Bayesian parameter inference

Jean-Michel Marin

University of Montpellier Faculty of Sciences

HAX918X / 2025-2026

- 1 The Bayesian paradigm
- Bayesian estimates
- Conjugate prior
- Moninformative prior
- Jeffreys prior
- Bayesian Credible Intervals

The interpretation of probability depends on different epistemological frameworks, and there are several philosophical approaches to understanding what probability means

long-run frequency of events / degree of belief in a hypothesis

Subjectivism

Harold Jeffreys (1891-1989)

Frank Plumpton Ramsey (1903-1930)

Bruno de Finetti (1906-1985)

Leonard Jimmie Savage (1921-1971)



Subjectivism in probability refers to the interpretation that probabilities represent personal beliefs or degrees of certainty about uncertain events

It contrasts with objective interpretations, viewing probability as a subjective measure based on individual judgment rather than frequency or inherent properties

Subjectivists argue that probabilities can vary between individuals based on their information and perspective

Bayesian probability, which updates beliefs as new information becomes available, is a key framework in this approach

R.A. Fisher (1890-1962) and Harold Jeffreys (1891-1989)

The Jeffreys-Fisher conflict is a most important episode in the recent history of scientific ideas. Both were themselves eminent mathematical scientists: Fisher a statistician and geneticist, Jeffreys a geophysicist

The controversy between R.A. Fisher and Harold Jeffreys was centered around differing interpretations of probability and statistical inference

Interpretation of Probability

Fisher (Frequentist Approach): Fisher viewed probability as a long-run frequency of events. According to the frequentist interpretation, probability is defined through repeated experiments: the probability of an event is the proportion of times it occurs in a large number of trials

Jeffreys (Bayesian Approach): Jeffreys, on the other hand, embraced a subjective interpretation of probability. He believed that probability measures the degree of belief in a hypothesis, which can be updated with new evidence

Hypothesis Testing

Fisher developed the method of significance testing, focusing on p-values as a measure of evidence against the null hypothesis. He believed that if the p-value was below a threshold (often 0.05), the null hypothesis could be rejected. He did not require an explicit alternative hypothesis, focusing more on the ability to reject or retain the null hypothesis

Jeffreys was critical of p-values and focused on Bayes factors, which compare the likelihood of the data under two competing hypotheses (the null and the alternative). Bayesian hypothesis testing incorporates prior information and gives a more direct comparison between the two hypotheses

Use of Prior Information

Fisher was critical of using prior probabilities, especially when they were subjective or arbitrary. He believed that statistical inference should be based on the data at hand without needing prior beliefs

Jeffreys, in contrast, argued that it was natural and necessary to incorporate prior information into statistical analysis. His approach, known as Bayesian inference, updates prior beliefs with data to produce a posterior probability, using Bayes' theorem

Likelihood Principle

Fisher strongly promoted the likelihood principle, arguing that the likelihood function (the probability of observing the data given different parameter values) was key in statistical inference. He introduced the idea of maximum likelihood estimation (MLE) to find the parameter values that maximize the likelihood of the observed data

Jeffreys agreed that the likelihood function is important but argued that it should be integrated with prior probabilities to produce posterior probabilities for parameters, a key feature of Bayesian statistics

Core of the Controversy

Fisher's frequentist approach was based on using data to make conclusions about hypotheses without reference to prior beliefs

Jeffreys' Bayesian approach involved incorporating prior knowledge and using Bayes' theorem to update the probability of a hypothesis given new data

Both approaches have strengths and limitations, and modern statistics has, in many cases, integrated elements of both Fisherian and Bayesian thinking

Fisher's frequentist methods have become standard in many scientific disciplines, especially in hypothesis testing and estimation techniques

Jeffreys' Bayesian ideas gained more prominence in the latter half of the 20th century, especially with the rise of computational power, allowing for more complex Bayesian analyses

Bayes theorem = Inversion of probabilities

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

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

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

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

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

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

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



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

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

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

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

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 $\pi(\theta)$

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

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

$$\begin{split} \pi(\mu|\mathscr{D}_n) &\propto \pi(\mu) \ f(\mathscr{D}_n|\mu) \\ &\propto \pi(\mu) \ \ell(\mu|\mathscr{D}_n) \\ &\propto \text{exp}\{-\mu^2/2\sigma^2\} \ \text{exp} \left\{-n(\bar{x}-\mu)^2/2\sigma^2\right\} \\ &\propto \text{exp} \left\{-(n+1)\mu^2/2\sigma^2 + 2n\mu\bar{x}/2\sigma^2\right\} \\ &\propto \text{exp} \left\{-(n+1)[\mu - n\bar{x}/(n+1)]^2/2\sigma^2\right\} \end{split}$$

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

When
$$\sigma^2$$
 is unknown, $\theta=(\mu,\sigma^2),$ if $\mu|\sigma^2\sim\mathcal{N}\left(0,\sigma^2\right)$ and $\sigma^2\sim\mathcal{I}(1,1),$ then $\pi((\mu,\sigma^2)|\mathcal{D}_n)\propto\pi(\sigma^2)\times\pi(\mu|\sigma^2)\times f(\mathcal{D}_n|\mu,\sigma^2)$
$$\propto (\sigma^{-2})^{1/2+2}\exp\left\{-(\mu^2+2)/2\sigma^2\right\}\mathbf{1}_{\sigma^2>0}$$

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

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

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

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

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

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:

$$\mathbb{E}_{\theta|\mathscr{D}_n}^{\pi}\left(\mathsf{L}\left(\theta,\hat{\theta}(\mathscr{D}_n)\right)\right)$$

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,

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

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

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

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

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 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

Justifications

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

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

$f(x \theta)$	$\pi(\theta)$	$\pi(\theta x)$
Negative Binomial	Beta	
$\mathcal{N}eg(m, heta)$	$\mathcal{B}e(lpha,eta)$	$\mathcal{B}e(\alpha+m,\beta+x)$
Multinomial	Dirichlet	
$\mathcal{M}_k(heta_1,\ldots, heta_k)$	$\mathcal{D}(\alpha_1,\ldots,\alpha_k)$	$\mathcal{D}(\alpha_1+x_1,\ldots,\alpha_k+x_k)$
Normal	Gamma	
$\mathcal{N}(\mu, 1/\theta)$	$\mathcal{G}a(lpha,eta)$	$\mathcal{G}(\alpha + 0.5, \beta + (\mu - x)^2/2)$

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

Generalized Bayesian estimators with improper prior distributions

Noninformative prior

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

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, $\pi(\theta) \propto 1/\theta \times 1_{\theta>0}$

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

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

where |I| denotes the determinant of the matrix I

Jeffreys prior

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

The Fisher information matrix leads to the Jeffreys prior

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

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

$$\sigma^2 | \mathscr{D}_n \sim \mathscr{I} \mathscr{G} \left(n/2, s^2/2 \right)$$

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

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.

A standard credible set corresponds to the values of $\boldsymbol{\theta}$ with the highest posterior values,

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

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

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

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

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

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

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