

# Computational Cognitive Science

## Lecture 8: Model comparison and selection 2

Benjamin Peters

School of Informatics

University of Edinburgh

October 11, 2024

# Readings

- Chapter 10.6 of F&L
- Chapter 11 of F&L

Optional:

- “Bayesian hypothesis testing for psychologists: A tutorial on the Savage–Dickey method” (2010) by Wagenmakers et al. ([Link](#))

# Model comparison

- We discussed some methods for comparing models using likelihoods under MLEs, and predictive accuracy

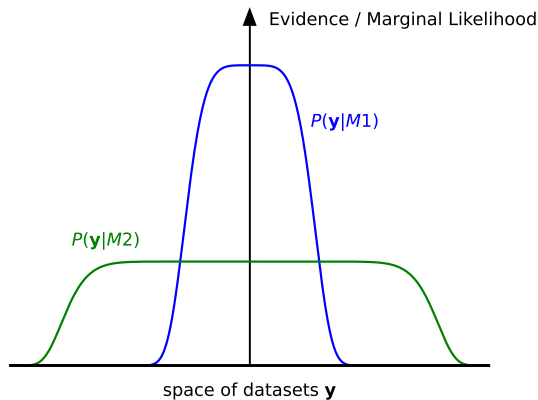
# Bayesian model comparison

Last time we discussed approaches to model comparison that don't involve marginal likelihoods.

Today we'll talk about using and approximating marginal likelihoods:

$$p(\mathbf{y}|\mathcal{M}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M})p(\boldsymbol{\theta}|\mathcal{M})d\boldsymbol{\theta}$$

# Bayesian model comparison



- More complex models spread their probability across more possible outcomes.
- Data  $\mathbf{y}$  that both models can predict, tends to have higher marginal likelihood (evidence) under the simpler model.

## Bayesian model comparison

Even if we have marginal likelihoods under several models, we can't compute  $P(\mathcal{M}|\mathbf{y})$  ( $P(\mathcal{M}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathcal{M})P(\mathcal{M})}{p(\mathbf{y})}$ ):

We would need to sum over all possible models and data distributions for all of them. Instead:

- Given any two models  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , we can compute the ratio of their posterior probabilities:

$$\frac{P(\mathcal{M}_1|\mathbf{y})}{P(\mathcal{M}_2|\mathbf{y})} = \frac{P(\mathbf{y}|\mathcal{M}_1)P(\mathcal{M}_1)}{P(\mathbf{y}|\mathcal{M}_2)P(\mathcal{M}_2)}$$

The intractable/unknown normalization terms cancel out, and we get the *relative* probabilities of our models.

# Bayesian model comparison

There is unlikely to be consensus about  $P(\mathcal{M})$ ; in practice people use *Bayes factors* (BF):

$$\underbrace{\frac{P(\mathcal{M}_1|\mathbf{y})}{P(\mathcal{M}_2|\mathbf{y})}}_{\text{posterior odds}} = \underbrace{\frac{P(\mathbf{y}|\mathcal{M}_1)}{P(\mathbf{y}|\mathcal{M}_2)}}_{\text{BF}} \cdot \underbrace{\frac{P(\mathcal{M}_1)}{P(\mathcal{M}_2)}}_{\text{prior odds}}$$

- BF is the posterior ratio given equal prior probabilities
- How strongly you'd have to prefer a model *a priori* in order to (still) favor it *a posteriori*

Lots of opinions about what constitutes a “convincing” Bayes factor

# Estimating marginal likelihoods

If we have no closed-form solution for a marginal likelihood, what can we do?

We have several options, including:

- ① Numerical integration
- ② Importance sampling
- ③ ~~Harmonic mean estimation~~
- ④ Transdimensional MCMC
- ⑤ Savage-Dickey density ratio
- ⑥ BIC



# Numerical integration

Use a general-purpose algorithm to integrate a function within a hypercube.

- Easy!
- Requires bounds on the high-density parts of the space
- Intractable in high-dimensional spaces
- Risks missing narrow peaks

## Numerical integration: Example

We can use standard probability densities to test our methods, since we know their integral (over **data**, not parameters) is 1.

Recall that the normal density function is  $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\mu-x)^2}{2\sigma^2}}$ .

Let's approximate the integral of  $e^{-\frac{(\mu-x)^2}{2\sigma^2}}$ , which is the reciprocal of the normalizing constant  $\frac{1}{\sqrt{2\pi\sigma^2}}$ .

## Numerical integration: Example

```
library(cubature)
sd=10;mu=.5
unnGauss <- function(x) {exp(-(mu-x)^2/(2*sd^2))}
adaptIntegrate(unnGauss,c(-1E3),c(1E3))
```

Result:

```
$integral
[1] 25.06628

> sqrt(2*pi*sd^2)
[1] 25.06628
```

(Also see F&L listing 11.1)

# Importance sampling

Suppose we have:

- $p(\boldsymbol{\theta}|\mathcal{M})$  (prior over parameters)
- $p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M})$  (likelihood)

and we want the marginal likelihood:

- $$p(\mathbf{y}|\mathcal{M}) = \int_{\boldsymbol{\theta}} p(\mathbf{y}|\boldsymbol{\theta}, \mathcal{M})p(\boldsymbol{\theta}|\mathcal{M})d\boldsymbol{\theta}$$

as well as expected values for some psychologically interpretable parameters:

- $E[\boldsymbol{\theta}|\mathbf{y}, \mathcal{M}]$

# Importance sampling

We can estimate these quantities using *importance sampling*:  
(omitting  $\mathcal{M}$ )

- 1 Draw  $J$  samples from a **normalized** *proposal distribution*  
 $\theta^{(j)} \sim g(\theta)$
- 2 Weight each sample:  $w_j = \frac{p(\mathbf{y}|\theta^{(j)})p(\theta^{(j)})}{g(\theta^{(j)})}$
- 3 The expectation of  $\theta$  is approximately  $\frac{\sum_j w_j \theta^{(j)}}{\sum_j w_j}$
- 4 The marginal likelihood is approximately  $\frac{1}{J} \sum_j w_j$

If the variance of the weights is low, these are *probably* trustworthy estimates.

# Simple Monte Carlo integration

A special case of importance sampling:

- 1 Sample from the prior (usually easy)
- 2 Weight by likelihood:  $w_j = \frac{p(\mathbf{y}|\boldsymbol{\theta}^{(j)})p(\boldsymbol{\theta}^{(j)})}{p(\boldsymbol{\theta}^{(j)})} = p(\mathbf{y}|\boldsymbol{\theta}^{(j)})$

Easy and acceptable if you think the samples will cover high-density areas of the posterior.

Better: Find a (normalized) proposal function that resembles your posterior.

# Importance sampling

```
impSamp <- function(targD,ef) {  
  nSamps = 40000 # The more the better  
  # Using a student's t distribution, df=1  
  proposals <- rt(nSamps,1)  
  pDens <- dt(proposals,1)  
  unnP <- targD(proposals)  
  w <- unnP/pDens  
  print(paste("Expected value of target function:",  
    sprintf("%.2.3f",sum(w*ef(proposals))/sum(w))))  
  print(paste("Average importance weight:",  
    sprintf("%.2.3f",sum(w)/nSamps)))  
}
```

## Gaussian example

```
sd=.05;mu=5;
unnGauss <- function(x) {exp(-(mu-x)^2/(2*sd^2))}
# Real
print(sprintf("Real: %.3f",sqrt(2*pi*sd^2)))
# Numerical
adaptIntegrate(unnGauss,c(-1E3),c(1E3))
# Importance
impSamp(unnGauss,function(x) x)
```



## Gaussian example

```
sd=.05;mu=5
```

```
Real normalization constant Z: 0.125
```

```
Cubature estimate of Z: 0 (oops)
```

```
Importance sampling:
```

```
  Estimated mean: 5.000
```

```
  Avg. importance weight (Z): 0.132 (close)
```

## Importance sampling

In the Gaussian example, we estimated the integral of the unnormalized Gaussian, by integrating over  $x$ .

$$\int e^{-\frac{(\mu-x)^2}{2\sigma^2}} dx$$

If we were interested in the marginal likelihood, we would propose  $\mu$  and  $\sigma$  (i.e., the parameters  $\theta = [\mu, \sigma]$ ) rather than  $x$ . For this we would need proposal distributions and priors for both parameters.

$$p(x) = \int \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\mu-x)^2}{2\sigma^2}} p(\mu, \sigma) d\mu d\sigma$$

# Importance sampling

- General-purpose Monte Carlo method for approximating parameter distributions
- Can exploit knowledge about high-density regions of posterior
- Can compute expectations of functions of params
- Requires good proposals
  - Increasingly so as dimensionality goes up
- In these cases, additional tricks may be necessary, e.g.,
  - **annealed** importance sampling ([link](#))
  - Inference trees ([link](#))

# Harmonic mean estimation

See “The Harmonic Mean of the Likelihood: Worst Monte Carlo Method Ever” by Radford Neal. ([link](#))

Excerpts:

- “abysmal performance in most real problem[s]”
- “the total unsuitability of the harmonic mean estimator should have been apparent within an hour of its discovery”

Don't use it.

# Transdimensional MCMC

- Use Markov Chain Monte Carlo, combining multiple models into a single overarching one
- Nice in principle
- Often difficult / fiddly in practice
- Out of scope for this course

# Savage-Dickey density ratio

- Efficiently compare nested probabilistic models via Bayes Factor
- Effectively a better and Bayesian alternative to the likelihood-ratio test
  - See recommended reading to learn more

# BIC

$$BIC = -2 \cdot LL + K \log(N)$$

where  $N$  is the number of data points and  $K$  is the number of parameters.

- Motivated by model comparison per se, not prediction
- Approximates the log marginal likelihood under a set of assumptions
- Can be understood as a “minimum description length” approach
- Like AIC, a model-comparison method that boils down to MLE likelihoods and counting parameters
- Like AIC, rests on assumptions; guarantees are asymptotic
- Easy!
- Safer than AIC if arguing for a more complex model

# Identifiability

“Models are unidentifiable when there is no unique mapping between any possible data pattern and a corresponding set of parameter estimates”

*Strict* identifiability is a high bar. More pragmatically:

- We should be able to identify parameters *we care about*
- There may be some patterns of data that don't lead to unique parameter estimates



# Identifiability

If likelihood doesn't depend strongly on parameters:

- Effectively a simpler model
- Parameters may not be identifiable in practice

# Identifiability

Weak relationships between parameters and data are not all bad:

- Less likely to have unwanted flexibility
- Sometime we want parameters *not* to matter
  - “Nuisance parameters”
  - Ideally we integrate them out and forget about them

# Identifiability

However:

- If parameters have important psychological interpretations, they should be identifiable
- Might be a sign that the experiment isn't adequate

# Identifiability

- ① Sometimes non-identifiability is inherent in a model; function from parameters to data isn't invertible
  - E.g., using mean response times to infer parameters of a Weibull distribution
- ② Sometimes identification is impossible in practice, because data are too sparse or noisy

# Identifiability

If identifiability is important, we can perform an identifiability “sanity check” before collecting data.

This is sometimes possible to do mathematically, e.g., Jacobian rank (See F&L 10.6.1), but **simulation**-based approaches are often easier and more useful.

# Fake data simulation

- ① Define your model(s) and decide on how you will estimate parameters and compare models
- ② Choose some theoretically interesting and plausible hypothetical models and/or parameters
- ③ Simulate data for your experiment based on each of those hypotheticals
  - Inspect your simulated data. Do they look implausible? If so, revisit steps 1-2
- ④ Compare models/fit parameters given your simulated data
  - Do you come to the right conclusion?
  - If not, you need to fix your models, methods, and/or experiments

Can also serve as a power analysis.

# Summary

- Many methods for approximating the marginal likelihood
- Easy cases:
  - low-dimensional models (numerical integration, importance sampling)
  - nested models (Savage-Dickey)
  - conjugate priors (didn't discuss)
- Hard case: High-dimensional, non-conjugate, non-nested
  - Well-tuned importance sampling
  - Annealed importance sampling (out of scope)
  - Transdimensional MCMC (out of scope)
- If you care about parameter identifiability, check!