

Part 1: Recap on the Bayesian paradigm

Jean-Michel Marin

Université de Montpellier
Institut Montpelliérain Alexander Grothendieck (IMAG)

September, 22-27, Fréjus, France

- 1 The Bayesian paradigm
- 2 Bayesian estimates
- 3 Conjugate prior
- 4 Noninformative prior
- 5 Jeffreys prior
- 6 Bayesian Credible Intervals
- 7 Bayesian model choice
- 8 Bayesian Model Averaging
- 9 Difficulties with the Bayesian paradigm

The Bayesian paradigm

Bayes theorem = Inversion of probabilities

If A and B are events such that $\mathbb{P}(B) \neq 0$,

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)} =$$
$$\frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(A)\mathbb{P}(B|A) + \mathbb{P}(\bar{A})\mathbb{P}(B|\bar{A})}$$

The Bayesian paradigm

Subjectivism

Frank Plumpton Ramsey (1903-1930)

Bruno de Finetti (1906-1985)

Leonard Jimmie Savage (1921-1971)

The Bayesian paradigm

Given an iid sample $\mathcal{D}_n = (x_1, \dots, x_n)$ from a density $f(x|\theta)$, depending upon an unknown parameter $\theta \in \Theta$, the associated likelihood function is

$$\ell(\theta|\mathcal{D}_n) = \prod_{i=1}^n f(x_i|\theta)$$

The Bayesian paradigm

When \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$, we get

$$\begin{aligned}\ell(\theta|\mathcal{D}_n) &= \prod_{i=1}^n \exp\{-(x_i - \mu)^2/2\sigma^2\} / \sqrt{2\pi}\sigma \\ &\propto \exp\left\{-\sum_{i=1}^n (x_i - \mu)^2/2\sigma^2\right\} / \sigma^n \\ &\propto \exp\left\{-\left(n\mu^2 - 2n\bar{x}\mu + \sum_{i=1}^n x_i^2\right)/2\sigma^2\right\} / \sigma^n \\ &\propto \exp\left\{-[n(\mu - \bar{x})^2 + s^2]/2\sigma^2\right\} / \sigma^n,\end{aligned}$$

\bar{x} denotes the empirical mean and $s^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

The Bayesian paradigm

In the Bayesian approach θ is considered as a random variable

The Bayesian paradigm

In the Bayesian approach θ is considered as a random variable

In some sense, the likelihood function is transformed into a *posterior* distribution, which is a valid probability distribution on Θ

$$\pi(\theta|\mathcal{D}_n) = \frac{\ell(\theta|\mathcal{D}_n)\pi(\theta)}{\int \ell(\theta|\mathcal{D}_n)\pi(\theta) d\theta}$$

The Bayesian paradigm

In the Bayesian approach θ is considered as a random variable

In some sense, the likelihood function is transformed into a *posterior* distribution, which is a valid probability distribution on Θ

$$\pi(\theta|\mathcal{D}_n) = \frac{\ell(\theta|\mathcal{D}_n)\pi(\theta)}{\int \ell(\theta|\mathcal{D}_n)\pi(\theta) d\theta}$$

$\pi(\theta)$ is called the *prior* distribution and it has to be chosen to start the analysis

The Bayesian paradigm

The posterior density is a probability density on the parameter, which does not mean the parameter θ need be a genuine random variable

The Bayesian paradigm

The posterior density is a probability density on the parameter, which does not mean the parameter θ need be a genuine random variable

This density is used as an inferential tool, not as a truthful representation

The Bayesian paradigm

Two motivations:

- ▶ the prior distribution summarizes the *prior information* on θ . However, the choice of $\pi(\theta)$ is often decided on practical grounds rather than strong subjective beliefs
- ▶ the Bayesian approach provides a fully probabilistic framework for the inferential analysis, with respect to a reference measure $\pi(\theta)$

The Bayesian paradigm

Suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n

When σ^2 is known, if $\mu \sim \mathcal{N}(0, \sigma^2)$, then

$$\begin{aligned}\pi(\mu|\mathcal{D}_n) &\propto \pi(\mu) \ell(\theta|\mathcal{D}_n) \\ &\propto \exp\{-\mu^2/2\sigma^2\} \exp\{-n(\bar{x} - \mu)^2/2\sigma^2\} \\ &\propto \exp\{-(n+1)\mu^2/2\sigma^2 + 2n\mu\bar{x}/2\sigma^2\} \\ &\propto \exp\{-(n+1)[\mu - n\bar{x}/(n+1)]^2/2\sigma^2\}\end{aligned}$$

The Bayesian paradigm

Suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n

When σ^2 is known, if $\mu \sim \mathcal{N}(0, \sigma^2)$, then

$$\begin{aligned}\pi(\mu|\mathcal{D}_n) &\propto \pi(\mu) \ell(\theta|\mathcal{D}_n) \\ &\propto \exp\{-\mu^2/2\sigma^2\} \exp\{-n(\bar{x} - \mu)^2/2\sigma^2\} \\ &\propto \exp\{-(n+1)\mu^2/2\sigma^2 + 2n\mu\bar{x}/2\sigma^2\} \\ &\propto \exp\{-(n+1)[\mu - n\bar{x}/(n+1)]^2/2\sigma^2\}\end{aligned}$$

$$\mu|\mathcal{D}_n \sim \mathcal{N}(n\bar{x}/(n+1), \sigma^2/(n+1))$$

The Bayesian paradigm

When σ^2 is unknown, $\theta = (\mu, \sigma^2)$, if $\mu|\sigma^2 \sim \mathcal{N}(0, \sigma^2)$ and $\sigma^2 \sim \mathcal{IG}(1, 1)$, then $\pi((\mu, \sigma^2)|\mathcal{D}_n) \propto \pi(\sigma^2) \times \pi(\mu|\sigma^2) \times f(\mathcal{D}_n|\mu, \sigma^2)$

$$\propto (\sigma^{-2})^{1/2+2} \exp\{-(\mu^2 + 2)/2\sigma^2\} \mathbf{1}_{\sigma^2 > 0}$$

$$(\sigma^{-2})^{n/2} \exp\{-(n(\mu - \bar{x})^2 + s^2) / 2\sigma^2\}$$

The Bayesian paradigm

When σ^2 is unknown, $\theta = (\mu, \sigma^2)$, if $\mu|\sigma^2 \sim \mathcal{N}(0, \sigma^2)$ and $\sigma^2 \sim \mathcal{IG}(1, 1)$, then $\pi((\mu, \sigma^2)|\mathcal{D}_n) \propto \pi(\sigma^2) \times \pi(\mu|\sigma^2) \times f(\mathcal{D}_n|\mu, \sigma^2)$

$$\propto (\sigma^{-2})^{1/2+2} \exp\{-(\mu^2 + 2)/2\sigma^2\} \mathbf{1}_{\sigma^2 > 0}$$

$$(\sigma^{-2})^{n/2} \exp\{-(n(\mu - \bar{x})^2 + s^2)/2\sigma^2\}$$

$$\mu|\mathcal{D}_n, \sigma^2 \sim \mathcal{N}\left(\frac{n\bar{x}}{n+1}, \frac{\sigma^2}{n+1}\right)$$

$$\sigma^2|\mathcal{D}_n \sim \mathcal{IG}\left(\left\{1 + \frac{n}{2}\right\}, \left\{1 + \frac{s^2}{2} + \frac{n\bar{x}}{2(n+1)}\right\}\right)$$

The Bayesian paradigm

Variability in σ^2 induces more variability in μ , the marginal posterior in μ being then a Student's t distribution

The Bayesian paradigm

Variability in σ^2 induces more variability in μ , the marginal posterior in μ being then a Student's t distribution

$$\mu | \mathcal{D}_n \sim \mathcal{T} \left(n + 2, \frac{n\bar{x}}{n + 1}, \frac{2 + s^2 + (n\bar{x})/(n + 1)}{(n + 1)(n + 2)} \right)$$

Bayesian estimates

For a given loss function $L(\theta, \hat{\theta}(\mathcal{D}_n))$, we deduce a Bayesian estimate by minimizing the posterior expected loss:

$$\mathbb{E}_{\theta|\mathcal{D}_n}^{\pi} (L(\theta, \hat{\theta}(\mathcal{D}_n)))$$

Bayesian estimates

For a given loss function $L(\theta, \hat{\theta}(\mathcal{D}_n))$, we deduce a Bayesian estimate by minimizing the posterior expected loss:

$$\mathbb{E}_{\theta|\mathcal{D}_n}^{\pi} (L(\theta, \hat{\theta}(\mathcal{D}_n)))$$

To minimize the posterior expected loss is equivalent to minimize the Bayes risk, the frequentist risk integrated over the prior distribution

Bayesian estimates

For instance, for the L_2 loss function, the corresponding Bayes optimum is the expected value of θ under the posterior distribution,

$$\hat{\theta}(\mathcal{D}_n) = \int \theta \pi(\theta | \mathcal{D}_n) d\theta = \frac{\int \theta \ell(\theta | \mathcal{D}_n) \pi(\theta) d\theta}{\int \ell(\theta | \mathcal{D}_n) \pi(\theta) d\theta}$$

Bayesian estimates

When no specific penalty criterion is available, the posterior expectation is often used as a default estimator, although alternatives are also available. For instance, the *maximum a posteriori estimator* (MAP) is defined as

$$\hat{\theta}(\mathcal{D}_n) \in \operatorname{argmax}_{\theta} \pi(\theta|\mathcal{D}_n)$$

Bayesian estimates

When no specific penalty criterion is available, the posterior expectation is often used as a default estimator, although alternatives are also available. For instance, the *maximum a posteriori estimator* (MAP) is defined as

$$\hat{\theta}(\mathcal{D}_n) \in \operatorname{argmax}_{\theta} \pi(\theta|\mathcal{D}_n)$$

Similarity of with the maximum likelihood estimator: the influence of the prior distribution $\pi(\theta)$ on the estimate progressively disappears as the number of observations n increases

Conjugate prior

The selection of the prior distribution is an important issue in Bayesian statistics

Conjugate prior

The selection of the prior distribution is an important issue in Bayesian statistics

When prior information is available about the data or the model, it can be used in building the prior

Conjugate prior

The selection of the prior distribution is an important issue in Bayesian statistics

When prior information is available about the data or the model, it can be used in building the prior

In many situations, however, the selection of the prior distribution is quite delicate

Conjugate prior

The selection of the prior distribution is an important issue in Bayesian statistics

When prior information is available about the data or the model, it can be used in building the prior

In many situations, however, the selection of the prior distribution is quite delicate

Since the choice of the prior distribution has a considerable influence on the resulting inference, this inferential step must be conducted with the utmost care

Conjugate prior

Conjugate priors are such that the prior and posterior densities belong to the same parametric family

Conjugate prior

Conjugate priors are such that the prior and posterior densities belong to the same parametric family

An advantage when using a conjugate prior, is that one has to select only a few parameters to determine the prior distribution

Conjugate prior

Conjugate priors are such that the prior and posterior densities belong to the same parametric family

An advantage when using a conjugate prior, is that one has to select only a few parameters to determine the prior distribution

But the information known a priori may be either insufficient or incompatible with the structure imposed by conjugacy

Conjugate prior

Justifications

- ▶ Device of virtual past observations
- ▶ First approximations to adequate priors, backed up by robustness analysis
- ▶ But mostly... tractability and simplicity

Conjugate prior

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Normal $\mathcal{N}(\theta, \sigma^2)$	Normal $\mathcal{N}(\mu, \tau^2)$	$\mathcal{N}(\rho(\sigma^2\mu + \tau^2x), \rho\sigma^2\tau^2)$ $\rho^{-1} = \sigma^2 + \tau^2$
Poisson $\mathcal{P}(\theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + x, \beta + 1)$
Gamma $\mathcal{G}(\nu, \theta)$	Gamma $\mathcal{G}(\alpha, \beta)$	$\mathcal{G}(\alpha + \nu, \beta + x)$
Binomial $\mathcal{B}(n, \theta)$	Beta $\mathcal{Be}(\alpha, \beta)$	$\mathcal{Be}(\alpha + x, \beta + n - x)$

Conjugate prior

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Negative Binomial $Neg(m, \theta)$	Beta $Be(\alpha, \beta)$	$Be(\alpha + m, \beta + x)$
Multinomial $\mathcal{M}_k(\theta_1, \dots, \theta_k)$	Dirichlet $\mathcal{D}(\alpha_1, \dots, \alpha_k)$	$\mathcal{D}(\alpha_1 + x_1, \dots, \alpha_k + x_k)$
Normal $\mathcal{N}(\mu, 1/\theta)$	Gamma $\mathcal{G}a(\alpha, \beta)$	$\mathcal{G}a(\alpha + 0.5, \beta + (\mu - x)^2/2)$

Noninformative prior

Conjugate priors are nice to work with, but require hyperparameters's determination

Noninformative prior

Conjugate priors are nice to work with, but require hyperparameters's determination

One can opt for a completely different perspective and rely on so-called *noninformative* priors that aim at attenuating the impact of the prior on the resulting inference

Noninformative prior

Conjugate priors are nice to work with, but require hyperparameters's determination

One can opt for a completely different perspective and rely on so-called *noninformative* priors that aim at attenuating the impact of the prior on the resulting inference

These priors are fundamentally defined as coherent extensions of the uniform distribution

Noninformative prior

For unbounded parameter spaces, the densities of noninformative priors actually may fail to integrate to a finite number and they are defined instead as positive measures

Noninformative prior

For unbounded parameter spaces, the densities of noninformative priors actually may fail to integrate to a finite number and they are defined instead as positive measures

Generalized Bayesian estimators with improper prior distributions

Noninformative prior

Location models $x|\theta \sim f(x - \theta)$ are usually associated with flat priors $\pi(\theta) \propto 1$

Noninformative prior

Location models $x|\theta \sim f(x - \theta)$ are usually associated with flat priors $\pi(\theta) \propto 1$

Scale models $x|\theta \sim \frac{1}{\theta} f\left(\frac{x}{\theta}\right)$ are usually associated with the log-transform of a flat prior, that is, $\pi(\theta) \propto 1/\theta \times \mathbf{1}_{\theta>0}$

Jeffreys prior

In a more general setting, the noninformative prior favored by most Bayesians is the so-called **Jeffreys prior** which is related to the Fisher information matrix

$$I_x^F(\theta) = -\mathbb{E} \left(\frac{\partial^2 \log f(x|\theta)}{(\partial\theta)^2} \right)$$

by

$$\pi^J(\theta) \propto \sqrt{|I_x^F(\theta)|} \times \mathbf{1}_{\theta \in \Theta},$$

where $|I|$ denotes the determinant of the matrix I

Jeffreys prior

Suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$

Jeffreys prior

Suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$

The Fisher information matrix leads to the Jeffreys prior

$$\pi^J(\mu, \sigma^2) \propto 1/\{(\sigma^2)\}^{3/2} \mathbf{1}_{\sigma^2 > 0}$$

Jeffreys prior

Suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$

The Fisher information matrix leads to the Jeffreys prior

$$\pi^J(\mu, \sigma^2) \propto 1/\{(\sigma^2)\}^{3/2} \mathbf{1}_{\sigma^2 > 0}$$

$$\mu | \sigma^2, \mathcal{D}_n \sim \mathcal{N}(\bar{x}, \sigma^2/n)$$

$$\sigma^2 | \mathcal{D}_n \sim \mathcal{IG}(n/2, s^2/2)$$

$$\mu | \mathcal{D}_n \sim \mathcal{T}(n, \bar{x}, s^2/n^2)$$

Bayesian Credible Intervals

Since the Bayesian approach processes θ as a random variable, a natural definition of a confidence region on θ is to determine $C(\mathcal{D}_n)$ such that

$$\pi(\theta \in C(\mathcal{D}_n) | \mathcal{D}_n) = 1 - \alpha$$

where α is a predetermined level

Bayesian Credible Intervals

Since the Bayesian approach processes θ as a random variable, a natural definition of a confidence region on θ is to determine $C(\mathcal{D}_n)$ such that

$$\pi(\theta \in C(\mathcal{D}_n) | \mathcal{D}_n) = 1 - \alpha$$

where α is a predetermined level

The integration is done over the parameter space, rather than over the observation space

Bayesian Credible Intervals

Since the Bayesian approach processes θ as a random variable, a natural definition of a confidence region on θ is to determine $C(\mathcal{D}_n)$ such that

$$\pi(\theta \in C(\mathcal{D}_n) | \mathcal{D}_n) = 1 - \alpha$$

where α is a predetermined level

The integration is done over the parameter space, rather than over the observation space

The quantity $1 - \alpha$ thus corresponds to the probability that a random θ belongs to this set $C(\mathcal{D}_n)$, rather than to the probability that the random set contains the true value of θ

Bayesian Credible Intervals

Given this drift in the interpretation of a confidence set is called a *credible set* by Bayesians.

Bayesian Credible Intervals

Given this drift in the interpretation of a confidence set is called a *credible set* by Bayesians.

A standard credible set corresponds to the values of θ with the highest posterior values,

$$C(\mathcal{D}_n) = \{\theta; \pi(\theta|\mathcal{D}_n) \geq k_\alpha\}$$

where k_α is determined by the coverage constraint

Bayesian Credible Intervals

Given this drift in the interpretation of a confidence set is called a *credible set* by Bayesians.

A standard credible set corresponds to the values of θ with the highest posterior values,

$$C(\mathcal{D}_n) = \{\theta; \pi(\theta|\mathcal{D}_n) \geq k_\alpha\}$$

where k_α is determined by the coverage constraint

This region is called the **Highest Posterior Density** (HPD) region

Bayesian Credible Intervals

Once again, suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$

$$\mu | \sigma^2, \mathcal{D}_n \sim \mathcal{N}(\bar{x}, \sigma^2/n)$$

$$\sigma^2 | \mathcal{D}_n \sim \mathcal{IG}(n/2, s^2/2)$$

$$\mu | \mathcal{D}_n \sim \mathcal{T}(n, \bar{x}, s^2/n^2)$$

Bayesian Credible Intervals

Once again, suppose \mathcal{D}_n is a normal $\mathcal{N}(\mu, \sigma^2)$ sample of size n and $\theta = (\mu, \sigma^2)$

$$\mu | \sigma^2, \mathcal{D}_n \sim \mathcal{N}(\bar{x}, \sigma^2/n)$$

$$\sigma^2 | \mathcal{D}_n \sim \mathcal{IG}(n/2, s^2/2)$$

$$\mu | \mathcal{D}_n \sim \mathcal{T}(n, \bar{x}, s^2/n^2)$$

Therefore, the credible interval of probability $1 - \alpha$ on μ is

$$[\bar{x} - t_{1-\alpha/2, n} \sqrt{(n-1)s^2/n^2}, \bar{x} + t_{1-\alpha/2, n} \sqrt{(n-1)s^2/n^2}]$$

Bayesian model choice

When are comparing models with indices $k = 1, 2, \dots, J$, we introduce a model indicator \mathfrak{M} taking values in $\{1, 2, \dots, J\}$ and representing the index of the “true” model

Bayesian model choice

When are comparing models with indices $k = 1, 2, \dots, J$, we introduce a model indicator \mathfrak{M} taking values in $\{1, 2, \dots, J\}$ and representing the index of the “true” model

If $\mathfrak{M} = k$, the data \mathcal{D}_n are generated from a statistical model \mathfrak{M}_k with likelihood $\ell_k(\theta_k|\mathcal{D}_n)$ and parameter $\theta_k \in \Theta_k$

Bayesian model choice

Bayes procedures will depend on the posterior probabilities in the model space

$$\mathbb{P}^\pi(\mathfrak{M} = k | \mathcal{D}_n)$$

Bayesian model choice

The prior π is defined over the collection of model indices $\{1, 2, \dots, J\}$, and, conditionally on the model index \mathfrak{M} , on the corresponding parameter space $\Theta_{\mathfrak{M}}$

Bayesian model choice

The prior π is defined over the collection of model indices $\{1, 2, \dots, J\}$, and, conditionally on the model index \mathfrak{M} , on the corresponding parameter space Θ_k

Choice of the prior model probabilities $\mathbb{P}^\pi(\mathfrak{M} = k)$

- ▶ in some cases, there is experimental or subjective evidence about those probabilities,
- ▶ typically, we are forced to settle for equal weights $\mathbb{P}^\pi(\mathfrak{M} = k) = 1/J$

A key quantity, the integrated likelihood, also called the evidence

$$\mathbb{P}^\pi(\mathfrak{M} = \mathbf{k} | \mathcal{D}_n) \propto \mathbb{P}^\pi(\mathfrak{M} = \mathbf{k}) \int \ell_{\mathbf{k}}(\theta_{\mathbf{k}} | \mathcal{D}_n) \pi_{\mathbf{k}}(\theta_{\mathbf{k}}) d\theta_{\mathbf{k}}$$

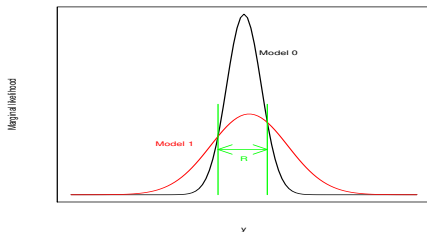
Bayesian model choice

A key quantity, the integrated likelihood, also called the evidence

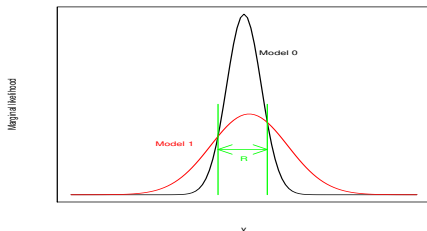
$$\mathbb{P}^\pi(\mathfrak{M} = k | \mathcal{D}_n) \propto \mathbb{P}^\pi(\mathfrak{M} = k) \int \ell_k(\theta_k | \mathcal{D}_n) \pi_k(\theta_k) d\theta_k$$

$\mathbb{P}^\pi(\mathfrak{M} = k | \mathcal{D}_n)$ is the core object in Bayesian model choice, the default procedure is to select the model with the highest posterior probability

Why Bayesian inference embodies Occam's razor?

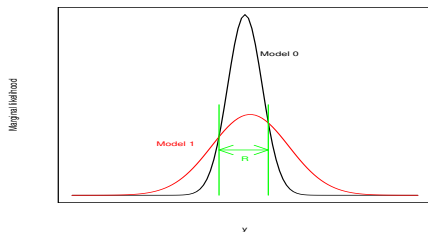


Why Bayesian inference embodies Occam's razor?



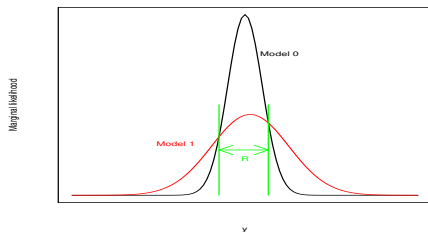
A simple model, like Model 0, makes only a limited range of predictions; a more powerful model, like Model 1, that has, for example, more free parameters, is able to predict a greater variety of data sets

Bayesian model choice



Suppose that equal prior probabilities have been assigned to the two models. Then, if the data set falls in region R, the less powerful model will be the more probable model

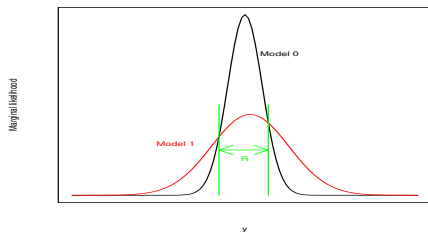
Bayesian model choice



Suppose that equal prior probabilities have been assigned to the two models. Then, if the data set falls in region R , the less powerful model will be the more probable model

The marginal likelihood corresponds to a penalized likelihood!

Bayesian model choice



Suppose that equal prior probabilities have been assigned to the two models. Then, if the data set falls in region R , the less powerful model will be the more probable model

The marginal likelihood corresponds to a penalized likelihood!

The BIC information criterium comes from an asymptotic Laplace approximation of the evidence

Bayesian model choice

Bayesian test and Bayesian model choice: the same problem

Bayesian model choice

Bayesian test and Bayesian model choice: the same problem

For instance, given a single observation $x \sim \mathcal{N}(\mu, \sigma^2)$ from a normal model where σ^2 is known

Bayesian model choice

Bayesian test and Bayesian model choice: the same problem

For instance, given a single observation $x \sim \mathcal{N}(\mu, \sigma^2)$ from a normal model where σ^2 is known

If $\mu \sim \mathcal{N}(\xi, \tau^2)$, the posterior distribution $\mu|x \sim \mathcal{N}(\xi(x), \omega^2)$ with

$$\xi(x) = \frac{\sigma^2 \xi + \tau^2 x}{\sigma^2 + \tau^2} \quad \text{and} \quad \omega^2 = \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2}$$

Bayesian model choice

If the question of interest is to decide whether μ is negative or positive, we can directly compute

$$\begin{aligned}\mathbb{P}^\pi(\mu < 0|x) &= \mathbb{P}^\pi\left(\frac{\mu - \xi(x)}{\omega} < \frac{-\xi(x)}{\omega}\right) \\ &= \Phi(-\xi(x)/\omega)\end{aligned}$$

where Φ is the normal cdf

Bayesian model choice

If the question of interest is to decide whether μ is negative or positive, we can directly compute

$$\begin{aligned}\mathbb{P}^\pi(\mu < 0|x) &= \mathbb{P}^\pi\left(\frac{\mu - \xi(x)}{\omega} < \frac{-\xi(x)}{\omega}\right) \\ &= \Phi(-\xi(x)/\omega)\end{aligned}$$

where Φ is the normal cdf

This computation does not seem to follow from the principles we just stated but it is only a matter of perspective

Bayesian model choice

We can derive the priors on both models from the original prior

Bayesian model choice

We can derive the priors on both models from the original prior

Deriving this posterior probability indeed means that, a priori, μ is negative with probability $\mathbb{P}^\pi(\mu < 0) = \Phi(-\xi/\tau)$ and that, in this model, the prior on μ is the truncated normal

$$\pi_1(\mu) = \frac{\exp\{-(\mu - \xi)^2/2\tau^2\}}{\sqrt{2\pi}\tau\Phi(-\xi/\tau)} \mathbb{I}_{\mu < 0}$$

while μ is positive with probability $\Phi(\xi/\tau)$ and, in this second model, the prior on μ is the truncated normal

$$\pi_2(\mu) = \frac{\exp\{-(\mu - \xi)^2/2\tau^2\}}{\sqrt{2\pi}\tau\Phi(\xi/\tau)} \mathbb{I}_{\mu > 0}$$

Bayesian model choice

The Bayes factor

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\mathbb{P}^{\pi}(\mathfrak{M} = 2 | \mathcal{D}_n) / \mathbb{P}^{\pi}(\mathfrak{M} = 1 | \mathcal{D}_n)}{\mathbb{P}^{\pi}(\mathfrak{M} = 2) / \mathbb{P}^{\pi}(\mathfrak{M} = 1)}$$

The Bayes factor

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\mathbb{P}^{\pi}(\mathfrak{M} = 2 | \mathcal{D}_n) / \mathbb{P}^{\pi}(\mathfrak{M} = 1 | \mathcal{D}_n)}{\mathbb{P}^{\pi}(\mathfrak{M} = 2) / \mathbb{P}^{\pi}(\mathfrak{M} = 1)}$$

While this quantity is a simple one-to-one transform of the posterior probability, it can be used for Bayesian model choice without first resorting to a determination of the prior weights of both models

The Bayes factor

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\mathbb{P}^{\pi}(\mathfrak{M} = 2 | \mathcal{D}_n) / \mathbb{P}^{\pi}(\mathfrak{M} = 1 | \mathcal{D}_n)}{\mathbb{P}^{\pi}(\mathfrak{M} = 2) / \mathbb{P}^{\pi}(\mathfrak{M} = 1)}$$

While this quantity is a simple one-to-one transform of the posterior probability, it can be used for Bayesian model choice without first resorting to a determination of the prior weights of both models

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\int_{\Theta_2} \ell_2(\theta_2 | \mathcal{D}_n) \pi_2(\theta_2) d\theta_2}{\int_{\Theta_1} \ell_1(\theta_1 | \mathcal{D}_n) \pi_1(\theta_1) d\theta_1} = \frac{m_2(\mathcal{D}_n)}{m_1(\mathcal{D}_n)}$$

The Ban on Improper Priors

Looking at the expression of the Bayes factor,

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\int_{\Theta_2} \ell_2(\theta_2|\mathcal{D}_n)\pi_2(\theta_2) d\theta_2}{\int_{\Theta_1} \ell_1(\theta_1|\mathcal{D}_n)\pi_1(\theta_1) d\theta_1}$$

it is clear that, when either π_1 or π_2 are improper, it is impossible to normalise the improper measures in a unique manner

The Ban on Improper Priors

Looking at the expression of the Bayes factor,

$$B_{21}^{\pi}(\mathcal{D}_n) = \frac{\int_{\Theta_2} \ell_2(\theta_2 | \mathcal{D}_n) \pi_2(\theta_2) d\theta_2}{\int_{\Theta_1} \ell_1(\theta_1 | \mathcal{D}_n) \pi_1(\theta_1) d\theta_1}$$

it is clear that, when either π_1 or π_2 are improper, it is impossible to normalise the improper measures in a unique manner

Therefore, the Bayes factor becomes completely arbitrary since it can be multiplied by one or two arbitrary constants

Bayesian model choice

Since improper priors are an essential part of the Bayesian approach, there are many proposals found in the literature to overcome this ban

Bayesian model choice

Since improper priors are an essential part of the Bayesian approach, there are many proposals found in the literature to overcome this ban

Most of those proposals rely on a device that transforms the improper prior into a proper probability distribution by exploiting a fraction of the data \mathcal{D}_n

Bayesian model choice

Since improper priors are an essential part of the Bayesian approach, there are many proposals found in the literature to overcome this ban

Most of those proposals rely on a device that transforms the improper prior into a proper probability distribution by exploiting a fraction of the data \mathcal{D}_n

The variety of available solutions is due to the many possibilities of removing the dependence on the choice of the portion of the data used in the first step

Bayesian model choice

Since improper priors are an essential part of the Bayesian approach, there are many proposals found in the literature to overcome this ban

Most of those proposals rely on a device that transforms the improper prior into a proper probability distribution by exploiting a fraction of the data \mathcal{D}_n

The variety of available solutions is due to the many possibilities of removing the dependence on the choice of the portion of the data used in the first step

The resulting procedures are called *pseudo-Bayes factors*, although some may actually correspond to true Bayes factors

Bayesian model choice

There is a major exception to this ban on improper priors that we can exploit

Bayesian model choice

There is a major exception to this ban on improper priors that we can exploit

If both models under comparison have parameters that have similar enough meanings to share the same prior distribution, as for instance a measurement error σ^2 , then the normalisation issue vanishes

Bayesian model choice

There is a major exception to this ban on improper priors that we can exploit

If both models under comparison have parameters that have similar enough meanings to share the same prior distribution, as for instance a measurement error σ^2 , then the normalisation issue vanishes

Note that we are not assuming that parameters are *common* to both models and thus that we do not contradict the earlier warning about different parameters to different models

Bayesian Model Averaging

The posterior probabilities in the model space can be used to average over the decisions coming from different models

Bayesian Model Averaging

The posterior probabilities in the model space can be used to average over the decisions coming from different models

Suppose that we are interested in the prediction of z and that, for model k , the predictive distribution of z is $g_k(z|\mathcal{D}_n)$

Bayesian Model Averaging

The posterior probabilities in the model space can be used to average over the decisions coming from different models

Suppose that we are interested in the prediction of z and that, for model k , the predictive distribution of z is $g_k(z|\mathcal{D}_n)$

The average predictive of z is

$$\sum_{k=1}^J \mathbb{P}^\pi(\mathfrak{M} = k|\mathcal{D}_n) g_k(z|\mathcal{D}_n)$$

Difficulties with the Bayesian paradigm

Prior difficulties:

- ▶ When we have prior informations, how to choose the prior distributions on the parameters of each model in a compatible way? What about the prior distribution in the models's space?
- ▶ When we do not have any prior information, **we can not use improper prior distribution**. Indeed, in that case, the models's posterior probabilities are only defined up to some arbitrary constants. How to choose the various prior distributions?

Difficulties with the Bayesian paradigm

Computational difficulties:

- ▶ How to approximate the various posterior probabilities?
- ▶ How to approximate the evidences?
- ▶ When the number of models in consideration is huge, how to explore the models's space?