



# Model Comparison

**Matthew Talluto**

August 18, 2017



UNIVERSITÉ DE  
SHERBROOKE

*Douter de tout ou tout croire sont deux solutions également commodes,  
qui nous dispensent de réfléchir.*

–Henri Poincaré

# Introduction to Model Comparison

Why compare models?

# Introduction to Model Comparison

Why compare models?

- All models are imperfect
- How good is our model *given the modelling goals?*

## Comparing models

Before beginning, evaluate the goals of the comparison

- Predictive performance
- Hypothesis testing
- Reduction of overfitting

If you are asking yourself, “should I use A/B/DIC?”

Remember Betteridge’s law. . .

## Comparing models

Before beginning, evaluate the goals of the comparison

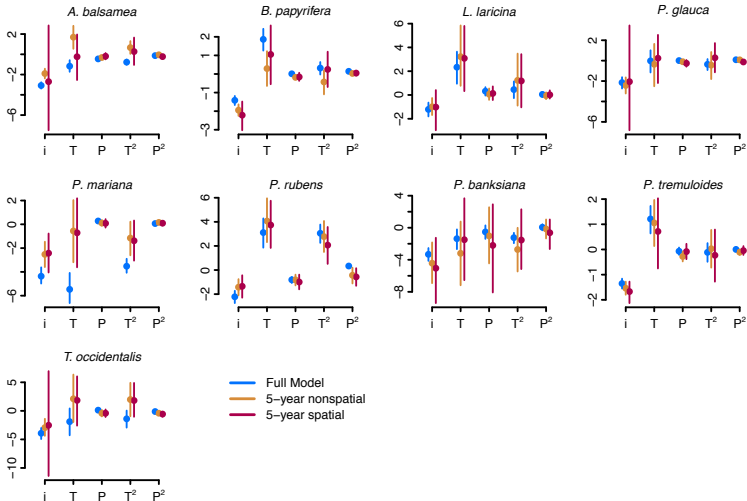
- Predictive performance
- Hypothesis testing
- Reduction of overfitting

If you are asking yourself, “should I use A/B/DIC?”

Remember Betteridge’s law. . .

*Any headline that ends in a question mark can be answered with the word “NO”*

# Informal model comparison



# Comparison through evaluation

If the goal is predictive performance, evaluate directly.

- Cross-validation
- k-fold cross validation

Cost: can be computationally intensive (especially for Bayesian). But you are already paying this cost (you ARE evaluating your models, right?)



# Comparison through evaluation

If the goal is predictive performance, evaluate directly.

- Cross-validation
- k-fold cross validation

Cost: can be computationally intensive (especially for Bayesian). But you are already paying this cost (you ARE evaluating your models, right?)

Requires selecting an evaluation score

- ROC/TSS (classification)
- RMSE (continuous)
- Goodness of fit
- ...

# Bayesian predictive performance

Consider a regression model

$$\begin{aligned}\text{pr}(\theta|y, x) &\propto \text{pr}(y, x, |\theta)\text{pr}(\theta) \\ y &\sim \mathcal{N}(\alpha + \beta x, \sigma)\end{aligned}$$

From a new value  $\hat{x}$  we can compute a posterior prediction  $\hat{y} = \alpha + \beta x$

## Bayesian predictive performance

We can then compute the *log posterior predictive density* (lppd):

$$\text{lppd} = \text{pr}(\hat{y}|\theta)$$

# Bayesian predictive performance

We can then compute the *log posterior predictive density* (lppd):

$$\text{lppd} = \text{pr}(\hat{y}|\theta)$$

Where is the prior?

# Bayesian predictive performance

We want to summarize lppd taking into account:

- an entire set of prediction points  $\hat{x} = \{x_1, x_2, \dots, x_n\}$
- the entire posterior distribution of  $\theta$ 
  - (or, realistically, a set of  $S$  draws from the posterior distribution)

## Bayesian predictive performance

We want to summarize lppd taking into account:

- an entire set of prediction points  $\hat{x} = \{x_1, x_2, \dots, x_n\}$
- the entire posterior distribution of  $\theta$ 
  - (or, realistically, a set of  $S$  draws from the posterior distribution)

$$\text{lppd} = \sum_{i=1}^n \log \left( \frac{1}{S} \sum_{s=1}^S \text{pr}(\hat{y} | \theta^s) \right)$$

## Bayesian predictive performance

We want to summarize lppd taking into account:

- an entire set of prediction points  $\hat{x} = \{x_1, x_2, \dots, x_n\}$
- the entire posterior distribution of  $\theta$ 
  - (or, realistically, a set of  $S$  draws from the posterior distribution)

$$\text{lppd} = \sum_{i=1}^n \log \left( \frac{1}{S} \sum_{s=1}^S \text{pr}(\hat{y} | \theta^s) \right)$$

To compare two competing models  $\theta_1$  and  $\theta_2$ , simply compute  $\text{lppd}_{\theta_1}$  and  $\text{lppd}_{\theta_2}$ , the “better” model (for prediction) is the one with a larger lppd.

## Information criteria

What do we do when  $\theta_1$  and  $\theta_2$  are very different?



## Information criteria

What do we do when  $\theta_1$  and  $\theta_2$  are very different?

Considering the lpd (using the calibration data), it can be proven, when  $\theta_2$  is *strictly nested* within  $\theta_1$ , that  $\text{lpd}_{\theta_1} > \text{lpd}_{\theta_2}$ .

## Information criteria

What do we do when  $\theta_1$  and  $\theta_2$  are very different?

Considering the lpd (using the calibration data), it can be proven, when  $\theta_2$  is *strictly nested* within  $\theta_1$ , that  $\text{lpd}_{\theta_1} > \text{lpd}_{\theta_2}$ .

Thus, we require a method for penalizing the larger (or more generally, more flexible) model to avoid simply overfitting, especially when validation data are unavailable.

$$\text{AIC} = 2k - 2 \log \text{pr}(x|\hat{\theta})$$

- $\text{pr}(x|\hat{\theta}) = \max(\text{pr}(x|\theta))$  and  $k$  is the number of parameters.
- AIC increases as the model gets worse or the number of parameters gets larger
- $-2 \log \text{pr}(x|\hat{\theta})$  is sometimes referred to as *deviance*

$$\text{AIC} = 2k - 2 \log \text{pr}(x|\hat{\theta})$$

- $\text{pr}(x|\hat{\theta}) = \max(\text{pr}(x|\theta))$  and  $k$  is the number of parameters.
- AIC increases as the model gets worse or the number of parameters gets larger
- $-2 \log \text{pr}(x|\hat{\theta})$  is sometimes referred to as *deviance*

What is the number of parameters in a hierarchical model?

$$D(\theta) = -2 \log(\text{pr}(x|\theta))$$

$$D(\theta) = -2 \log(\text{pr}(x|\theta))$$

We still penalize the model based on complexity, but we must estimate how many *effective* parameters there are:

$$p_D = \mathbb{E}[D(\theta)] - D(\mathbb{E}[\theta])$$

$$D(\theta) = -2 \log(\text{pr}(x|\theta))$$

We still penalize the model based on complexity, but we must estimate how many *effective* parameters there are:

$$p_D = \mathbb{E}[D(\theta)] - D(\mathbb{E}[\theta])$$

$$\text{DIC} = D(\mathbb{E}[\theta]) + 2p_D$$

## Pros:

- Easy to estimate
- Widely used and understood
- Effective for a variety of models regardless of nestedness or model size

## Cons

- Not Bayesian
- Assume  $\theta \sim \mathcal{MN}$
- Modest computational cost



## Bayes factor

Consider two competing models  $\theta_1$  and  $\theta_2$

In classical likelihood statistics, we can compute the likelihood ratio:

$$LR = \frac{MLE(X|\theta_1)}{MLE(X|\theta_2)}$$

Consider two competing models  $\theta_1$  and  $\theta_2$

In classical likelihood statistics, we can compute the likelihood ratio:

$$LR = \frac{\text{MLE}(X|\theta_1)}{\text{MLE}(X|\theta_2)}$$

A fully Bayesian approach is to take into account the entire posterior distribution of both models:

$$K = \frac{\text{pr}(\theta_1|X)}{\text{pr}(\theta_2|X)}$$

For a single posterior estimate of each model:

$$\begin{aligned}K &= \frac{\text{pr}(\theta_1|X)}{\text{pr}(\theta_2|X)} \\ &= \frac{\text{pr}(X|\theta_1)\text{pr}(\theta_1)}{\text{pr}(X|\theta_2)\text{pr}(\theta_2)}\end{aligned}$$

## Bayes factor

To account for the entire distribution:

$$\begin{aligned} K &= \frac{\int \text{pr}(\theta_1|X)d\theta_1}{\int \text{pr}(\theta_2|X)d\theta_2} \\ &= \frac{\int \text{pr}(X|\theta_1)\text{pr}(\theta_1)d\theta_1}{\int \text{pr}(X|\theta_2)\text{pr}(\theta_2)d\theta_2} \end{aligned}$$

## And others

- Bayesian model averaging
- Reversible jump MCMC

```
library(mcmc)
suppressMessages(library(bayesplot))

logposterior <- function(params, dat)
{
  if(params[2] <= 0)
    return(-Inf)

  mu <- params[1]
  sig <- params[2]

  lp <- sum(dnorm(dat, mu, sig, log=TRUE)) +
    dnorm(mu, 16, 0.4, log = TRUE) +
    dgamma(sig, 1, 0.1, log = TRUE)
  return(lp)
}
```

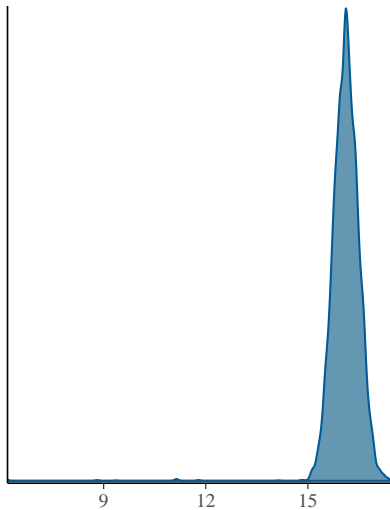
```
X <- c(15, 19.59, 15.06, 15.71, 14.65, 21.4, 17.64, 18.31,  
      15.12, 14.40)  
inits <- c(5, 2)  
tuning <- c(1.5, 0.5)  
  
model <- metrop(logposterior, initial = inits,  
               nbatch = 10000, dat = X, scale = tuning)  
model$accept
```

```
## [1] 0.2326
```

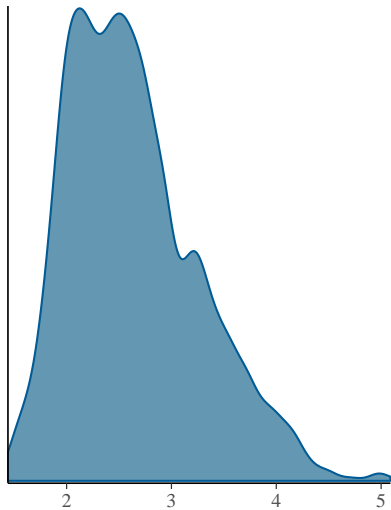
```
colnames(model$batch) = c('mu', 'sigma')  
colMeans(model$batch)
```

```
##          mu          sigma  
## 16.114213  2.635871
```

mu

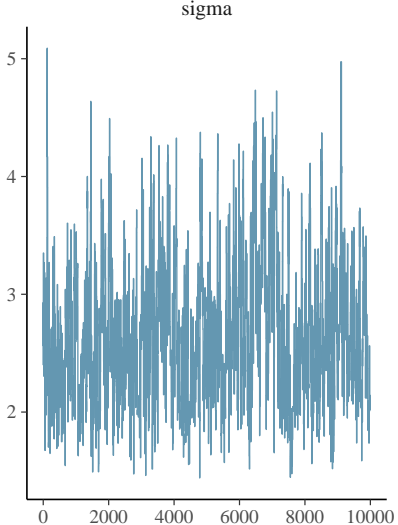
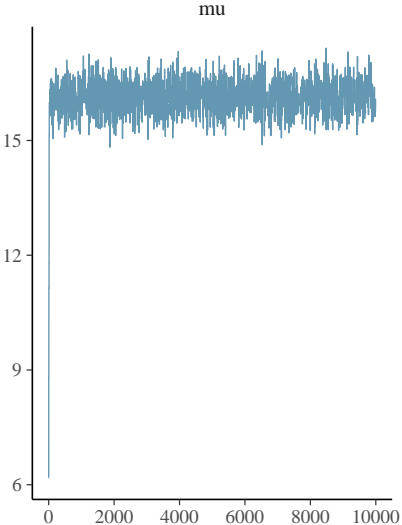


sigma





# Software



## Other software

- mcmc
- LaplacesDemon
- JAGS
- Stan