



Valencia 2010

# Testing Indistinguishable Hypotheses?

★

Jan-Willem Romeijn  
Faculty of Philosophy  
University of Groningen

# **Contents**

<b>1</b>	<b>Conjunction fallacy</b>	<b>3</b>
<b>2</b>	<b>Model selection</b>	<b>8</b>
<b>3</b>	<b>Comparing causal models</b>	<b>13</b>
<b>4</b>	<b>Convergence measure</b>	<b>17</b>
<b>5</b>	<b>Future research</b>	<b>20</b>

# 1 Conjunction fallacy

Linda is 31 years old, unmarried, assertive, and intelligent. She studied philosophy and wrote her thesis on social issues and justice. She was active in the campaign against the war in Iraq.

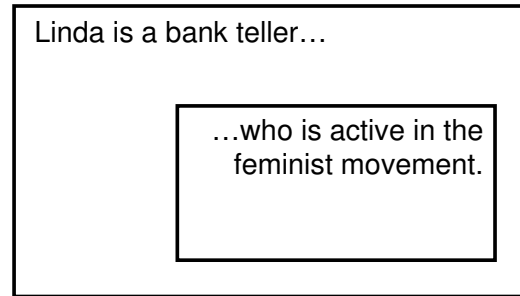


Which of the following two statements is more probable?

1. Linda is a bank teller.
2. Linda is a bank teller who is active in the feminist movement.

## Probability and set inclusion

Statisticians will not give the answer that is often given by the probabilistically naive.



Probability is a measure over sets, and the set of worlds in which Linda is a feminist bank teller is strictly included in the set of worlds in which she is a bank tellers.

## Bayesian model selection

The idea behind BMS is to assign probabilities to models, and compare these models by means of the marginal, or average, likelihoods:

$$\frac{P(M_1|x)}{P(M_0|x)} = \frac{P(x|M_1) P(M_1)}{P(x|M_0) P(M_0)}$$

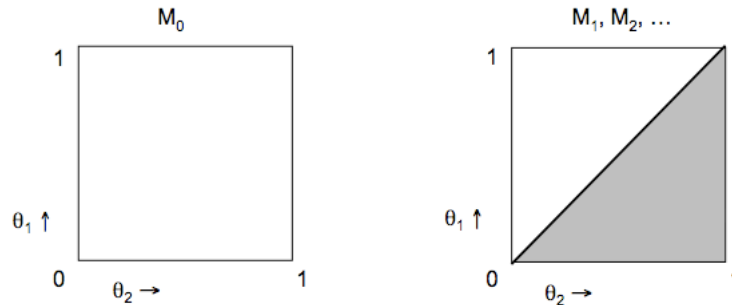
where

$$P(x|M_i) = \int_{M_i} P(x|\theta)P(\theta)d\theta.$$

For want of a prior over models, the likelihood ratio is often taken as the sole guide to choosing the model.

## Model selection and Linda

It seems that Bayesian model selection is prone to committing the conjunction fallacy.



The probability of the constrained model may get larger than the probability of the unconstrained model. What do those probabilities refer to?

## **Bayesian information criterion**

The BIC is an approximation to the marginal log-likelihood of a model in the long run:

$$P(x|M_i) \sim -\log P(x|\hat{\theta}_x) + d \log n$$

in which  $\hat{\theta}_x$  is the maximum likelihood (ML) estimation of  $\theta$  for the data  $x$ , and  $n$  is the size of the data set  $x$ . It seems that the BIC is prone to committing the same conjunction fallacy as BMS.

## 2 Model selection

Are other model selection tools also prone to the conjunction fallacy? And if not, can we understand BIC along the lines of the other ICs?

**AIC** We estimate the expected distance between the truth and the ML-estimation in the model.

**DIC** We estimate the expected predictive accuracy of the ML-estimation in the model.

These selection tools employ some penalty term, typically the number of model parameters  $d$ , as a measure of complexity. Note that prima facie this makes choosing between inequality constrained models different.

## AIC

The general idea of the AIC is that we determine the estimated distance between the true chance  $\theta^*$  and the chance  $\theta$  estimated in the model.

$$\Delta(\theta^*, \theta) \sim E_{P(x|\theta^*)}[-\log P(x|\theta)]$$

In this expression  $\theta$  can still range over the whole model, while distances are of course between two points. This is resolved by taking an expectation towards a dummy data set  $y$ .

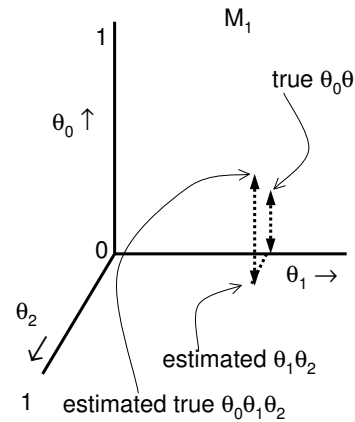
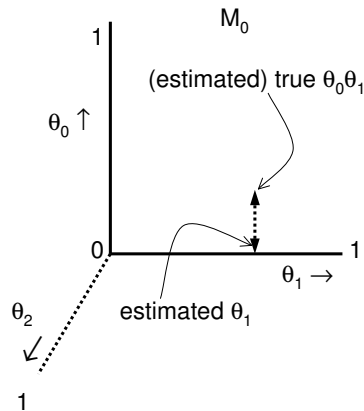
$$\Delta(\theta^*, \theta) \sim E_{P(y|\theta^*)} [E_{P(x|\theta^*)}[-\log P(x|\hat{\theta}_y)]]$$

In integrating this data set out, the number of dimensions appears:

$$\Delta(\theta^*, \theta) \sim E_{P(x|\theta^*)}[-\log P(x|\hat{\theta}_x)] + d \sim \log P(x|\hat{\theta}_x) + d$$

## AIC continued

Formulas are perhaps not the best way of explaining the procedure. . . Here is an artist's impression of how the distance to an unknown truth  $\theta^*$  is estimated.



## DIC

In the DIC we determine a Bayesian expectation for the parameter  $\bar{\theta}_y$ , and we see how well this expected value predicts data from the true distribution,  $P(x|\theta^*)$ .

$$E_{P(x|\theta^*)}[-\log P(x|\bar{\theta}_y)] \sim -\log P(y|\bar{\theta}_y) + d$$

The predictive accuracy is measured by the logarithmic loss function:

$$-\log P(x|\theta)$$

This is also only term that matters to the distance function featuring in the AIC. Although different in interpretation, the loss function expressing predictive accuracy and the distance to the truth may be used interchangeably.

## **BMS**

Clearly AIC and DIC do not commit the conjunction fallacy. But can we frame the BIC as the other ICs? There is some similarity in the use of average likelihoods:

$$-\log P(y|M) \sim E_{P(\theta)}[-\log P(y|\theta)].$$

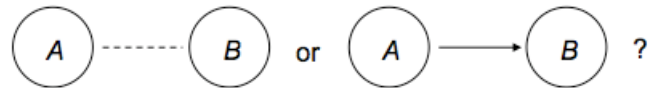
But this is a far cry from presenting BMS on a par with the other ICs.

- The true parameter value  $\theta^*$  is not alluded to in the BMS.
- BMS does not rest on an expression involving all possible data  $x$  but only the data actually obtained.

In the following we are going to bite the bullet and see what else the probabilities featuring in BMS might mean.

### 3 Comparing causal models

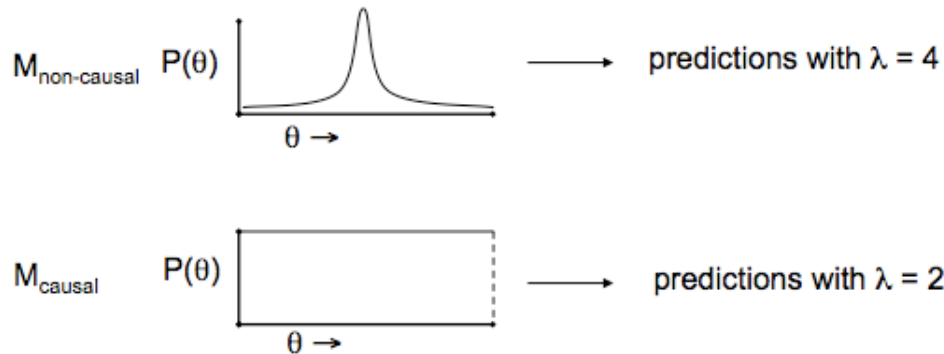
Say that we compare the following causal hypotheses about two binary variables:



The non-causal model has a so-called simplex as its parameter space, with four independent probabilities adding up to one. The causal model is represented by a parameter for the probability of A, and two parameters for the conditional probabilities of B on A.

## Different uniform priors

The different parameterisations, associated with a different causal picture, motivate different uniform priors over what is essentially the same parameter space.



## Different predictive properties

Depending on the prior, we find predictions that approach the true parameter value more or less quickly.

$$P(y_{t+1} = 1 | y_1 \cdots y_t) = \frac{\lambda}{t + \lambda} \frac{1}{2} + \frac{t}{t + \lambda} \frac{t_1}{t}.$$

The same expression also captures the expected values for the parameter, because

$$P(y_{t+1} = 1 | \theta, y_1 \cdots y_t) = \theta \quad P(y_{t+1} = 1 | y_1 \cdots y_t) = \int \theta P(\theta) d\theta = \bar{\theta}_y.$$

The marginal likelihoods thus capture how the expected value for  $\theta$  approaches the true value.

### An example

For the causal hypotheses at hand, we can derive the following likelihood ratio for the two models as follows:

$$BF = \frac{P(y_1 \cdots y_t | M_{\text{causal}})}{P(y_1 \cdots y_t | M_{\text{non-causal}})} = \frac{6(t_0 + 1)(t_1 + 1)}{(t + 2)(t + 3)}$$

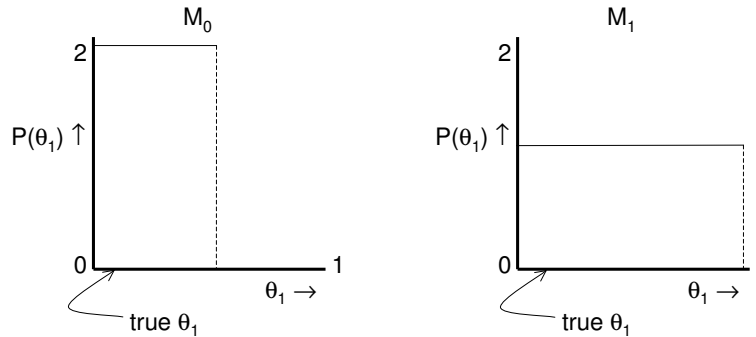
We find the following interesting points:

Number of observations $t$	Interval in which $BF > 1$
$< 12$	-
12	$\frac{1}{2}$
48	$[\frac{1}{4}, \frac{3}{4}]$
$\infty$	$[\frac{1}{2} - \frac{1}{2\sqrt{3}}, \frac{1}{2} + \frac{1}{2\sqrt{3}}]$

This means that with less than 12 observations, the causal model always performs better! What is going on?

## 4 Convergence measure

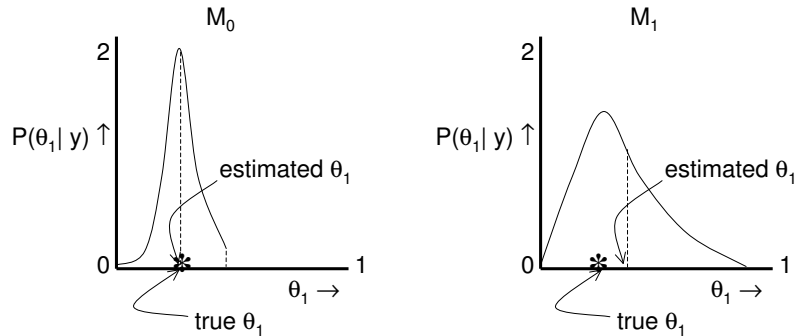
In the foregoing example, the Bayes-factor measures how fast the expected value of the parameter  $\bar{\theta}_y$  approaches the true value  $\theta^*$ .



We can view BMS in this way more generally: we are comparing the priors on their convergence properties.

## Choosing the best expectation

Both models generate an expected value for the parameter value  $\bar{\theta}_y$ .



If  $\theta^*$  lies within the restricted region, the expected value of the restricted model will be closer to the true value.

## **What marginal likelihood measures**

Marginal likelihood thus combine the aforementioned aspects of model selection in a particular way.

- The Bayes' factors measure the online predictive performance of the models. This performance is determined by
  - the probability distributions in the model, and
  - the prior probability over the model.
- The online predictive performance is a measure for the distance to the true parameter value: it expresses the speed of approaching the maximum likelihood estimation.

However, the performance and distance are determined by the data that we have obtained, and not by all data that we could have got under the true distribution.

## 5 Future research

This leaves many questions unanswered. . .

- Can we derive a BIC for equi-dimensional models that only differ in the prior defined over them? For inequality-constrained priors the answer is affirmative, but the general case is unclear.
- Following the BIC, we can still view BMS as generating a kind of penalty term for complexity, although different from the penalty that the ICs give. But what exactly does it penalise for?
- We cannot interpret posterior model probabilities (PMPs) as probabilities that the true parameter value is included in the region. Is the concept of PMP terminally ill?