^1 x f(x | \theta^*) \stackrel{Eq. (1)}{=} 1\theta^* + 0(1 - \theta^*) = \theta^*. \quad (17)$$

14 This implies that the mean estimator $E(\bar{X})$ is unbiased; however, to prove that \bar{X} is consistent
 15 we also need to show that the anticipated outcomes of the estimator to concentrate near θ^* as
 16 n increases.

17 To this end we need to study the variance of T . In particular, to calculate $\text{Var}(\bar{X})$ we
 18 once again exploit the i.i.d. assumption and the fact that the variance of the prototypical X
 19 is given by $\text{Var}(X) = \theta^*(1 - \theta^*)$:

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) \stackrel{i.i.d.}{=} \frac{n \text{Var}(X)}{n^2} = \frac{\theta^*(1 - \theta^*)}{n}. \quad (18)$$

20 Thus, the variance of \bar{X} shrinks to zero as the number of trials n grows indefinitely. This

1 implies that the chance that \bar{X} yield an estimate close to the true value goes to one and we
 2 therefore call \bar{X} consistent.⁴ ◇

3 Recall that when a participant's true latent ability is given by $\theta^* = 0.4$, there is a slim
 4 but non-zero chance that the participant never responds correctly in the animal discrimination
 5 task ($y = 0$; see Fig. 4). Such an extreme outcome will yield the estimate $\bar{x} = 0$, which is
 6 inaccurate as it is far from the true value $\theta^* = 0.4$. Hence, probabilistic models do not afford
 7 iron-clad guarantees: unrepresentative data will lead to inaccurate estimates.

8 However, Eq. (18) does justify the idea that an empirical scientist can obtain a more
 9 accurate measurement of a participant's latent ability θ by increasing n . In other words,
 10 consistency of \bar{X} , within the Bernoulli model, implies that the chance of \bar{X} yielding an
 11 inaccurate estimate goes to zero, whenever the number of trials n is sufficiently large.

12 **Step (iii): Derive sampling distribution** Clearly, we cannot detain a participant
 13 and get her to sit through a million trials or more. Hence, to quantify “sufficiently large” we
 14 first obtain the sampling distribution of the estimator.

15 **Example 1** (Sampling Distribution of the Sample Mean within the Bernoulli Model, Contin-
 16 ued). *Within the Bernoulli model we can consider \bar{X} as a rescaling of the number of successes*
 17 *in n trials, since $\bar{X} = \frac{1}{n}Y$, where $Y = \sum_{i=1}^n X_i$. Consequently, by Eq. (4) we know that the*
 18 *sampling distribution of \bar{X} is essentially a binomial with outcomes rescaled by the factor*
 19 *$\frac{1}{n}$. For instance, with $n = 10$ the estimator can take on the estimates (i.e., the outcomes)*
 20 *$0, 0.1, 0.2, \dots, 0.9, 1.0$ which corresponds to $0, 1, 2, \dots, 9, 10$ number of successes y .* ◇

21 This sampling distribution quantifies how both the true value of the parameter θ and
 22 the number of trials n affect the probable outcomes of the estimator \bar{X} , schematically:
 23 $(\theta, n) \xrightarrow{\text{Bin}(\theta, n)} \bar{x}$.⁵ As experimenters we cannot control θ , but we can lower the chances of
 24 inaccurate estimates by increasing n . **Step (iv): Derive n** In the following example we

⁴Formally, we call an estimator T for θ consistent, if for every possible true value θ^* , the estimator converges in probability, i.e., $\lim_{n \rightarrow \infty} P(|t(\bar{x}) - \theta^*| \geq \epsilon) = 0$, where n refers to the number of trials involved in a sequence of outcomes \bar{x} .

⁵Note the direction, knowing the true value θ^* allows us to deduce the chances with which the estimates occur, and *not* the other way around. A sampling distribution does *not* give us the probabilities of θ given an observed estimate.

1 show how the sampling distribution can be used to pre-experimentally deduce the number of
 2 trials n for a pre-specified level of accuracy.

3 **Example 1** (Determining the Number of Trials for the Sample Mean within the Bernoulli
 4 Model, Continued). *We wish to determine the number of trials n for an experiment \vec{X} with*
 5 $X \sim \text{Ber}(\theta)$ *such that the chance of obtaining an estimate that is more than 0.1 distance*
 6 *away from the true value is no larger than, say, $\alpha = 0.25$, regardless of the exact true value of*
 7 θ . *This means that we require n such that $P(\bar{X} \in [\theta^* \pm 0.1]) \geq 0.75$ for every possible true*
 8 θ^* . *To translate this targeted level of accuracy into a sample size n we use properties of*
 9 *the sampling distribution of \bar{X} as follows:*

10 *From consistency we know that the sampling distribution will increasingly concentrated*
 11 *its estimates around the true value as n grows, since $\text{Var}(\bar{X}) = \frac{\text{Var}(X)}{n}$. Hence, when the*
 12 *single-trial variance $\text{Var}(X)$ is already low to begin with we require fewer trials to the target*
 13 *level of accuracy. For instance, when $\theta^* = 0.99$, almost all outcomes will be 1 and the resulting*
 14 *estimate \bar{x} is also approximately 1 which is close to the true value.*

15 *On the other hand, when the single-trial variance is large, we require more trials to sat-*
 16 *isfy the aforementioned requirement. The goal is, thus, to tame the variance of the estimator*
 17 *and we should therefore set out to minimize the largest variation $\text{Var}(X) = \theta^*(1 - \theta^*)$, which*
 18 *occurs at $\theta^* = 0.5$. By plugging in $\theta = 0.5$ in the sampling distribution Eq. (4) we deduce that*
 19 *we require the participant to sit through at least $n = 25$ trials for the specified requirement to*
 20 *hold.*

21 *In other words, suppose the true value is $\theta^* = 0.5$ and we were to replicate the experi-*
 22 *ment \vec{X} say, $m = 100$ times yielding as many estimates $\bar{x}_1, \dots, \bar{x}_m$. Further suppose that each*
 23 *replication experiment consists of $n = 25$ trials. Then we expect that 75 of the 100 replicated*
 24 *estimates will fall between 0.4 and 0.6, satisfying the requirement $P(\bar{X} \in [\theta^* \pm 0.1]) \geq 0.75$.*

25 \diamond

26 **Intermediate Conclusion** In the above example we introduced concepts such as a “true
 27 value” θ^* for the parameter, consistency, and the sampling distribution of an estimator and
 28 showed how these concepts can be used in a pre-experimental analysis. In particular, we

1 showed that a specific estimator, \bar{X} , is consistent for θ , by studying the mean and the variance
 2 of this estimator. Furthermore, we showed how consistency of the mean estimator can be
 3 exploited to determine the number of trials n for an experiment at a pre-specified level of
 4 accuracy by controlling the largest variance at $\theta^* = 0.5$. Note that $\text{Var}(\bar{X}) \stackrel{\text{Eq. (10)}}{=} \frac{1}{nI(\theta)}$ and
 5 this worst-case scenario was therefore already hinted at in Fig. 4, as the data are then expected
 6 to be least informative. In the remainder of the section we further explore the connection of
 7 frequentist estimation theory and Fisher information by comparing two estimators within the
 8 normal model.

9 **Evaluating the Performance of an Estimator Within More Complex Models**

10 **Step (i): Choose an estimator** In more complex models, there are often several viable
 11 estimators for a parameter of interest. For instance, if the data are modeled according to
 12 a normal distribution, i.e., $X \sim \mathcal{N}(\mu, \sigma^2)$, we can then use both the mean \bar{X} and median
 13 estimator M for μ . This is due to the fact that μ is not only the population mean, but
 14 also the population median. Furthermore, estimates of the mean and the median estimator
 15 will typically differ. From a design point of view, we prefer the estimator that requires the
 16 least number of samples to distill a certain amount of information about μ . **Step (ii)+(iii):**
 17 **Assessing consistency and deriving the sampling distributions** First we show that
 18 within the normal model both the \bar{X} and M are consistent for μ by considering the means
 19 and variances of these estimators.

20 **Example 2** (Mean and Variances of the Sample Mean and Sample Median within the Normal
 21 Model, Continued). *When $X \sim \mathcal{N}(\mu, \sigma^2)$ the mean estimator \bar{X} has mean $E(\bar{X}) = \mu^*$ and
 22 variance $\text{Var}(\bar{X}) = \frac{\sigma^{*2}}{n}$, whatever the true values of μ and σ^2 . This can be shown as in
 23 Eq. (16) and Eq. (18). Moreover, a sum of normals is also normal, from which it follows
 24 that the sampling distribution of the mean estimator is also normal, $\bar{X} \sim \mathcal{N}(\mu^*, \frac{\sigma^{*2}}{n})$, this is
 25 a particular special case.*

26 *More commonly, however, we cannot derive the sampling distribution of an estima-*
 27 *tor exactly, that is, for finite n and we then resort to so-called large-sample approximations.*
 28 *In particular, instead of deriving the sampling of the median estimator exactly we call upon*

1 Corollary 21.5 from van der Vaart (1998), which tells us that the asymptotic sampling dis-
 2 tribution of median estimator is also normal, that is, $M \rightsquigarrow \mathcal{N}\left(\mu^*, \frac{2\pi\sigma^{*2}}{4n}\right)$, where the symbol
 3 \rightsquigarrow conveys that this equality only holds for large enough n . Note that the variances of both
 4 estimators \bar{X}, M tend to zero with the number of trials n increases. \diamond

5 **Step (iv): Derive n** Thus, both \bar{X} and M are consistent estimators for μ , and we
 6 can, therefore, determine, for each of these estimators, the number of trials that are required
 7 to reach a targeted level of accuracy.

8 **Example 2** (Determining the Number of Trials n within the Normal Model When the Vari-
 9 ance is Known, Continued). We wish to determine the number of trials n for an experiment \bar{X}
 10 with $X \sim \mathcal{N}(\mu, 1)$ such that the chance of obtaining an estimate that is less than 0.2 distance
 11 away from the true value μ^* is $1 - \alpha = 95.8\%$, regardless of the exact true value of μ .

12 To this end we exploit the fact that the estimators have an (asymptotically) normal
 13 sampling distribution. More specifically, we know that a probability of 95.8% of a normal
 14 distribution corresponds to two standard deviations of \bar{X} away from the mean of the estimator,
 15 the true value μ^* . Hence, we therefore translate the requirement $P(\bar{X} \in [\mu^* \pm 0.2]) \leq 0.042$
 16 into $\text{Var}(\bar{X}) \leq 0.01$. By solving $\text{Var}(\bar{X}) = \frac{\sigma^2}{n} = \frac{1}{n} = 0.01$ we see that we require at least
 17 $n = 100$ trials if we choose the mean estimator \bar{X} to estimate μ .

18 Similarly, as the asymptotic distribution of the median estimator M is also normal, we
 19 use the same translation of the targeted accuracy $1 - \alpha = 0.958$ in terms of the quantiles of a
 20 normal distribution with an asymptotic variance as derived in the previous example. Solving
 21 $\text{Var}(M) \rightsquigarrow \frac{2\pi}{4n} = 0.01$, we deduce that we require at least $n = 157$ trials to estimate μ whenever
 22 we use M . Hence, we therefore prefer the mean estimator \bar{X} over the median estimator M :
 23 the mean estimator \bar{X} is able to extract more information about μ from $n = 100$ trials than
 24 the median estimator. \diamond

25 Note that we approximated the exact sampling distribution of the median estimator
 26 by its asymptotic version. This simplification is applied commonly to most estimators as the
 27 exact sampling distribution is typically intractable. In this example, we replaced the exact
 28 sampling distribution of M by a normal distribution that is fully specified by its mean and

1 variance, i.e., $M \rightsquigarrow \mathcal{N}\left(\mu^*, \frac{2\pi}{4n}\right)$. Subsequently, we translated the targeted accuracy of the
 2 estimator into the quantiles of this normal sampling distribution, $\mathcal{N}\left(\mu^*, \frac{2\pi}{4n}\right)$, to calculate the
 3 required sample size n .

4 Hence, this approximation of n can be applied to any estimator that has an asymptotic
 5 normal sampling distribution for which we have to identify the corresponding asymptotic
 6 mean and variance. In particular, this holds true for the MLE as shown below. Recall that
 7 \bar{X} is also the MLE for μ whenever the data are modeled according to a normal distribution.

8 **Computational Advantage of the MLE** In general, the MLE $\hat{\theta}$ for θ within regular
 9 models⁶ is: **Step (ii)** consistent and **Step (iii)** has an asymptotically normal sampling
 10 distribution with an asymptotic mean given by the true value of the parameter, thus, it
 11 is “aimed” at the right value and with an asymptotic variance given by the inverse Fisher
 12 information, i.e.,

$$\hat{\theta} \rightsquigarrow \mathcal{N}\left(\theta^*, \frac{I^{-1}(\theta^*)}{n}\right) \text{ or equivalently, } \hat{\theta} \rightsquigarrow \theta^* + \epsilon \text{ with } \epsilon \sim \mathcal{N}\left(0, \frac{I^{-1}(\theta^*)}{n}\right), \quad (19)$$

13 where $I^{-1}(\theta^*)$ is the inverse unit Fisher information at the true value θ^* of the parameter
 14 that governs the possible outcomes. The latter formulation of Eq. (19) can be thought of
 15 as a linear regression in which the MLE regresses to the true value of the parameter with
 16 an error determined by the inverse Fisher information. In particular, note that this error
 17 disappears when the sample size n grows indefinitely, which confirms the claim that the MLE
 18 is consistent.

19 Note that Eq. (19) is not that surprising for \bar{X} as the central limit theorem states that
 20 the asymptotic sampling distribution of the mean estimator converges to a normal distribution
 21 $\bar{X} \rightsquigarrow \mathcal{N}\left(E(X), \frac{\text{Var}(X)}{n}\right)$, where $E(X)$ and $\text{Var}(X)$ refer to the population mean and the
 22 population variance respectively. The statement of the central limit theorem, however, is
 23 only restricted to mean estimators, while Eq. (19) also generalizes to other estimators as well:

⁶Regular models are models $f(x|\theta)$ that depend smoothly on the parameter θ , such that the Fisher information is non-zero and bounded for every θ not on the boundary, see van der Vaart (1998, pp. 61-64).

1 **Example 2** (The MLE for σ^2 within the Normal Model, Continued). *In many cases, we*
 2 *have to infer both the population mean and variance. If we choose the MLE for σ^2 , that is*
 3 $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ *then Eq. (19) tells us that $\hat{\sigma}^2 \rightsquigarrow \mathcal{N}(\sigma^{2*}, \frac{2\sigma^{4*}}{n})$, regardless of the true*
 4 *value σ^{2*} of σ^2 .*

5 *In the particular case that the data are modeled according to a normal distribution,*
 6 $X \sim \mathcal{N}(\mu, \sigma^2)$, *we can actually derive the sampling distributions of $\hat{\sigma}^2$ exactly and compare*
 7 *it to the approximation given in Eq. (19). To do so, we use the fact that the sums of squares*
 8 *divided by σ^2 have a χ^2 -distribution. As a result of this calculation we conclude that $E(\hat{\sigma}^2) =$*
 9 $(1 - \frac{1}{n})\sigma^{2*}$ *and that $\text{Var}(\hat{\sigma}^2) = (1 - \frac{1}{n})\frac{2}{\sigma^{4*n}} = (1 - \frac{1}{n})\frac{I^{-1}(\sigma^{2*})}{n}$, regardless of the true values*
 10 *of μ and σ^2 .*

11 *Note that the MLE is biased, i.e., $E(\hat{\sigma}^2) = (1 - \frac{1}{n})\sigma^{2*}$ and that this bias of $\frac{1}{n}\sigma^{2*}$*
 12 *can be made arbitrarily small regardless of the true value of σ^2 by choosing n large enough.*
 13 *Similarly, by choosing n large enough we can get the variance of the MLE $\text{Var}(\hat{\sigma}^2)$ close to*
 14 *the inverse Fisher information. Hence, this confirms that the inverse Fisher information is*
 15 *indeed a quite accurate approximation of the exact variance of the MLE. \diamond*

16 This last example allows us to view Eq. (19) as an extension of the central limit theorem
 17 to the MLE with an asymptotical mean and variance given by the true value of the parameter
 18 and the inverse Fisher information, respectively. In the example above Eq. (19) we also saw
 19 the merits of choosing the estimator with the lower asymptotic variance as this lower variance
 20 leads to smaller sample size n we require a participant to sit through.

21 **Estimator Selection by Fisher Information** Apart from being computationally conve-
 22 nient, the use of the MLE as an estimator for the parameter can also be theoretically justified
 23 by the Hájek-LeCam convolution theorem (Hájek, 1970; Inagaki, 1970).⁷, which guarantees
 24 that the sample size n derived from the MLE is (asymptotic) minimal.

25 The Hájek-LeCam convolution theorem states that whenever we choose to model the
 26 relationship $f(x|\theta)$ between the data and the latent property θ such that the Fisher in-
 27 formation exists for every possible true value θ^* (not on the boundary), then every regular

⁷See Ghosh (1985) for a beautiful review in which this theorem is related to the more popular, but also more restricted Cramér-Rao information lower bound derived by Rao (1945), Cramér (1946), and Fréchet (1943).

1 estimator⁸, such as the median estimator M within the normal model and the MLE, can
 2 be decomposed into an independent sum of an arbitrary random variable W and a normally
 3 distributed random variable ϵ that has a variance given by the inverse Fisher information.
 4 That is, any regular estimator T can be written as:

$$T \rightsquigarrow W + \epsilon, \text{ where } \epsilon \sim \mathcal{N}\left(0, \frac{I^{-1}(\theta^*)}{n}\right). \quad (20)$$

5 when n is sufficiently large. Due to the independence of W and ϵ , this implies that the
 6 asymptotic variance of any regular estimator T is given by $\text{Var}(T) = \text{Var}(W) + \text{Var}(\epsilon) =$
 7 $\text{Var}(W) + \frac{I^{-1}(\theta^*)}{n}$. Furthermore, as the variance of W cannot be negative we know that the
 8 asymptotic variance of any regular estimator is therefore *at least* the inverse Fisher information
 9 at the true value of θ . Clearly, this is the case for the MLE, Eq. (19), as the true value θ^*
 10 takes on the role of W , which does not vary within the frequentists' framework, $\text{Var}(W) = 0$.
 11 Hence, of all regular estimators the MLE has the lowest asymptotical variance and the sample
 12 size determined based its sampling distribution is therefore the lowest. This means that the
 13 MLE is expected to be (asymptotically) efficient in extracting all the information about the
 14 parameter of interest from the anticipated data.

15 **Summary of Section 1** We discussed how the MLE in conjunction with the Fisher in-
 16 formation can be used to determine the number of trials n that are to be presented to a
 17 participant in order to estimate the latent properties θ from observed behavior \vec{x} within a
 18 statistical model $f(\vec{x} | \vec{\theta})$ given a targeted level of accuracy $1 - \alpha$.

19 This method of sample size determination was justified by the Hájek-LeCam convolution
 20 theorem, which states that the resulting sample size is (asymptotically) minimal, if we choose
 21 to model the relationship $f(x | \theta)$ between the data and the latent property θ smoothly, that
 22 is, whenever the Fisher information is non-zero and finite for every possible true value θ^* that
 23 does not lie on the boundary of the parameter space.

⁸Regular estimators are estimators that have a limiting distributions that depends smoothly on the parameters, see Bickel, Klaassen, Ritov, and Wellner (1993) for a challenging, almost agonizing account. Furthermore, see Beran (1995) for a discussion of the convolution theorem and the bootstrap.

2. Using Fisher Information to Define a Default Prior According to Jeffreys' Rule

Introduction In this section we discuss the role Fisher information plays in Bayesian statistics. The Bayesian estimation procedure proceeds as follows: **Step 0** Propose a model that contains a parameter θ . For instance, one could choose the Bernoulli model in Example 1 which serves as the running example for the next two steps.

Step 1 Extend the model by assigning the parameter a distribution, denoted $g(\theta)$, prior to seeing the data. We can use $g(\theta)$ to incorporate previous findings about θ . Alternatively, we can use $g(\theta)$ to incorporate expert knowledge, and we can refer to $g(\theta)$ as the expert's prior beliefs. This implies that we consider θ as an additional random variable for which we can make probabilistic statements. **Step 2** Employ Bayes' theorem (see Eq. (21) below) to update knowledge about the parameter θ in light of the observed data. Bayes' theorem indicates how data should modify prior beliefs about θ to yield a posterior belief, measured by a posterior distribution $g(\theta | \vec{X} = \vec{x})$, schematically: $\theta \xleftarrow{g(\theta | \vec{X} = \vec{x})} \vec{x}$. If needed, the posterior distribution can be summarized by a point estimate and a credible interval, which reflects the uncertainty of the estimate.

In sum, the goal of coherently updating information using the observed data about the parameter of interest reduces to an application of Bayes' theorem. To this end, we first have to extend the model with a prior distribution $g(\theta)$ which might be unavailable when we venture into novel research. To nonetheless take advantage of the Bayesian machinery, one might be tempted to construct a "flat" or uniform prior distribution that assigns equal probability to each parameter value.

We start this section with the adoption of this flat prior to elaborate on the role of Bayes' theorem in **Step 2**. We then show the flaws of flat priors as they lead to conclusions that depend on the arbitrary way in which the problem is represented. To overcome this limitation we can use Jeffreys' rule to specify priors that are translation invariant Jeffreys (1946). The concept of Fisher information forms a key component of Jeffreys' rule.

1 **From Prior to Posterior: The Principle of Insufficient Reason and Bayes' Theorem**

2 When there is a complete lack of knowledge about the parameter of interest it seems intuitively
 3 sensible to expand the model $f(\vec{x}|\theta)$ with a prior that assigns equal probability to every
 4 possible value of θ , a rule known as Laplace's (1749 – 1827) "principle of insufficient reason"
 5 (Laplace, 1986; Stigler, 1986). **Step 1** In particular, this suggests extending data that are
 6 modeled according to a Bernoulli distribution Eq. (2) with a uniform prior for the parameter
 7 of interest $\theta \sim U[0, 1]$. **Step 2** Equipped with this prior, we can now draw conclusions about
 8 θ conditioned on observed data \vec{x} using Bayes' theorem:

$$g(\theta | \vec{X} = \vec{x}) = \frac{f(\vec{x} | \theta)g(\theta)}{\int_{\Theta} f(\vec{x} | \theta)g(\theta) d\theta}, \quad (21)$$

9 where $g(\theta)$ is the prior density function for the parameter θ . The left-hand side of Eq. (21) is
 10 known as the posterior density of θ , and Eq. (21) is often verbalized as:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}. \quad (22)$$

11 Note that the marginal likelihood does not involve the parameter θ , and that it is given by a
 12 single number that ensures that the area under the posterior distribution equals one, see the
 13 right panel in Fig. 5. It is common practice, therefore, to write $g(\theta | \vec{X} = \vec{x}) \propto f(\vec{x} | \theta)g(\theta)$,
 14 which says that the posterior is proportional to the likelihood times the prior. Furthermore,
 15 note that the posterior distribution is a combination of what we knew before we saw the data
 16 (i.e., the information in the prior distribution), and what we have learned from the data in
 17 terms of likelihood.

18 **Example 1** (Uniform prior on θ , Continued). *We assume that no prior research has been*
 19 *conducted for the animal discrimination task and set $g(\theta) = 1$ for every possible value of*
 20 *the ability parameter θ in $[0, 1]$ in accordance with Laplace's principle. We then apply the*
 21 *following scheme*

$$\text{Prior for } \theta \xrightarrow{\text{Data}} \text{Posterior for } \theta \quad (23)$$

1 to update our knowledge of θ . Fig. 5 illustrates this updating scheme with data characterized
 2 by $y = 7$ successes out of $n = 10$ trials.

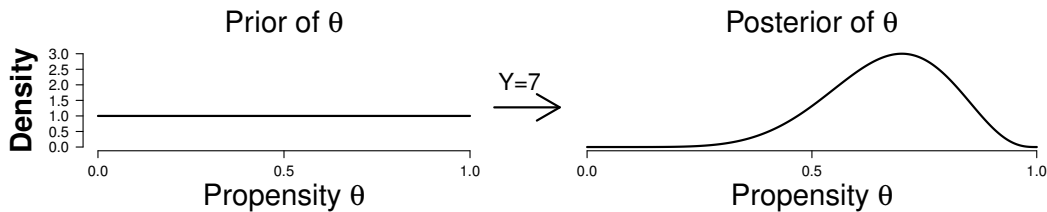


Figure 5. Bayesian updating based on observing $y = 7$ successes out of $n = 10$ Bernoulli trials. In the left panel, the flat prior distribution assigns equal probability to every possible value of θ . In the right panel, the posterior distribution is a compromise between the prior and the observed data (i.e., $y = 7$ and $n = 10$).

3 In this case, we have an analytically tractable posterior given by a beta distribution
 4 $g(\theta | y = 7) \sim \text{Beta}(8, 11)$ (e.g., Bayes & Price, 1763; Diaconis & Ylvisaker, 1979). If
 5 desired, the posterior distribution for a participant's latent animal discrimination ability θ
 6 can be summarized by a point estimate and a so-called credible interval. A commonly used
 7 point estimate is the posterior median, which is given by 0.68 and we can accompany it with
 8 a 95% credible interval that ranges from 0.39 to 0.89 to quantify our uncertainty about this
 9 estimate. Note the similarities between the posterior distribution in Fig. 5 and the likelihood
 10 function in Fig. 2; the maximum likelihood is in fact the mode of this posterior. \diamond

11 In more complicated data models, we might not be able to derive a posterior distribution
 12 analytically, as the integration involved in the marginal likelihood is typically hard to perform.
 13 However, due to the advent of computer-driven sampling methodology generally known as
 14 Markov chain Monte Carlo (MCMC: e.g., Gamerman & Lopes, 2006; Gilks, Richardson, &
 15 Spiegelhalter, 1996), we can directly sample sequences of values from the posterior distribution
 16 of interest, foregoing the need for closed-form analytic solutions (e.g., Lee & Wagenmakers,
 17 2013). In the following, we omit the derivations of the posteriors and instead provided R code

1 on the first author’s website.

2 **The Principle of Insufficient Reason and its Dependence on the Parametrization**

3 **Step 1** Laplace’s principle seems to yield priors that reflect complete ignorance about the
 4 parameter θ . In spite of its intuitive formulation, this is not true, as it depends on how a
 5 statistical problem is represented. To demonstrate this paradoxical result we exploit the sim-
 6 ilarity the Bernoulli model has to a coin flip experiment with a bent coin that has propensity
 7 $\theta = P(x = 1)$ to land heads instead of tails. Furthermore, we say that this propensity is due
 8 to the angle, denoted as ϕ , unbeknownst to the researcher, with which the coin is bent.

9 Moreover, suppose that the coin’s physical relation between the angle ϕ and the propen-
 10 sity θ is given by the hypothetical function $h(\phi) = \frac{1}{2} + \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3$.⁹ Thus, the generative model
 11 can be modified in terms of ϕ by the following chain: $\phi \xrightarrow{h} \theta \xrightarrow{f(\vec{x}|\theta)} \vec{x}$. In the absence of
 12 prior knowledge we could therefore model our ignorance about the angles ϕ according to a
 13 uniform prior on $(-\pi, \pi)$, thereby extending the Bernoulli model with a prior density given
 14 by $g(\phi) = \frac{1}{2\pi}$ for every possible value ϕ of the angle. **Step 2** This uniform prior on the
 15 angles of ϕ then induces a prior on θ for which we can subsequently apply Bayes’ theorem,
 16 as schematically described by Eq. (24):

$$\begin{array}{ccc}
 \text{Prior on } \phi & & \\
 h \downarrow & & (24) \\
 \text{Prior on } \theta \text{ induced from } \phi & \xrightarrow{\text{Data}} & \text{Posterior for } \theta \text{ induced from } \phi
 \end{array}$$

17 The results of this chain of induction is plotted in Fig. 6.

18 The resulting posterior can also be obtained by first updating the uniform prior on ϕ
 19 to a posterior on ϕ , which is subsequently transformed into a posterior for θ . Either way, the
 20 resulting posterior on θ will differ substantially from the previous result plotted in the right
 21 panel of Fig. 5. Hence, the principle that was set out to reflect ignorance depends greatly on
 22 how the problem is represented. In particular, a uniform prior on the coin’s angle ϕ yields a

⁹This function was chosen purely for mathematical convenience.

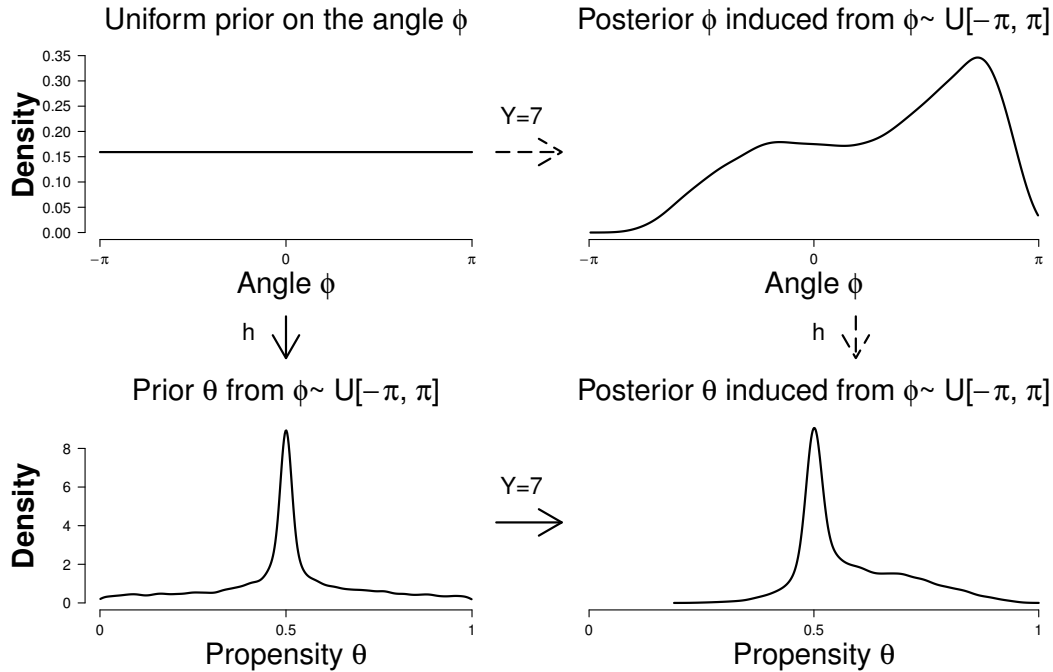


Figure 6. A uniform prior is not always uninformative: example of a bent coin landing heads in $y = 7$ out of $n = 10$ tosses. The top left panel is the starting point of our investigation with a uniform prior on the coin’s angle ϕ . The corresponding bottom left panel shows that this prior distribution is highly informative for the propensity of the coin as it induces a lot of mass near $\theta = 0.5$. This is also apparent from the bottom right panel, which shows the posterior for θ with a posterior mean at 0.53 and a 95%-credible interval that ranges from 0.42 to 0.84. The same posterior on θ is obtained if we proceed via an alternative route in which we first updated the prior on ϕ to a posterior on ϕ (top left panel) and then use the transformation h to obtain the posterior on θ .

- 1 highly informative prior in terms of the coin’s propensity θ . Similarly, a uniform prior on the
- 2 propensity induces a prior on ϕ that assigns relatively much mass near the endpoints $-\pi$ and
- 3 π .

4 It may not be clear, therefore, on what scale we should apply Laplace’s principle of
 5 insufficient reason: one can defend a uniform prior on propensity or on angle, but neither
 6 representation seems privileged in any particular way.

7 **Using Fisher Information to Construct a Translation Invariant Prior Step 1** This
 8 dependence on how the problem is parameterized has been put forward as a major objection
 9 to Bayesian inference (e.g., Edwards, 1992; Fisher, 1930).

10 One way to meet this objection is to use a rule for constructing prior distributions that

1 are translation invariant, i.e., that will lead to the same conclusion (i.e., the same posterior
 2 distribution) regardless of how the problem is represented. The following rule proposes a
 3 prior that is proportional to the square root of the Fisher information:

$$g(\theta) = \frac{\sqrt{I(\theta)}}{\int \sqrt{I(\theta)} d\theta} \propto \sqrt{I(\theta)}, \quad (25)$$

4 which is known as the prior obtained from Jeffreys' rule, or Jeffreys' prior (Jeffreys, 1946).
 5 We now apply this prior to the coin example both in terms of θ and ϕ to demonstrate that it
 6 is indeed translation invariant.

7 **Example 1** (Invariance of the Jeffreys' Prior within the Bernoulli Model, Continued). **Step 2**
 8 *To apply Jeffreys' rule, Eq. (25), to the propensity θ for data distributed according to a*
 9 *Bernoulli distribution, recall that the unit information for θ is given by $I(\theta) = \frac{1}{\theta(1-\theta)}$, from*
 10 *which it follows that*

$$g(\theta) = \frac{1}{\pi} \theta^{-0.5} (1 - \theta)^{-0.5}. \quad (26)$$

11 *The proof that $\int \sqrt{I(\theta)} d\theta = \pi$ can be found in the online appendix. Jeffreys' prior on the coin*
 12 *propensity θ is plotted in the bottom left panel of Fig. 7. We can then apply Bayes' theorem,*
 13 *using the data (i.e., $y = 7$ heads out of $n = 10$ tosses) to update this prior to a posterior on θ*
 14 *(bottom right panel) and subsequently use the backwards transform $k(\theta) = \phi$ to translate this*
 15 *result in terms of ϕ (top right panel).*

16 *A re-parametrization invariant rule implies that we could as well have started in terms*
 17 *of ϕ . The online appendix shows that the Fisher information for ϕ is given by $I(\phi) = \frac{9\phi^4}{\pi^6 - \phi^6}$,*
 18 *which yields the following prior*

$$g(\phi) = \frac{1}{\pi} \frac{3\phi^2}{\sqrt{\pi^6 - \phi^6}}. \quad (27)$$

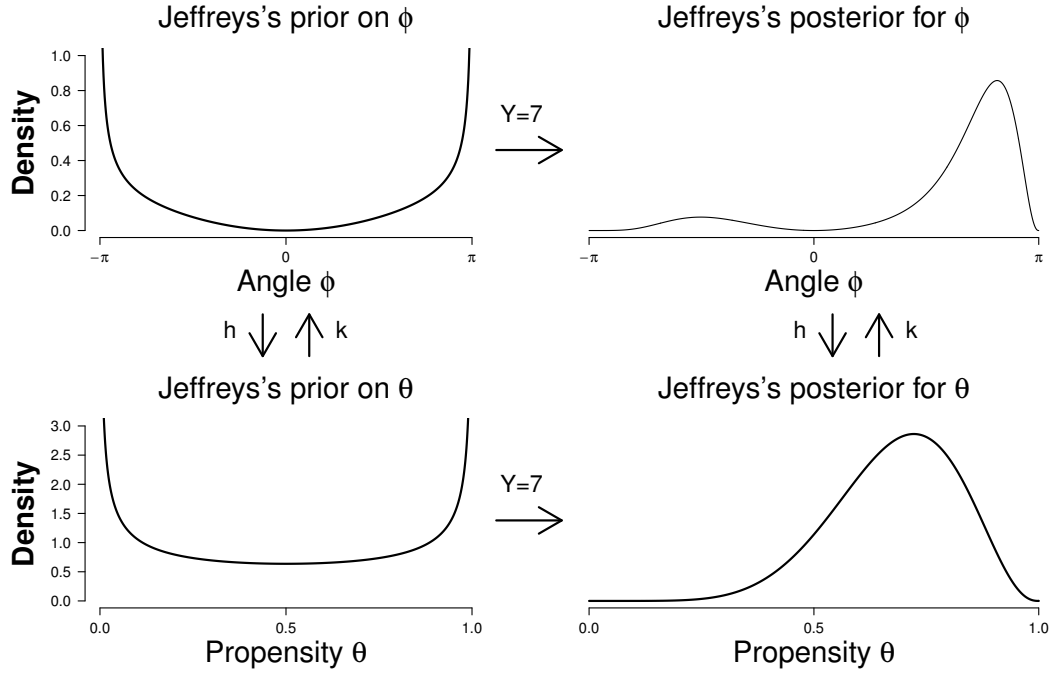


Figure 7. For priors constructed through Jeffreys' rule it does not matter whether the problem is represented in term of the coin's angles ϕ or its propensity θ . Thus, not only is the problem equivalent due to the transformations $\theta = h(\phi)$ and $\phi = k(\theta)$, the prior information is the same in both representations. This also holds for the posteriors. The posterior median for θ is given by 0.69 with a 95% credible interval of (0.40, 0.91).

1 Note that the normalization constant for this prior is also π , a property we elaborate on in
 2 the next section. The top left panel of Fig. 7 shows the Jeffreys' prior in terms of ϕ and the
 3 top right panel shows the resulting posterior. This posterior is identical to the one obtained
 4 from the previously described updating procedure that starts by postulating Jeffreys' prior on
 5 θ instead of on ϕ . ◇

6 The above example shows how the Jeffreys' prior is constructed from a rule such that
 7 the same posterior knowledge is reached regardless of the scale for the prior. That is, one
 8 draws the same conclusions about θ regardless of whether we (1) use Jeffreys' rule to construct
 9 a prior on θ and update with the observed data, or (2) use Jeffreys' rule to construct a prior
 10 on ϕ , update to a posterior distribution on ϕ , and then transform that posterior in terms of
 11 θ . In contrast, a flat prior on ϕ induces a posterior distribution on θ that differs substantially
 12 from the one obtained by assigning a flat prior on θ (i.e., compare the bottom right panel of
 13 Fig. 6 to Fig. 5).

1 **Summary of Section 2** In this section we showed how a model $f(\vec{x}|\theta)$ can be extended
 2 with default priors constructed according to Jeffreys' rule, a rule that uses Fisher information
 3 to define default priors that are translation invariant in the sense that they do not depend on
 4 the arbitrary scale on which the problem is defined.

5 One of the drawbacks of Jeffreys' rule is that it does not apply to multiple parameters
 6 simultaneously, for instance, when $\vec{\theta} = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$ within the normal model. Jeffreys (1961) suggests
 7 to apply this rule to each parameter individually, which he justified due to the fact that the
 8 Fisher information has zeros in the off-diagonal entries, see Eq. (15) and the discussion about
 9 orthogonality below. Another drawback is that Jeffreys' rule may lead to improper priors
 10 (i.e., priors which are non-normalized). For instance, the Jeffreys' prior for μ within a normal
 11 model is uniform over the real line as μ does not depend on σ . Surprisingly, however, this
 12 prior does lead to proper normalized posteriors.

13 3. Using Fisher Information to Measure Model Complexity

14 **Introduction** In this section we show how Fisher information can be used to measure
 15 model complexity from the perspective of data compression and the principle of minimum
 16 description length (Grünwald, 2007; Myung, Forster, & Browne, 2000; Pitt, Myung, & Zhang,
 17 2002).

18 Different theories are implemented as different statistical models, and model selection
 19 therefore helps to quantify the extent to which the data support one theory over another.
 20 However, the support in the data cannot be assessed by solely considering goodness-of-fit, as
 21 the ability to fit random data increases with model complexity, see (e.g., Roberts & Pashler,
 22 2000).

23 For this reason, every method for model selection needs to take into account the trade-off
 24 between goodness-of-fit and parsimony: more complicated models necessarily lead to better
 25 fits but may in fact over-fit the data. Such overly complicated models capture idiosyncratic
 26 noise rather than general structure, resulting in poor model generalizability (Myung, Forster,
 27 & Browne, 2000; Wagenmakers & Waldorp, 2006). The goal of most model selection methods,
 28 therefore, is to acknowledge the trade-off between goodness-of-fit and parsimony in some

1 principled way, allowing the selection of models that best predict new data coming from the
 2 same source.

3 Most popular amongst the many model selection methods are the penalized maximum
 4 likelihood criteria, including the Akaike information criterion (AIC; Akaike, 1974; Burnham
 5 & Anderson, 2002), the Bayesian information criterion (BIC; Raftery, 1995; Schwarz, 1978),
 6 and the Fisher information approximation based on the minimum description length principle
 7 (FIA and MDL; Rissanen, 1996; Grünwald, 2007; Pitt et al., 2002). All these methods have a
 8 single component that quantifies goodness-of-fit (i.e., through maximum likelihood) and one
 9 or more penalty components that discount the goodness-of-fit by the degree to which it was
 10 accomplished using a complex model. The methods differ in the way that they quantify the
 11 penalty term for complexity. Specifically, AIC, BIC, and FIA are formalized as follows:

$$\text{AIC} = -2 \log f(\vec{x} | \hat{\theta}) + 2d \quad (28)$$

$$\text{BIC} = -2 \log f(\vec{x} | \hat{\theta}) + d \log(n) \quad (29)$$

$$\text{FIA} = \underbrace{-\log f(\vec{x} | \hat{\theta})}_{\text{Goodness-of-fit}} + \underbrace{\frac{d}{2} \log \frac{n}{2\pi}}_{\text{Dimensionality}} + \underbrace{\log \left(\int_{\Theta} \sqrt{\det I(\theta)} d\theta \right)}_{\text{Geometric complexity}}, \quad (30)$$

12 where d refers to the number of free parameters, n refers to sample size, and $\hat{\theta}$ refers to the
 13 MLE. For all criteria, the model with the lowest criteria is best. Note that, in contrast to
 14 AIC and BIC, the FIA penalizes both for the number of models parameters and for their
 15 functional form. The latter is accomplished through the geometric complexity term, a term
 16 that is computed using the Fisher information (Myung, Balasubramanian, & Pitt, 2000).

17 The goal of this section is to visualize how Fisher information can be used to assess
 18 geometric complexity. These ideas are based on Kass (1989) and illustrated with a set of
 19 simple multinomial processing tree (MPT) models (e.g., Batchelder & Riefer, 1999; Wu,
 20 Myung, & Batchelder, 2010; Klauer & Kellen, 2011). For details about the philosophy of
 21 data compression as a method for model selection we refer the interested reader to the books
 22 by Grünwald, Myung, and Pitt (2005) and Grünwald (2007).

1 **Running Example with Three Outcomes** To demonstrate the role of Fisher information
 2 for model selection we use the following source-memory task: during an initial study phase,
 3 participants are presented with two list of words on a computer screen. List L is projected
 4 on the left-hand side and list R is projected on the right-hand side. In a later test phase,
 5 the participant is presented with two words, side by side, that can stem from either list,
 6 ll, lr, rl, rr . The participant is asked to categorize these pairs as follows:

- 7 1. Both words come from list L , i.e., ll ,
- 8 2. The words are mixed M , i.e., lr or rl ,
- 9 3. Both words come from list R , i.e., rr .

10 As before we assume that the participant is presented with n test pairs of equal difficulty,
 11 yielding the trial sequence $\vec{X} = (X_1, \dots, X_n)$ consisting of n i.i.d. copies of a prototypical
 12 trial X that has three outcomes. Below we propose three process models for X and show
 13 how MDL model selection using Fisher information can be used to measure the complexity
 14 of each model.

15 **General Model for a Random Variable with Three Outcomes** As a starting point,
 16 we discuss a general model for random variables X with three outcomes. This model, shown
 17 in the top left panel of Fig. 8, assumes that the participant categorizes words as pairs in a
 18 cascaded fashion: with probability a , the participant concludes that the words are mixed.
 19 With probability $1 - a$, both words are perceived to stem from the same list, and the partic-
 20 ipant subsequently decides whether the two words come from list L (with probability b) or
 21 list R (with probability $1 - b$).

22 The model parameters naturally depend on the nature of the stimulus; for mixed pairs,
 23 for instance, parameter a will be higher than for pure pairs. For the purposes of this tutorial
 24 we wish to keep the models as simple as possible, and hence we assume that the presented
 25 word pair stimulus subject to modeling is always “rr”.

26 In particular, when the participant’s parameters are given by $a = \frac{1}{3}, b = \frac{1}{2}$, then the
 27 model predicts that $P(X = L) = P(X = M) = P(X = R) = \frac{1}{3}$ meaning that the participant

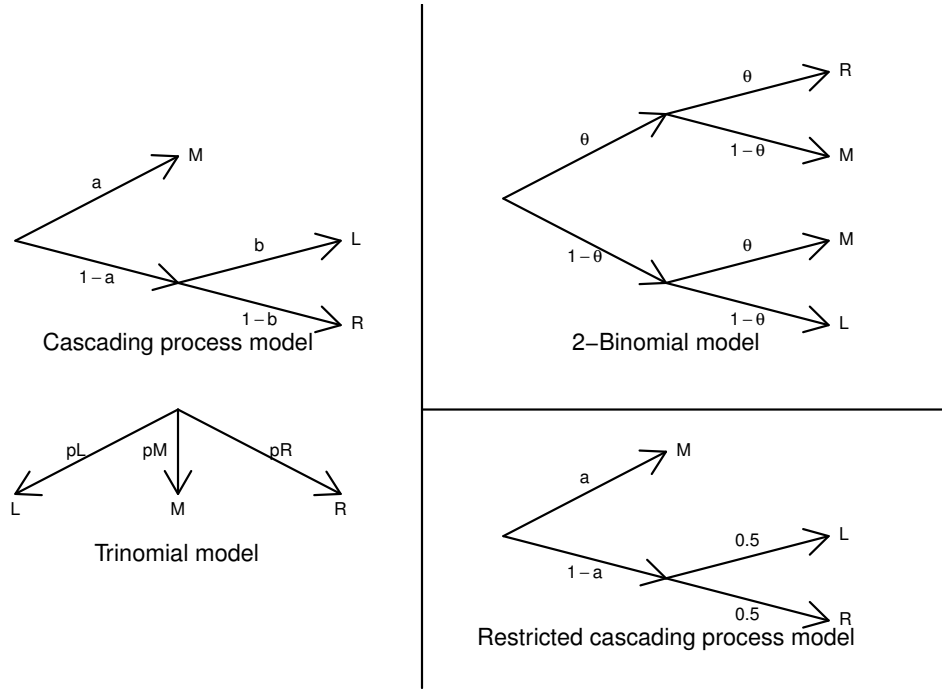


Figure 8. Three MPT models that describe a participant’s choices (i.e., L , M , or R) in the source-memory task described in the main text. The left panel shows two equivalent models: the cascading process model and the trinomial model. The top right panel shows the 2-binomial model and the bottom right panel shows a restricted version of the cascading process model with b fixed at $b = 0.5$. Note: in order to limit the number of model parameter, we assume that the presented word pair stimulus subject to modeling is always “rr”.

1 is to respond L, M, R with the same chances, and we therefore expect to see a data pattern
 2 that consists of an equal number of responses L, M , and R . In general, we can model the
 3 expected data patterns directly by re-parametrizing the cascading process model in terms of
 4 $pL = P(X = L) = (1-a)b, pM = P(X = M) = a$, and $pR = P(X = R) = (1-a)(1-b)$. This
 5 model is generally referred to as a (one-trial) trinomial model, $X \sim \text{Tri}(pL, pM, pR)$, which
 6 in fact has only two free parameters, pL, pM , as $pR = 1 - pL - pM$. This latter equation
 7 allows us to depict the data patterns that this model accommodates as a surface in three
 8 dimensions, see Fig. 9. As the components of each point $\vec{p} = (pL, pM, pR)$ on the surface sum
 9 to one, we may refer to each such point as an expected data pattern of the cascading process
 10 model or a pdf¹⁰ of the trinomial model.

11 A more natural representation of the trinomial model, which helps relate Fisher infor-

¹⁰This is clearly a probability mass function, since we are dealing with discrete outcomes. As the ideas elaborated here extend to pdfs as well, we do not distinguish the two here.

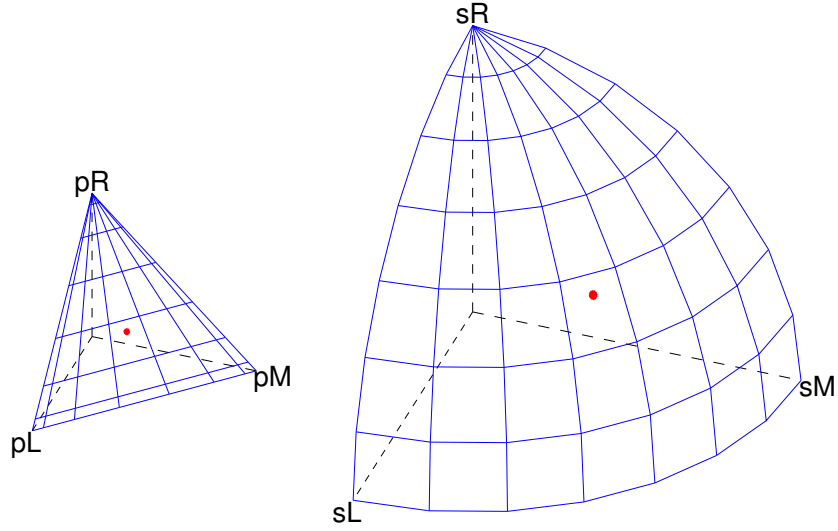


Figure 9. Expected data patterns from the trinomial model (left panel) and the transformed trinomial model (right panel). In the left panel, for instance, the left hand corner $\vec{p} = (1, 0, 0)$ represents the pdf that assigns all its mass to the L response and none to the other outcomes. The dot in the middle of the surface represents the pdf $\vec{p} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ that assigns equal probability to every outcome. Points outside the surface cannot be pdfs as they do not sum to one.

1 mation to model complexity, can be obtained by transforming the probabilities pL, pR, pM to
 2 two times the square roots of each probability: $sL = 2\sqrt{pL}, sM = 2\sqrt{pM}, sR = 2\sqrt{pR}$; the
 3 resulting surface is shown in the right panel of Fig. 9. Instead of forming a triangular surface,
 4 the transformed trinomial model now forms the positive octant of a sphere with radius two.
 5 In this representation all possible data patterns of the trinomial model are two units away
 6 from the origin, which can be easily verified by the Pythagorean theorem:

$$\|\vec{s}\| = \sqrt{sL^2 + sM^2 + sR^2} = \sqrt{4(pL + pM + pR)} = 2, \tag{31}$$

7 where we have written $\|\vec{s}\|$ for the length of an arbitrary point \vec{s} on the sphere. In conclusion,
 8 the trinomial model is a general model for X and we will use this fact to visualize the
 9 predictions from the more restricted 2-binomial model defined below.

1 **The Single-List Memorization Strategy** An alternative theory for the underlying cog-
 2 nitive process, shown in the top right panel of Fig. 8, assumes that the participant assesses
 3 the words individually but only with respect to a single list, say list R . In other words,
 4 with probability θ the participant decides that the first word comes from the right list. This
 5 decision process is repeated for the second word, after which the complete word pair can be
 6 categorized as L , M , or R .

7 This model has one parameter, θ , yielding expected data patterns given by $pL =$
 8 $(1 - \theta)^2$, $pM = 2\theta(1 - \theta)$ and $pR = \theta^2$. These probabilities correspond to a 2-binomial model
 9 which we define by $X \sim \text{Bin}(\theta, 2)$ where a single word coming from list R is seen as a success
 10 with probability θ . For the representation on the sphere, this yields:

$$\theta \mapsto \vec{s}(\theta) = \begin{pmatrix} sL(\theta) \\ sM(\theta) \\ sR(\theta) \end{pmatrix} = \begin{pmatrix} 2\sqrt{pL} \\ 2\sqrt{pM} \\ 2\sqrt{pR} \end{pmatrix} = \begin{pmatrix} 2(1 - \theta) \\ 2\sqrt{2\theta(1 - \theta)} \\ 2\theta \end{pmatrix} \quad (32)$$

11 As a 2-binomial distribution also has three outcomes, we can represent the corresponding
 12 probabilities within the trinomial model, see Fig. 10. Note that because the 2-binomial
 13 model has only one free parameter, its data patterns form a line instead of a surface.

14 **Expected Data Patterns and Model Specificity** The 2-binomial model can be sim-
 15 plified even further. For instance, we might entertain the hypothesis that the participant's
 16 responses are governed by a specific value for the free parameter, such as $\theta_0 = 0.5$. From this
 17 hypothesis, we expect to see 30 L , 60 M , and 30 R responses in $n = 120$ trials, hence, the data
 18 pattern $\vec{p} = (0.25, 0.50, 0.25)$.

19 Real data, however, will typically deviate around the expected data pattern even if
 20 the hypothesis $\theta_0 = 0.5$ holds true exactly. To see this, Fig. 11 shows data patterns of 200
 21 synthetic participants each completing 120 trials of the source-memory task with $\theta_0 = 0.5$
 22 according to the 2-binomial model.

23 Fig. 11 confirms the fact that even though the *expected* data patterns from the 2-

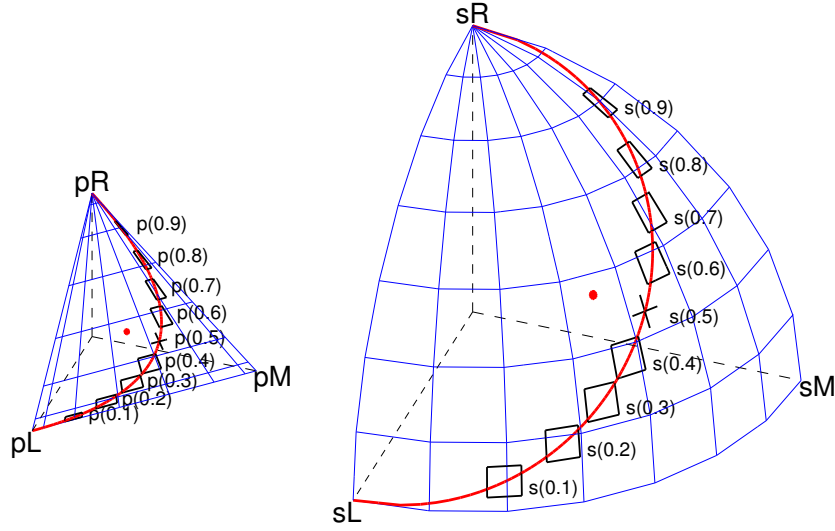


Figure 10. Expected data patterns from the 2-binomial model are represented by the curves within the trinomial model (left panel) and the transformed trinomial model (right panel). In the 2-binomial model $\theta = 0$ corresponds to a point on the surface floor, as we then expect only L responses. Data patterns corresponding to parameter values $0.1, 0.2, \dots, 0.8, 0.9$ are represented by rectangular symbols, except for the parameter value $\theta = 0.5$ which is represented by a cross. See Eq. (32) for the definition of $\vec{s}(\theta)$.

1 binomial model are confined to the line, random data patterns generated from such a model
 2 are not.

3 To illustrate how specific the 2-binomial model is in generating data, Fig. 12 shows the
 4 results from a similar simulation using other values of the θ parameter: $\theta = 0.1, 0.2, \dots, 0.9$.

5 Each dot in in Fig. 12 represents the data pattern of a synthetic participant completing
 6 120 trials and as the number of trials increases indefinitely we will see that the observed data
 7 pattern will coincide with the expected data pattern. This confirms the relation between the
 8 size of the expected data pattern, i.e., the curve in this case, and model specificity.

9 To contrast the specificity of the 2-binomial model to that of the trinomial model, note
 10 that the expected data patterns for the trinomial model are not restricted to the curve. In
 11 effect, the predictions of the trinomial model cover the whole sphere. This also implies that
 12 the trinomial model is less specific in its predictions, making it harder to falsify. This is why
 13 the trinomial model is said to be more complex than the 2-binomial model.

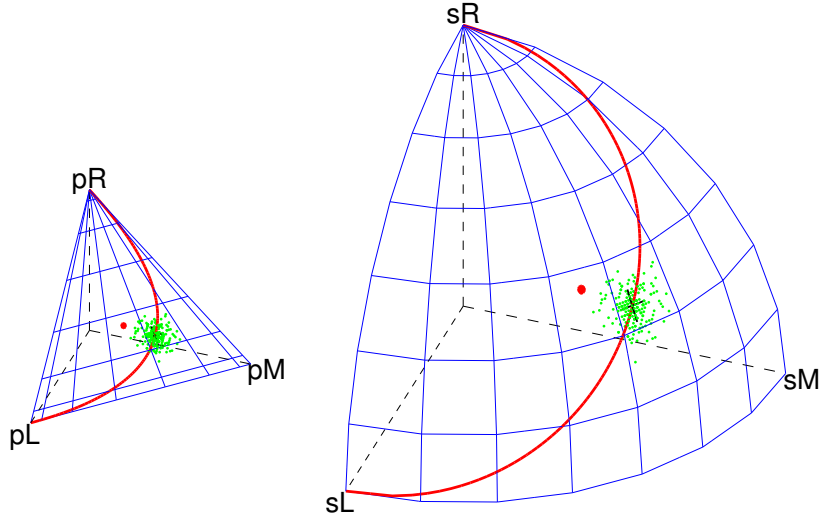


Figure 11. The curve represents the expected data patterns from the 2-binomial model and the cross represents the expected data pattern from to the restricted 2-binomial model with $\theta_0 = 0.5$. Each dot represents the data of a synthetic participant completing $n = 120$ trials.

1 **The Trade-off between Model Specificity and Goodness-of-Fit** What the trinomial
 2 model lacks in specificity is compensated by its ability to produce good fits for all data patterns
 3 with three outcomes. In contrast, when the data are simulated from the equiprobable pdf
 4 $X \sim \text{Tri}(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ we cannot yield good fits within the 2-binomial model as this pdf is not
 5 included on the line, see Fig. 10 - 12.

6 Thus, as the surface plot suggest, the 2-binomial model can never yield better fits
 7 than the trinomial model, as the former is nested within the latter. In other words, the
 8 maximum likelihood for the trinomial model is as least as large as the maximum likelihood
 9 for the 2-binomial model: $\log f(\vec{x} | \hat{\theta} \text{ within trinomial}) \geq \log f(\vec{x} | \hat{\theta} \text{ within 2-binomial})$. For
 10 this reason, model selection information criteria such as AIC and BIC do not consider only
 11 goodness-of-fit measures but also penalize for the number of free parameters. Within the MDL
 12 philosophy Eq. (30) such a penalty is incomplete because it ignores differences in model com-
 13 plexity due to the functional relationship between the parameters and the data, a relationship
 14 that can be measured using Fisher information.

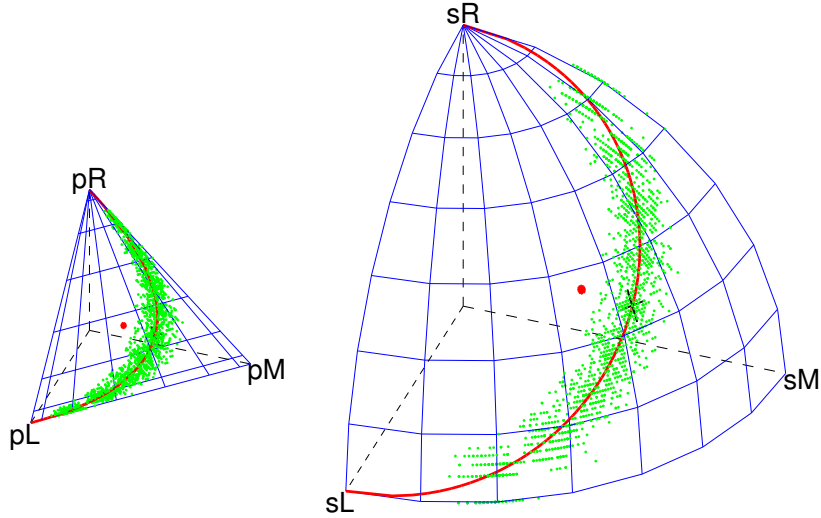


Figure 12. The curve represents the expected data patterns from the 2-binomial model and the cross represents the expected data pattern from to the restricted 2-binomial model with $\theta_0 = 0.5$. Each dot represents the data of a synthetic participant completing $n = 120$ trials. Data were simulated for each $\theta = 0.1, 0.2, \dots, 0.9$.

1 **Using Fisher Information to Measure the Size of a Model** In the previous section
 2 we related the complexity of a model to its expected data patterns. An intuitive measure for
 3 model complexity is therefore its size. In particular, we can use the length of the curve in
 4 Fig. 10 to measure the complexity of the 2-binomial model. Before we show that this length
 5 can be calculated from the Fisher information, recall that a single trial of a 2-binomial model
 6 consists of two independent Bernoulli flips. We therefore conclude that the unit information
 7 for the 2-binomial model is given by $I(\theta | 2\text{-binomial}) = \frac{2}{\theta(1-\theta)}$.

8 To calculate the length of the curve, we use the fact that we can linearly approximate,
 9 i.e., Taylor approximate, the parametrization from $\theta \mapsto \vec{s}(\theta)$. This strategy consists of two
 10 steps: (a) linear approximations and (b) summing the linear approximations over the domain
 11 of θ .

12 **Step a. Linear Approximation** In Fig. 13 we illustrate the first step by a linear approx-
 13 imation of the curve by considering three approximation points, namely: 0.05, 0.40, 0.75.¹¹

¹¹For technical reasons we avoid $\theta = 0$ for now.

For each of these approximation points we calculate $\vec{s}(\theta)$ and the key idea is to approximate

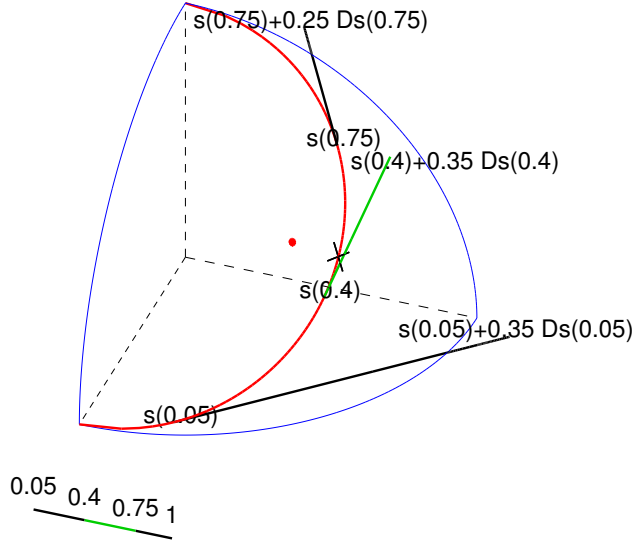


Figure 13. The size of a model can be calculated by linear approximations. The smooth curve represents the expected data patterns under the 2-binomial model. The three line segments below the bottom left corner show the partitioning of the subdomain $(0.05, 1)$ for the corresponding tangent vectors along the curve. The length of the curve can be approximated by summing the length of the three lines tangent to the curve. For clarity, we removed the grid lines that were drawn in Fig. 9, but we retained both the reference point corresponding to the data pattern $\text{Tri}(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and the plus sign for the expected data pattern for $\theta = 0.5$ within the 2-binomial model.

1
 2 the “gap” between, say, $\vec{s}(0.05)$ and $\vec{s}(0.4)$ by a linear extrapolation at $\vec{s}(0.05)$. Hence, for
 3 $\theta_0 = 0.05$ this yields

$$\vec{s}(\theta_0 + h) \approx \vec{s}(\theta_0) + hD\vec{s}(\theta_0), \tag{33}$$

4 where $D\vec{s}(\theta_0)$ denotes the 3-dimensional tangent vector at θ_0 and we used $h = 0.40 - 0.05 =$
 5 0.35 at $\theta_0 = 0.05$ in Fig. 13. This procedure is repeated for the other points $\theta = 0.40$ and
 6 $\theta = 0.75$ and we approximate the arc length by summing the length of the tangent vectors at
 7 each point. The accuracy of this approximation of the arc length increases with the number
 8 of approximation points. For this reason we desire an expression for the tangent vector at
 9 each arbitrarily chosen point θ , which can be obtained by taking the derivative with respect
 10 to θ for each component (i.e., possible outcomes L, M, R):

$$D\vec{s}(\theta) = \begin{pmatrix} \frac{d}{d\theta} sL \\ \frac{d}{d\theta} sM \\ \frac{d}{d\theta} sR \end{pmatrix} = \begin{pmatrix} \frac{d}{d\theta} \log(pL) \sqrt{pL} \\ \frac{d}{d\theta} \log(pM) \sqrt{pM} \\ \frac{d}{d\theta} \log(pR) \sqrt{pR} \end{pmatrix} = \begin{pmatrix} -2 \\ \frac{2-4\theta}{\sqrt{2\theta(1-\theta)}} \\ 2 \end{pmatrix} \quad (34)$$

1 where the second equality is due to the general relation between $\vec{s}(\theta)$ and $\vec{p}(\theta)$. The length of
 2 this tangent vector can be calculated by an application of the Pythagorean theorem, which
 3 yields:

$$\|D\vec{s}(\theta)\| = \sqrt{(-2)^2 + \left(\frac{2-4\theta}{\sqrt{2\theta(1-\theta)}}\right)^2 + 2^2} = \sqrt{\frac{2}{\theta(1-\theta)}} = \sqrt{I(\theta | 2\text{-binomial})}, \quad (35)$$

4 As θ was chosen arbitrarily we conclude that the tangent length at each $\theta \in (0, 1)$ within the
 5 2-binomial model is given by $\sqrt{I(\theta | 2\text{-binomial})}$, the square root of the Fisher information.

6 **Step b. Accumulating the Tangent Lengths** To establish the size of a model, we now
 7 have to sum the tangent lengths over all approximation points. This sum becomes more
 8 accurate as we increase the number of approximation points, eventuating into an integral.
 9 For the 2-binomial this yields

$$V_{2\text{Bin}} = \int_0^1 \sqrt{I(\theta | 2\text{-binomial})} d\theta = \sqrt{2} \int_0^1 \frac{1}{\sqrt{\theta(1-\theta)}} d\theta = \sqrt{2}\pi, \quad (36)$$

10 see the online appendix for the full calculation. Hence, the length of the curve is given by
 11 $V_{2\text{Bin}} = \int_{\Theta} \sqrt{I(\theta | 2\text{-binomial})} d\theta = \sqrt{2}\pi$, and we can use it as a measure for the specificity of
 12 the 2-binomial model, and, hence, its model complexity.

13 **Fisher Information as a Measure of Complexity for Models with Vector-valued**
 14 **Parameters** The ideas developed above can be easily generalized to models with vector-
 15 valued parameters such as the trinomial model with free parameters pL, pM or, equivalently,

1 the cascading process model with free parameters a, b . As there are two parameters to vary
 2 over we then have two tangent vectors which form a tangent surface. To compute the size of
 3 the model we then have to calculate the surface area of the Fisher information matrix, which
 4 is given by its determinant (for details see the online appendix). Hence,

$$V = \int \sqrt{\det(I(\vec{\theta}))} d\vec{\theta}, \tag{37}$$

5 where we have written $\vec{\theta}$ for a vector-valued parameter and wrote V for the volume of the
 6 model (Pitt et al., 2002). Fig. 14 shows the tangent surface at the equiprobable pdf $\vec{p}_e =$
 7 $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$: on the left parameterized in terms of pL, pM and on the right parameterized in
 8 terms of a, b .

9 The two models are equivalent as we can re-parameterize pL, pM to a, b and vice versa.
 10 This implies that the two models accommodate the same expected data patterns, hence, the
 11 two models are equal in size. In effect, we are free to choose the parametrization that allows
 12 us to calculate the size of the trinomial model conveniently as described by Eq. (37). What
 13 qualifies as convenient can be derived from the Fisher information matrix, which is given by:

$$I(pL, pM) = \frac{1}{1 - pL - pM} \begin{pmatrix} 1 - pM & 1 \\ 1 & 1 - pL \end{pmatrix}, \tag{38}$$

14 in terms of the parameters pL, pM , yielding $\det(I(pL, pM)) = \frac{pLpM - pL - pM}{(1 - pL - pM)^2}$. Consequently,
 15 the volume of the trinomial model is then given by

$$V_{\text{Tri}} = \int_0^1 \left(\int_0^{1-pL} \sqrt{\frac{pL \cdot pM - pL - pM}{(1 - pL - pM)^2}} dpM \right) dpL. \tag{39}$$

16 This integral is hard to compute as the inner integral depends on the outer integral, because
 17 the probabilities pL, pM and pR have to sum up to one. We prefer the Fisher information in
 18 the parametrization a, b as it leads to zero off-diagonal terms:

$$I(a, b) = \begin{pmatrix} \frac{1}{a(1-a)} & 0 \\ 0 & \frac{1-a}{b(1-b)} \end{pmatrix}. \quad (40)$$

1 Whenever the Fisher information matrix parametrized in terms of a, b has zero-off diagonals
 2 we then say that a, b are orthogonal parameters, see Fig. 14. Orthogonal parameters allow us
 3 to decouple the summing up to one restriction $pL + pR + pM = 1$ into two separate conditions
 4 $a \in (0, 1)$ and $b \in (0, 1)$, compare the integration bounds of Eq. (39) to those of Eq. (41). The
 5 volume of the cascading process model (thus, also of the trinomial model) is then given by

$$V_{\text{Tri}} = \int_0^1 \int_0^1 \frac{1}{\sqrt{ab(1-b)}} da db = \int_0^1 \frac{1}{\sqrt{a}} da \int_0^1 \frac{1}{\sqrt{b(1-b)}} db \stackrel{\text{Eq. (63)}}{=} \pi \int_0^1 \frac{1}{\sqrt{a}} da = 2\pi, \quad (41)$$

6 which equals an eighth of the surface area of a sphere of radius two, $\frac{1}{8}4\pi 2^2$, as one would
 7 expect. The size of the trinomial model can therefore be expressed as 2π .

8 **An Application of Fisher Information to Model Selection** To illustrate the merits
 9 of model selection using FIA over AIC and BIC, we introduce another model and compare
 10 it against the 2-binomial model. The new model is a version of the cascading process model
 11 where the parameter b is fixed to 0.5. Hence we refer to this model as the restricted cascading
 12 process model. According to the new model, participants only discriminate mixed from pure
 13 pairs of words and then randomly answer L or R . Fig. 15 shows the corresponding expected
 14 data patterns.

15 Note that the 2-binomial model and the restricted cascading process model are not
 16 nested and –because both models have only a single parameter– AIC and BIC discriminate
 17 the two models based on goodness-of-fit alone. Hence, when the observed data patterns lie
 18 near the right corner $P(X = M) = 1$, AIC and BIC will prefer the restricted cascading
 19 process model over the 2-binomial model. Conversely, when the observed data patterns lie
 20 near the bottom left corner $P(X = L) = 1$, AIC and BIC will prefer the 2-binomial model

1 over the restricted cascading process model.

2 Thus, AIC and BIC are unable to discriminate the 2-binomial model from the restricted
 3 cascading process model when the goodness-of-fit terms are close to each other, i.e., for data
 4 patterns near the point $\vec{p} = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4})$, where the two models intersect. Hence, whenever the
 5 likelihood and dimensions of the two models are similar, model selection by FIA then boils
 6 down to the assessment of the geometric complexity term, Eq. (30).

7 As noted above, the 2-binomial model has a volume term that equals $V_{2\text{Bin}} = \sqrt{2}\pi$,
 8 whereas the restricted cascading process model has a volume term given by $V_{\text{Restrict}} =$
 9 $\int_0^1 \frac{1}{\sqrt{a(1-a)}} da = \pi$ (for details see the online appendix). Thus, the restricted cascading process
 10 model is slightly less complex than the 2-binomial model, e.g., $\log(V_{\text{Restrict}}) < \log(V_{2\text{Bin}})$, and
 11 FIA will automatically prefer the simpler model whenever the goodness-of-fit is non-diagnostic
 12 (i.e., when the data are near the plus sign in Fig. 15).

13 **Summary of Section 3** In this section we showed how Fisher information can be used
 14 to quantify model complexity by measuring the size of model in the space of expected data
 15 patterns. This sophisticated conceptualization of model complexity takes into account the
 16 functional form of model parameters to strike an appropriate balance between parsimony and
 17 goodness-of-fit.

18 Concluding Comments

19 Fisher information is a central statistical concept that is of considerable relevance for
 20 mathematical psychologists. We illustrated the use of Fisher information in three different
 21 statistical paradigms: in the frequentist paradigm, Fisher information can be used to deter-
 22 mine the sample size required to estimate parameters at a target level of accuracy; in the
 23 Bayesian paradigm, Fisher information can be used to specify a default, translation-invariant
 24 prior distribution; finally, in the paradigm of information theory, data compression, and min-
 25 imum description length, Fisher information can be used to measure model complexity.

26 Our goal was to use concrete examples to provide more insight about Fisher information,
 27 something that may benefit psychologists who propose, develop, and compare mathematical

1 models for psychological processes. Our goal was not to provide a comprehensive or complete
2 treatment of all the uses of Fisher information throughout statistics. In our opinion, such a
3 treatment would require a book (e.g., Frieden, 2004) rather than a tutorial article.

4 Other usages of Fisher information are in the detection of model misspecification,
5 (Golden, 1995; Golden, 2000; Waldorp, Huizenga, & Grasman, 2005; Waldorp, 2009; Waldorp,
6 Christoffels, & van de Ven, 2011; White, 1982) and in the reconciliation of frequentist and
7 Bayesian estimation methods through the Bernstein–von Mises theorem (van der Vaart, 1998;
8 Bickel & Kleijn, 2012). In sum, Fisher information is a key concept in statistical modeling.
9 We hope to have provided an accessible and concrete tutorial article that explains the con-
10 cept and some of its uses for applications that are of particular interest to mathematical
11 psychologists.

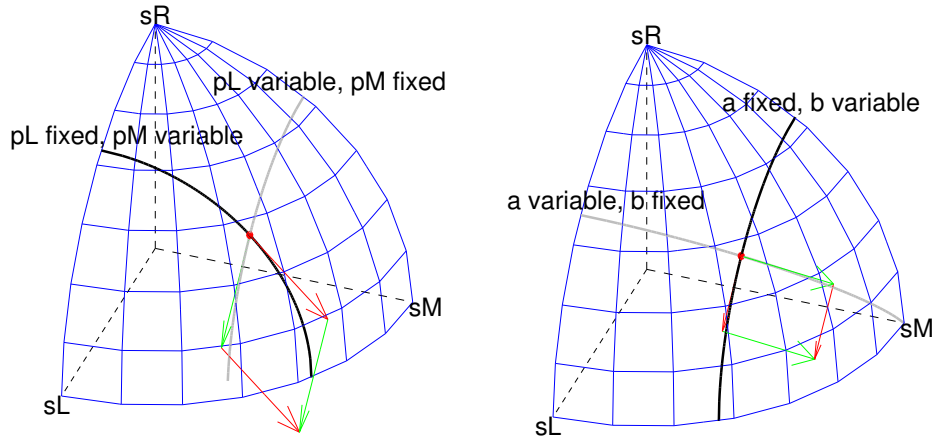


Figure 14. The size of the model is invariant under re-parametrization and we can therefore choose a parameterization which simplify the calculations. Recall that the size of an one-dimensional model, the arc length, is approximated by tangent lines, see Fig. 13. The size of a two-dimensional model on the other hand can be thought of an surface area that needs to be approximated by a tangent surface. This tangent surface can be decomposed into two tangent lines that can be derived by letting each parameter vary while the other stays fixed. The arrows in the left panel span up the tangent surface of the trinomial model with parameters pL, pM, pR at the dot that is represented by $pL = pM = pR = \frac{1}{3}$. The arrow that is directed to the left is tangent to the curve when pL is allowed to vary freely between zero and $1 - pM$, while $pM = \frac{1}{3}$ is set fixed. Similarly, the arrow that is directed to the right is tangent to the curve when pM is allowed to vary freely between zero and $1 - pR$, while $pR = \frac{1}{3}$ is set fixed. Analogously, the arrows in the right panel span up the tangent surface of the cascading process models with parameters a, b at the dot that is represented by $a = \frac{1}{3}, b = \frac{1}{2}$. The arrow that is directed to the left is tangent to the curve when a is allowed to vary freely between zero and one, while $b = \frac{1}{2}$ is set fixed. Similarly, the arrow that is directed to the right is tangent to the curve when b is allowed to vary freely between zero and one, while $a = \frac{1}{2}$ is set fixed. Note that the area of this tangent surface is easier to calculate due to the fact that the arrows are orthogonal to each other, which allows us simplify the calculations see Eq. (41).

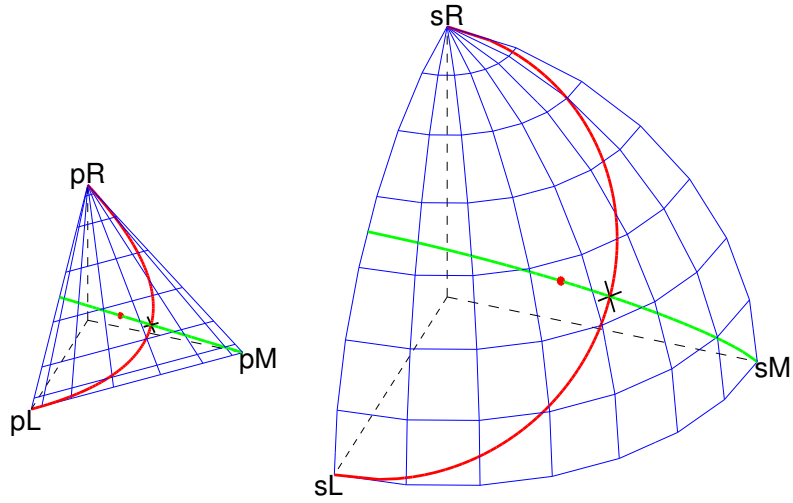


Figure 15. Expected data patterns from the 2-binomial model versus the restricted cascading process model. The latter model is represented by the line that traverses the surface starting from the bottom right corner. Note that it accommodates both the dot, i.e., the data pattern for which $pL = pR = pM = \frac{1}{3}$, and the cross that corresponds to $\theta = 0.5$ within the 2-binomial model. As before, the 2-binomial model is represented by the line that curves upwards starting from the bottom left corner.

References

- 1
- 2 Akaike, H. (1974). A new look at the statistical model identification. *Automatic Control, IEEE*
3 *Transactions on*, 19(6), 716–723.
- 4 Batchelder, W. H., & Riefer, D. M. (1980). Separation of storage and retrieval factors in free recall of
5 clusterable pairs. *Psychological Review*, 87, 375–397.
- 6 Batchelder, W. H., & Riefer, D. M. (1999). Theoretical and empirical review of multinomial process
7 tree modeling. *Psychonomic Bulletin & Review*, 6, 57–86.
- 8 Bayes, T., & Price, R. (1763). An essay towards solving a problem in the doctrine of chances. by
9 the late Rev. Mr. Bayes, FRS communicated by Mr. Price, in a letter to John Canton, AMFRS.
10 *Philosophical Transactions (1683-1775)*, 370–418.
- 11 Beran, R. (1995). The role of Hájek’s convolution theorem in statistical theory. *Kybernetika*, 31(3),
12 221–237.
- 13 Bickel, P. J., Klaassen, C. A., Ritov, Y., & Wellner, J. A. (1993). *Efficient and adaptive estimation*
14 *for semiparametric models*. Johns Hopkins University Press Baltimore.
- 15 Bickel, P. J., & Kleijn, B. J. K. (2012). The semiparametric Bernstein–von Mises theorem. *The Annals*
16 *of Statistics*, 40(1), 206–237.
- 17 Burnham, K. P., & Anderson, D. R. (2002). *Model selection and multimodel inference: A practical*
18 *information–theoretic approach (2nd ed.)*. New York: Springer Verlag.
- 19 Chechile, R. A. (1973). *The relative storage and retrieval losses in short-term memory as a function*
20 *of the similarity and amount of information processing in the interpolated task*. Unpublished
21 doctoral dissertation, University of Pittsburgh.
- 22 Cramér, H. (1946). Methods of mathematical statistics. *Princeton University Press*, 23.
- 23 DasGupta, A. (2011). *Probability for statistics and machine learning*. Springer.
- 24 Diaconis, P., & Ylvisaker, D. (1979). Conjugate priors for exponential families. *The Annals of*
25 *Statistics*, 7(2), 269–281.
- 26 Edwards, A. W. F. (1992). *Likelihood*. Baltimore, MD: The Johns Hopkins University Press.
- 27 Fisher, R. A. (1912). On an absolute criterion for fitting frequency curves. *Messenger of Mathematics*,
28 41, 155–160.

- 1 Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transac-*
2 *tions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical*
3 *Character*, 222, 309–368.
- 4 Fisher, R. A. (1925). Theory of statistical estimation. In *Mathematical Proceedings of the Cambridge*
5 *Philosophical Society* (Vol. 22, pp. 700–725).
- 6 Fisher, R. A. (1930). Inverse probability. In *Mathematical Proceedings of the Cambridge Philosophical*
7 *Society* (Vol. 26, pp. 528–535).
- 8 Fréchet, M. (1943). Sur l’extension de certaines évaluations statistiques au cas de petits échantillons.
9 *Revue de l’Institut International de Statistique*, 182–205.
- 10 Frieden, B. R. (2004). *Science from Fisher information: A unification*. Cambridge University Press.
- 11 Gamerman, D., & Lopes, H. F. (2006). *Markov chain Monte Carlo: Stochastic simulation for Bayesian*
12 *inference*. Boca Raton, FL: Chapman & Hall/CRC.
- 13 Ghosh, J. (1985). Efficiency of estimates—part I. *Sankhyā: The Indian Journal of Statistics, Series A*,
14 310–325.
- 15 Gilks, W. R., Richardson, S., & Spiegelhalter, D. J. (Eds.). (1996). *Markov chain Monte Carlo in*
16 *practice*. Boca Raton (FL): Chapman & Hall/CRC.
- 17 Golden, R. M. (1995). Making correct statistical inferences using the wrong probability model. *Journal*
18 *of Mathematical Psychology*, 39, 3–20.
- 19 Golden, R. M. (2000). Statistical tests for comparing possibly misspecified and nonnested models.
20 *Journal of Mathematical Psychology*, 44(1), 153–170.
- 21 Grünwald, P. (2007). *The minimum description length principle*. Cambridge, MA: MIT Press.
- 22 Grünwald, P., Myung, I. J., & Pitt, M. A. (Eds.). (2005). *Advances in minimum description length:*
23 *Theory and applications*. Cambridge, MA: MIT Press.
- 24 Hájek, J. (1970). A characterization of limiting distributions of regular estimates. *Zeitschrift für*
25 *Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 14(4), 323–330.
- 26 Huzurbazar, V. S. (1949). On a property of distributions admitting sufficient statistics. *Biometrika*,
27 36(1-2), 71–74.
- 28 Inagaki, N. (1970). On the limiting distribution of a sequence of estimators with uniformity property.
29 *Annals of the Institute of Statistical Mathematics*, 22(1), 1–13.

- 1 Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proceedings of*
2 *the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007), 453–461.
- 3 Jeffreys, H. (1961). *Theory of probability* (3rd ed.). Oxford, UK: Oxford University Press.
- 4 Kass, R. E. (1989). The geometry of asymptotic inference. *Statistical Science*, 188–219.
- 5 Klauer, K. C., & Kellen, D. (2011). The flexibility of models of recognition memory: An analysis by
6 the minimum-description length principle. *Journal of Mathematical Psychology*, 55(6), 430–450.
- 7 Laplace, P. S. (1986). Memoir on the probability of the causes of events. *Statistical Science*, 1(3),
8 364–378.
- 9 Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian modeling for cognitive science: A practical course*.
10 Cambridge University Press.
- 11 Myung, I. J. (2003). Tutorial on maximum likelihood estimation. *Journal of Mathematical Psychology*,
12 47, 90–100.
- 13 Myung, I. J., Balasubramanian, V., & Pitt, M. (2000). Counting probability distributions: Differential
14 geometry and model selection. *Proceedings of the National Academy of Sciences*, 97(21), 11170–
15 11175.
- 16 Myung, I. J., Forster, M. R., & Browne, M. W. (2000). Model selection [Special issue]. *Journal of*
17 *Mathematical Psychology*, 44(1–2).
- 18 Pitt, M., Myung, I. J., & Zhang, S. (2002). Toward a method of selecting among computational
19 models of cognition. *Psychological Review*, 109(3), 472–491.
- 20 Raftery, A. E. (1995). Bayesian model selection in social research. In P. V. Marsden (Ed.), *Sociological*
21 *methodology* (pp. 111–196). Cambridge: Blackwells.
- 22 Rao, C. R. (1945). Information and accuracy attainable in the estimation of statistical parameters.
23 *Bulletin of the Calcutta Mathematical Society*, 37(3), 81–91.
- 24 Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85, 59–108.
- 25 Rissanen, J. (1996). Fisher information and stochastic complexity. *IEEE Transactions on Information*
26 *Theory*, 42, 40–47.
- 27 Roberts, S., & Pashler, H. (2000). How persuasive is a good fit? A comment on theory testing in
28 psychology. *Psychological Review*, 107, 358–367.
- 29 Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6, 461–464.

- 1 Stevens, S. S. (1957). On the psychophysical law. *Psychological Review*, *64*(3), 153–181.
- 2 Stigler, S. (1986). Laplace’s 1774 memoir on inverse probability. *Statistical Science*, *1*(3), 359–363.
- 3 van der Vaart, A. W. (1998). *Asymptotic statistics*. Cambridge University Press.
- 4 Wagenmakers, E.-J., & Waldorp, L. (2006). Model selection: Theoretical developments and applica-
5 tions [Special issue]. *Journal of Mathematical Psychology*, *50*(2).
- 6 Waldorp, L. (2009). Robust and unbiased variance of GLM coefficients for misspecified autocorrelation
7 and hemodynamic response models in fMRI. *International Journal of Biomedical Imaging*, *2009*,
8 723912.
- 9 Waldorp, L., Christoffels, I., & van de Ven, V. (2011). Effective connectivity of fMRI data using
10 ancestral graph theory: Dealing with missing regions. *NeuroImage*, *54*(4), 2695 – 2705.
- 11 Waldorp, L., Huizenga, H., & Grasman, R. (2005). The Wald test and Cramér–Rao bound for
12 misspecified models in electromagnetic source analysis. *IEEE Transactions on Signal Processing*,
13 *53*(9), 3427–3435.
- 14 White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica*, *50*(1), 1–25.
- 15 Wu, H., Myung, I. J., & Batchelder, W. H. (2010). Minimum description length model selection of
16 multinomial processing tree models. *Psychonomic Bulletin & Review*, *17*, 275–286.

Appendix A

Calculating the MLE and Fisher information for the normal distribution

For the normal distribution we have the following log-likelihood functions.

$$\log f(\vec{x} | \vec{\theta}) = -\frac{n}{2} \log(2\pi v) - \frac{1}{2v} \sum_{i=1}^n (X_i - \mu)^2, \quad (42)$$

where we have written $v = \sigma^2$ to avoid confusion of differentiation with respect to the standard deviation. To derive the maximum likelihood functions we are required to solve the likelihood equations. Hence, find that parameter value such that the following partial derivatives are zero:

$$\frac{\partial}{\partial \mu} \log f(\vec{X} | \hat{\mu}_n, v) = \frac{1}{2v} \sum_{i=1}^n (X_i - \mu) \quad (43)$$

$$\frac{\partial}{\partial v} \log f(\vec{X} | \hat{\mu}_n, v) = -\frac{n}{2v} + \frac{1}{2v^2} \sum_{i=1}^n (X_i - \mu)^2 \quad (44)$$

To calculate the Fisher information we need to calculate the second derivatives, the Hessian matrix, and integrate these with respect to \vec{X} :

$$\frac{\partial^2}{\partial \mu \partial \mu} \log f(\vec{X} | \mu, v) = -\frac{1}{v}, \quad (45)$$

$$\frac{\partial^2}{\partial \mu \partial v} \log f(\vec{X} | \mu, v) = -\frac{1}{2v^2} \sum_{i=1}^n (X_i - \mu), \quad (46)$$

$$\frac{\partial^2}{\partial v \partial v} \log f(\vec{X} | \mu, v) = \frac{n}{2v^2} - \frac{1}{v^3} \sum_{i=1}^n (X_i - \mu)^2. \quad (47)$$

We use these equations to show:

Theorem 3 (Maximum likelihood estimators for normal distributions). *If $\vec{X} = (X_1, \dots, X_n)$ are i.i.d. with $X_1 \sim \mathcal{N}(\mu, \sigma^2)$, then the MLEs are given by $\hat{\mu}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, $\hat{\sigma}_n^2 = S_n =$*

1 $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. Furthermore, the expected Fisher information matrix is then given by

$$I_n(\vec{\theta}) = n \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}, \quad (48)$$

2 where the matrix on the right-hand side represents the unit information of a normally dis-
 3 tributed random variable \vec{X} . ◇

4 *Proof.* For the MLE of μ , we set (43) to zero, which yields $\sum_{i=1}^n X_i = n\mu$. Hence, $\hat{\mu}_n$ is given
 5 by the sample mean. We use this knowledge to derive the MLE for the variance v . Setting
 6 (44) to zero yields $n = \frac{1}{v} \sum_{i=1}^n (X_i - \mu)^2$ after we multiplied both sides by $2v$ and plugging in
 7 the MLE for μ yields S_n .

8 To compute the Fisher information matrix, we have to integrate the negatives of (45),
 9 (46) and (47) with respect to \vec{X} , see (9). First note that (45) does not depend on \vec{X} so it is
 10 straightforward to see that $I(\mu) = \frac{1}{v}$. Furthermore, as (46) defines an uneven function about
 11 μ , it is as much positive as negative on each side of μ , it will integrate to zero. By i.i.d. we
 12 know that the integration on the right-hand side of (47) times -1 yields $\frac{nv}{v^3} - \frac{n}{2v^2} = n\frac{1}{2v^2}$. □

Appendix B

Calculating the mean and variance of $\hat{\sigma}_n^2$ for the normal distribution

14 **Theorem 4.** If $\vec{X} = (X_1, \dots, X_n)$ are i.i.d. with $X_1 \sim \mathcal{N}(\mu, \sigma^2)$ with an MLE $\hat{\sigma}_n^2$ and the
 15 standard sample variance $S_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$, we have:

$$E(\hat{\sigma}_n^2) = \left(1 - \frac{1}{n}\right) E(S_{n-1}^2) \text{ and } E(S_{n-1}^2) = \sigma^2 \quad (49)$$

$$\text{Var}(\hat{\sigma}_n^2) = \left(1 - \frac{1}{n}\right)^2 \text{Var}(S_{n-1}^2) < \text{Var}(S_{n-1}^2) = \frac{2\sigma^4}{n-1} \quad (50)$$

16 *Proof.* First note that as n is known $1 - \frac{1}{n}$ can be considered as a constant and that $\hat{\sigma}_n^2 =$
 17 $\frac{n-1}{n} S_{n-1}^2$. This explains the factors on the left-hand side of both results. Furthermore, the
 18 result (49) follows directly from the linearity of the expectation and the fact that the sample

1 variance S_{n-1}^2 is unbiased, $E(S_{n-1}^2) = \sigma^2$. By Basu's theorem, (Ch. 18 DasGupta, 2011)
 2 we can show that $\frac{n-1}{\sigma^2}S_{n-1}^2 \sim \chi^2(n-1)$ is distributed according to a chi-squared distribution
 3 with $n-1$ degrees of freedom and thus a variance of $2(n-1)$. Hence,

$$2(n-1) = \text{Var}\left(\frac{n-1}{\sigma^2}S_{n-1}^2\right) = \frac{n-1}{\sigma^2}\text{Var}(S_{n-1}^2), \quad (51)$$

4 and the result (50) then follows directly from multiplying (51) with $\frac{\sigma^2}{n-1}$ on both sides. \square

Appendix C

Angle coin

6 Let $P(X=1) = \theta$ and ϕ is the angle of a bent coin with respect to the horizontal axis.

$$h : \Phi \rightarrow \Theta : \phi \mapsto \frac{1}{2} + \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3 \quad \text{with} \quad \frac{d\theta}{d\phi} = \frac{3}{2\pi^3}\phi^2 \quad (52)$$

7 and an inverse function

$$k = h^{-1} : \Theta \rightarrow \Phi : \theta \mapsto \begin{cases} -\pi(1-2\theta)^{\frac{1}{3}} & \text{with } \frac{d\phi}{d\theta} = \frac{2\pi}{3(1-2\theta)^{\frac{2}{3}}} \quad \text{when } 0 \leq \theta < 0.5 \\ 0 & \text{with } \frac{d\phi}{d\theta} = \frac{2\pi}{3(1-2\theta)^{\frac{2}{3}}} \quad \text{when } \theta = 0.5 \\ \pi(2\theta-1)^{\frac{1}{3}} & \text{with } \frac{d\phi}{d\theta} = \frac{2\pi}{3(2\theta-1)^{\frac{2}{3}}} \quad \text{when } 0.5 < \theta \leq 1 \end{cases} \quad (53)$$

8 For $I(\phi)$, write

$$\log p(x|\phi) = x \log \left(\frac{1}{2} + \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3 \right) + (1-x) \log \left(\frac{1}{2} - \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3 \right) \quad (54)$$

9 Furthermore,

$$\frac{d}{d\phi} \log p(x | \phi) = \frac{3}{\pi^3} \left(x \frac{\phi^2}{1 + \left(\frac{\phi}{\pi}\right)^3} - (1-x) \frac{\phi^2}{1 - \left(\frac{\phi}{\pi}\right)^3} \right) \quad (55)$$

1 Moreover, (this is where it gets exciting)

$$\frac{d^2}{d\phi^2} \log p(x | \phi) = \frac{3}{\pi^3} \left(x \frac{2\phi - \frac{\phi^4}{\pi^3}}{\left(1 + \left(\frac{\phi}{\pi}\right)^3\right)^2} - (1-x) \frac{2\phi + \frac{\phi^4}{\pi^3}}{\left(1 - \left(\frac{\phi}{\pi}\right)^3\right)^2} \right) \quad (56)$$

2 Hence,

$$I(\phi) = \frac{9\phi^4}{\pi^6 - \phi^6} \quad (57)$$

3 Or, based on the fact that the Fisher information in basis θ is given by $I(\theta) = \frac{1}{\theta(1-\theta)}$:

$$I(\phi) = \left(\frac{d\theta}{d\phi} \right)^2 I(\theta(\phi)) = \left(\frac{3\phi^2}{2\pi^3} \right)^2 \frac{1}{\left(\frac{1}{2} + \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3\right) \left(\frac{1}{2} - \frac{1}{2} \left(\frac{\phi}{\pi}\right)^3\right)} \quad (58)$$

4 To calculate the normalizing constant for Jeffreys' prior we have to integrate $\int_{-\pi}^{\pi} \sqrt{I(\phi)} d\phi$.

$$\int_{-\pi}^{\pi} \sqrt{I(\phi)} d\phi = \int_{-\pi}^{\pi} \sqrt{\frac{9\phi^4}{\pi^6 - \phi^6}} d\phi \quad (59)$$

$$= \int_{-\pi}^{\pi} \frac{3\phi^2}{\sqrt{\pi^6 - \phi^6}} \frac{1}{\sqrt{\pi^6}} d\phi \quad (60)$$

$$= \int_{-\pi}^{\pi} \frac{1}{\sqrt{1 - \left(\frac{\phi}{\pi}\right)^6}} \frac{3\phi^2}{\pi^3} d\phi. \quad (61)$$

5 If we substitute $s = \left(\frac{\phi}{\pi}\right)^3$ we then have $ds = \frac{3\phi^2}{\pi^3} d\phi$. Furthermore, when $\phi = \pi$ this means

6 that $s = 1$ and we get $s = -1$ as a lower bound. Hence,

$$\int_{-\pi}^{\pi} \sqrt{I(\phi)} d\phi = \int_{-1}^1 \frac{1}{\sqrt{1-s^2}} ds, \quad (62)$$

$$= \sin^{-1}(s) \Big|_{s=-1}^{s=1} = \frac{\pi}{2} - \frac{-\pi}{2} = \pi, \quad (63)$$

1 where we have used a standard calculus result to derive the identity with \sin^{-1} . Similarly,
 2 the normalizing constant for the Jeffreys prior in θ basis is given by

$$\int_0^1 \sqrt{I(\theta)} d\theta = \int_0^1 \frac{1}{\sqrt{\theta(1-\theta)}} d\theta. \quad (64)$$

3 Substituting $\theta = \sin^2(z)$ we have $d\theta = 2 \sin(z) \cos(z)$ yields the result. Note that Jeffreys'
 4 prior on θ is in fact Beta(0.5, 0.5) and since a beta prior leads to a beta posterior we know
 5 that the posterior for θ is then given by:

$$\frac{1}{B(7.5, 3.5)} \theta^{6.5} (1-\theta)^{2.5}. \quad (65)$$

6 Similarly, we then have this yields the following

$$\frac{1}{B(7.5, 3.5)} \left(\frac{1}{2} + \frac{1}{2} \left(\frac{\phi}{\pi} \right)^3 \right)^{6.5} \left(\frac{1}{2} - \frac{1}{2} \left(\frac{\phi}{\pi} \right)^3 \right)^{2.5} \frac{3\phi^2}{2\pi^3}, \quad (66)$$

7 for ϕ .

8

Appendix D

Computations for Section 3

9 **Fisher information for the trinomial model** By definition we can derive the Fisher
 10 information by the second partial derivatives of the log-likelihood functions. The first partial
 11 derivatives for each single outcome are given by:

$$\begin{aligned}
\frac{\partial}{\partial pL} \log P(X = L) &= \frac{1}{pL} & \frac{\partial}{\partial pM} \log P(X = L) &= 0 \\
\frac{\partial}{\partial pL} \log P(X = M) &= 0 & \frac{\partial}{\partial pM} \log P(X = M) &= \frac{1}{pM} \\
\frac{\partial}{\partial pL} \log P(X = R) &= \frac{-1}{1-pL-pM} & \frac{\partial}{\partial pM} \log P(X = R) &= \frac{-1}{1-pL-pM}
\end{aligned} \tag{67}$$

1 The non-zero unmixed second partial derivatives for each single outcome are then given by

$$\begin{aligned}
\frac{\partial^2}{\partial pL^2} \log P(X = L) &= \frac{-1}{pL^2} & \frac{\partial^2}{\partial pM^2} \log P(X = M) &= \frac{-1}{pM^2} \\
\frac{\partial^2}{\partial pL^2} \log P(X = R) &= \frac{-1}{(1-pL-pM)^2} & \frac{\partial^2}{\partial pM^2} \log P(X = R) &= \frac{-1}{(1-pL-pM)^2}
\end{aligned} \tag{68}$$

2 Hence, taking the expectation over minus the left vector of (68) yields

$$I_{1,1} = I_{pL,pL} = \frac{1}{pL^2}pL + \frac{1}{(1-pL-pM)^2}(1-pL-pM) = \frac{1-pM}{1-pL-pM} \tag{69}$$

3 Hence, taking the expectation over minus the right vector of (68) yields

$$I_{2,2} = I_{pM,pM} = \frac{1}{pM^2}pM + \frac{1}{(1-pL-pM)^2}(1-pL-pM) = \frac{1-pL}{1-pL-pM} \tag{70}$$

4 The only non-zero mixed second partial derivatives with respect to both arguments pL, pM
5 is given by

$$\frac{\partial^2}{\partial pL \partial pM} \log P(X = R) = \frac{-1}{(1-pL-pM)^2} \tag{71}$$

6 which yields

$$I_{2,1} = I_{1,2} = I_{pL,pM} = \frac{1}{(1-pL-pM)^2}(1-pL-pM) \tag{72}$$

- 1 **Fisher information for the cascading process model** By definition we can derive the
 2 Fisher information by the second partial derivatives of the log-likelihood functions. The first
 3 partial derivatives for each single outcome are given by:

$$\begin{aligned}
 \frac{\partial}{\partial a} \log P(X = L) &= \frac{-1}{1-a} & \frac{\partial}{\partial b} \log P(X = L) &= \frac{1}{b} \\
 \frac{\partial}{\partial a} \log P(X = M) &= \frac{1}{a} & \frac{\partial}{\partial b} \log P(X = M) &= 0 \\
 \frac{\partial}{\partial a} \log P(X = R) &= \frac{-1}{1-a} & \frac{\partial}{\partial b} \log P(X = R) &= \frac{-1}{1-b}
 \end{aligned} \tag{73}$$

- 4 The unmixed second partial derivatives with respect to both a, b for each single outcome are
 5 then given by

$$\begin{aligned}
 \frac{\partial^2}{\partial a^2} \log P(X = L) &= \frac{-1}{(1-a)^2} & \frac{\partial^2}{\partial b^2} \log P(X = L) &= \frac{-1}{b^2} \\
 \frac{\partial^2}{\partial a^2} \log P(X = M) &= \frac{-1}{a^2} & \frac{\partial^2}{\partial b^2} \log P(X = M) &= 0 \\
 \frac{\partial^2}{\partial a^2} \log P(X = R) &= \frac{-1}{(1-a)^2} & \frac{\partial^2}{\partial b^2} \log P(X = R) &= \frac{-1}{(1-b)^2}
 \end{aligned} \tag{74}$$

- 6 Hence, taking the expectation over minus the left vector of (74) yields

$$I_{1,1} = I_{a,a} = \frac{1}{(1-a)^2}(1-a)b + \frac{1}{a^2}a + \frac{1}{(1-a)^2}(1-a)(1-b) = \frac{1}{a(1-a)} \tag{75}$$

- 7 Hence, taking the expectation over minus the right vector of (74) yields

$$I_{2,2} = I_{b,b} = \frac{1}{(b)^2}(1-a)b + \frac{1}{(1-b)^2}(1-a)(1-b) = \frac{1-a}{1-b} \tag{76}$$

- 8 All mixed second partial derivatives are zero. Hence, $I_{1,2} = I_{2,1} = 0$