

Bayesian modelling

Léo Belzile

Table of contents

Welcome	1
1 Bayesics	3
1.1 Probability and frequency	3
1.2 Posterior distribution	4
1.3 Posterior predictive distribution	10
1.4 Summarizing posterior distributions	12
2 Priors	19
2.1 Conjugate priors	19
2.2 Uninformative priors	23
2.3 Expert knowledge	25
References	27

Welcome

This book is a web complement to MATH 80601A *Bayesian modelling*, a graduate course offered at HEC Montréal.

These notes are licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License and were last compiled on Saturday, September 02 2023.

The objective of the course is to provide a hands on introduction to Bayesian data analysis. The course will cover the formulation, evaluation and comparison of Bayesian models through examples and real-data applications.

1 Bayesics

The Bayesian paradigm is an inferential framework that is used widespread in data science. Numerical challenges that prevented it's widespread adoption until the 90's, when the Markov chain Monte Carlo revolution allowed models estimation.

Bayesian inference, which builds on likelihood-based inference, offers a natural framework for prediction and for uncertainty quantification. The interpretation is more natural than that of classical frequentists methods, and it is more easy to generalized models to complex settings, using notably hierarchical constructions. The main source of controversy is the role of the prior distribution, which allows one to incorporate subject-matter expertise but leads to different inferences being drawn by different practitioners; this subjectivity is not to the taste of many and has been the subject of many controversies.

The Bayesian paradigm includes multiples notions that are not covered in undergraduate introductory courses. The purpose of this chapter is to introduce these concepts and put them in perspective; the reader is assumed to be familiar with basics of likelihood-based inference. We begin with a discussion of the notion of probability, then define priors, posterior distributions, marginal likelihood and posterior predictive distributions. We focus on the interpretation of posterior distributions and explain how to summarize the posterior, leading leading to definitions of high posterior density region, credible intervals, posterior mode for cases where we either have a (correlated) sample from the posterior, or else have access to the whole distribution. Several notions, including sequentiality, prior elicitation and estimation of the marginal likelihood, are mentioned in passing. A brief discussion of Bayesian hypothesis testing (and alternatives) is presented.

1.1 Probability and frequency

In classical (frequentist) parametric statistic, we treat observations \mathbf{Y} as realizations of a distribution whose parameters θ are unknown. All of the information about parameters is encoded by the likelihood function.

The interpretation of probability in the classical statistic is in terms of long run frequency, which is why we term this approach frequentist statistic. Think of a fair die: when we state that values $\{1, \dots, 6\}$ are equiprobable, we mean that repeatedly tossing the die

1 Bayesics

should result, in large sample, in each outcome being realized roughly $1/6$ of the time (the symmetry of the object also implies that each facet should be equally likely to lie face up). This interpretation also carries over to confidence intervals: a $(1 - \alpha)$ confidence interval either contains the true parameter value or it doesn't, so the probability level $(1 - \alpha)$ is only the long-run proportion of intervals created by the procedure that should contain the true fixed value, not the probability that a single interval contains the true value. This is counter-intuitive to most.

In practice, the true value of the parameter θ vector is unknown to the practitioner, thus uncertain: Bayesians would argue that we should treat the latter as a random quantity rather than a fixed constant. Since different people may have different knowledge about these potential values, the prior knowledge is a form of **subjective probability**. For example, if you play cards, one person may have recorded the previous cards that were played, whereas other may not. They thus assign different probability of certain cards being played. In Bayesian inference, we consider θ as random variables to reflect our lack of knowledge about potential values taken. Italian scientist Bruno de Finetti, who is famous for the claim “Probability does not exist”, stated in the preface of Finetti (1974):

Probabilistic reasoning — always to be understood as subjective — merely stems from our being uncertain about something. It makes no difference whether the uncertainty relates to an unforeseeable future, or to an unnoticed past, or to a past doubtfully reported or forgotten: it may even relate to something more or less knowable (by means of a computation, a logical deduction, etc.) but for which we are not willing or able to make the effort; and so on [...] The only relevant thing is uncertainty — the extent of our knowledge and ignorance. The actual fact of whether or not the events considered are in some sense *determined*, or known by other people, and so on, is of no consequence.

On page 3, de Finetti continues (Finetti 1974)

only subjective probabilities exist — i.e., the degree of belief in the occurrence of an event attributed by a given person at a given instant and with a given set of information.

1.2 Posterior distribution

We consider a parametric model with parameters θ defined on $\Theta \subseteq \mathbb{R}^p$. In Bayesian learning, we adjoin to the likelihood $\mathcal{L}(\theta; \mathbf{y}) \equiv p(\mathbf{y} | \theta)$ a **prior** function $p(\theta)$ that reflects the prior knowledge about potential values taken by the p -dimensional parameter vector, before

observing the data \mathbf{y} . The prior makes θ random and the distribution of the parameter reflects our uncertainty about the true value of the model parameters.

In a Bayesian analysis, observations are random variables but inference is performed conditional on the observed sample values. By Bayes' theorem, our target is therefore the posterior density $p(\theta | \mathbf{y})$, defined as

$$\underbrace{p(\theta | \mathbf{y})}_{\text{posterior}} = \frac{\underbrace{p(\mathbf{y} | \theta)}_{\text{likelihood}} \times \underbrace{p(\theta)}_{\text{prior}}}{\underbrace{\int p(\mathbf{y} | \theta) p(\theta) d\theta}_{\text{marginal likelihood } p(\mathbf{y})}}. \quad (1.1)$$

The posterior $p(\theta | \mathbf{y})$ is proportional, as a function of θ , to the product of the likelihood and the prior function.

For the posterior to be **proper**, we need the product of the prior and the likelihood on the right hand side to be integrable as a function of θ over the parameter domain Θ . The integral in the denominator, termed marginal likelihood and denoted $p(\mathbf{y}) = E_{\theta}\{p(\mathbf{y} | \theta)\}$. The denominator of Equation 1.1 is a normalizing constant, making the posterior a valid density.

If θ is low dimensional, numerical integration such as quadrature methods can be used to compute the marginal likelihood.

To fix ideas, we consider next a simple one-parameter model where the marginal likelihood can be computed explicitly.

Example 1.1 (Binomial model with beta prior). Consider a binomial likelihood with probability of success p and n trials, $Y \sim \text{Bin}(n, p)$. If we take a beta prior, $p \sim \text{Be}(\alpha, \beta)$ and observe k successes, the posterior is

$$\begin{aligned} p(\theta | y = k) &\propto \binom{n}{k} p^k (1-p)^{n-k} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \\ &\propto p^{k+\alpha-1} (1-p)^{n-k+\beta-1} \end{aligned}$$

and is

$$\int_0^1 p^{k+\alpha-1} (1-p)^{n-k+\beta-1} dp = \frac{\Gamma(k + \alpha) \Gamma(n - k + \beta)}{\Gamma(n + \alpha + \beta)},$$

a Beta function. Since we need only to keep track of the terms that are function of the parameter p , we could recognize directly that the posterior distribution is $\text{Be}(k + \alpha, n - k + \beta)$ and deduce the normalizing constant from there.

1 Bayesics

If $Y \sim \text{Bin}(n, p)$, the expected number of success is np and the expected number of failures $n(1 - p)$ and so the likelihood contribution, relative to the prior, will dominate as the sample size n grows.

Another way to see this is to track moments (expectation, variance, etc.) The Beta distribution, whose density is $f(x; \alpha, \beta) \propto x^{\alpha-1}(1-x)^{\beta-1}$, has expectation $\alpha/(\alpha + \beta)$ and variance $\alpha\beta/\{(\alpha + \beta)^2(\alpha + \beta + 1)\}$. The posterior mean is

$$\mathbb{E}(p \mid y) = w \frac{y}{n} + (1 - w) \frac{a}{a + b}, \quad w = \frac{n}{n + a + b},$$

a weighted average of the maximum likelihood estimator and the prior mean. We can think of the parameter α (respectively β) as representing the prior number of success (resp. failures).

Figure 1.1 shows three different posterior distributions with different beta priors: the first prior, which favors values closer to $1/2$, leads to a more peaked posterior density, contrary to the second which is symmetric, but concentrated toward more extreme values near endpoints of the support. The rightmost panel is truncated: as such, the posterior is zero for any value of p beyond $1/2$ and so the posterior mode may be close to the endpoint of the prior. The influence of such a prior will not necessarily vanish as sample size and should be avoided, unless there are compelling reasons for restricting the domain.

Remark (Proportionality). Any term appearing in the likelihood times prior function that does not depend on parameters can be omitted since they will be absorbed by the normalizing constant. This makes it useful to compute normalizing constants or likelihood ratios.

Remark. An alternative parametrization for the beta distribution sets $\alpha = \mu\kappa$, $\beta = (1 - \mu)\kappa$ for $\mu \in (0, 1)$ and $\kappa > 0$, so that the model is parametrized directly in terms of mean μ , with κ capturing the dispersion.

Remark. A density integrates to 1 over the range of possible outcomes, but there is no guarantee that the likelihood function, as a function of θ , integrates to one over the parameter domain Θ .

For example, the binomial likelihood with n trials and k successes satisfies

$$\int_0^1 \binom{n}{k} p^k (1 - p)^{n-k} dp = \frac{1}{n + 1}.$$

Moreover, the binomial distribution is discrete (supported on $0, \dots, n$), whereas the likelihood is continuous as a function of the probability of success, as evidenced by Figure 1.2

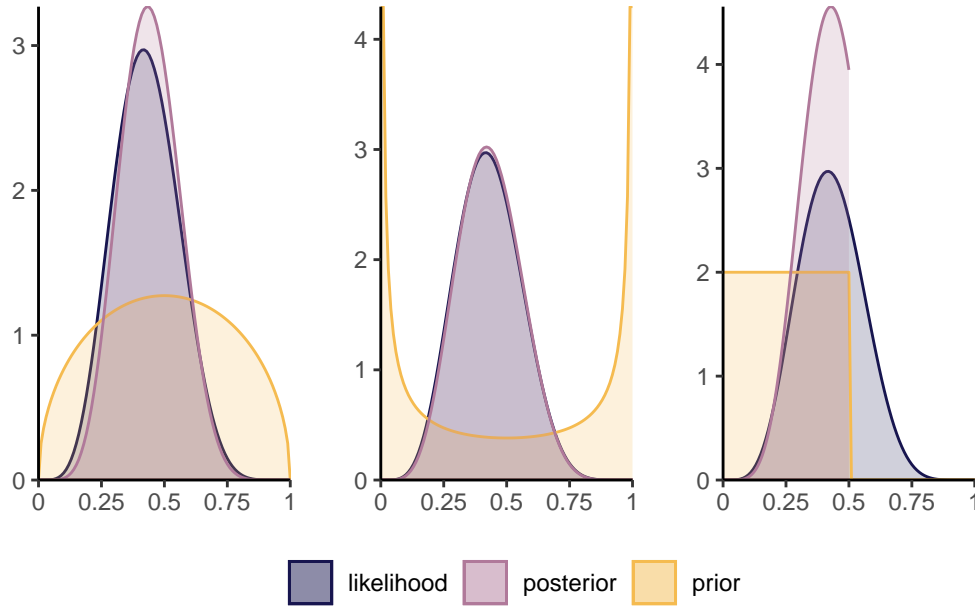


Figure 1.1: Scaled binomial likelihood for six successes out of 14 trials, with Beta(3/2, 3/2) prior (left), Beta(1/4, 1/4) (middle) and truncated uniform on $[0, 1/2]$ (right), with the corresponding posterior distributions.

Proposition 1.1 (Sequentiality and Bayesian updating). *The likelihood is invariant to the order of the observations if they are independent. Thus, if we consider two blocks of observations \mathbf{y}_1 and \mathbf{y}_2*

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1, \mathbf{y}_2) = p(\boldsymbol{\theta} \mid \mathbf{y}_1)p(\boldsymbol{\theta} \mid \mathbf{y}_2),$$

so it makes no difference if we treat data all at once or in blocks. More generally, for data exhibiting spatial or serial dependence, it makes sense to consider rather the conditional (sequential) decomposition

$$f(\mathbf{y}; \boldsymbol{\theta}) = f(\mathbf{y}_1; \boldsymbol{\theta})f(\mathbf{y}_2; \boldsymbol{\theta}, \mathbf{y}_1) \cdots f(\mathbf{y}_n; \boldsymbol{\theta}, \mathbf{y}_1, \dots, \mathbf{y}_{n-1})$$

where $f(\mathbf{y}_k; \mathbf{y}_1, \dots, \mathbf{y}_{k-1})$ denotes the conditional density function given observations $\mathbf{y}_1, \dots, \mathbf{y}_{k-1}$.

By Bayes' rule, we can consider updating the posterior by adding terms to the likelihood, noting that

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1, \mathbf{y}_2) \propto p(\mathbf{y}_2 \mid \mathbf{y}_1, \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathbf{y}_1)$$

1 Bayesics

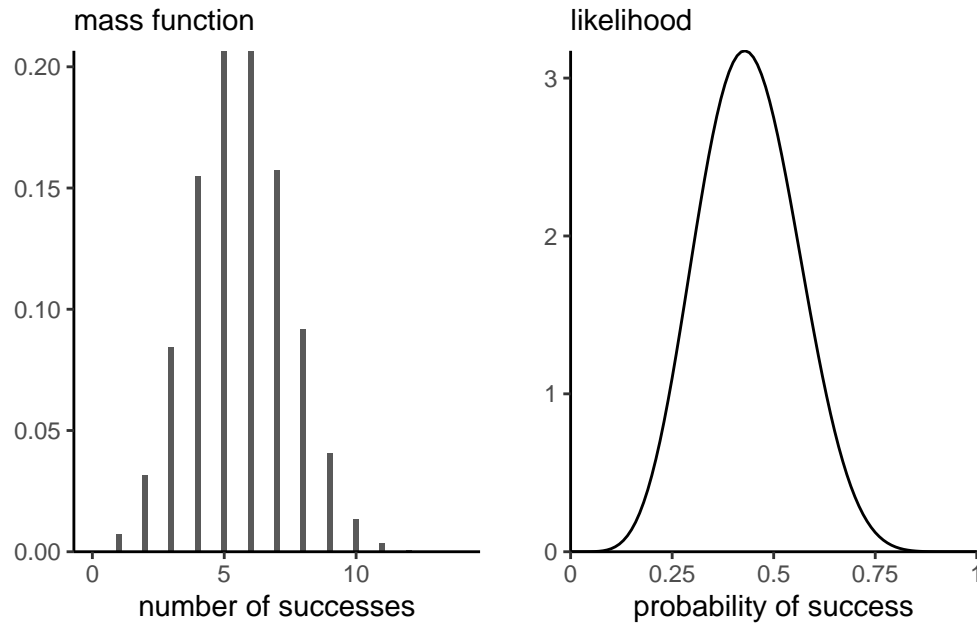


Figure 1.2: Binomial mass function (left) and scaled likelihood function (right).

which amounts to treating the posterior $p(\theta \mid y_1)$ as a prior. If data are exchangeable, the order in which observations are collected and the order of the belief updating is irrelevant to the full posterior. Figure 1.3 shows how the posterior becomes gradually closer to the scaled likelihood as we increase the sample size, and the posterior mode moves towards the true value of the parameter (here 0.3).

Example 1.2. While we can calculate analytically the value of the normalizing constant for the beta-binomial model, we could also for arbitrary priors use numerical integration or Monte Carlo methods in the event the parameter vector θ is low-dimensional.

While estimation of the normalizing constant is possible in simple models, the following highlights some challenges that are worth keeping in mind. In a model for discrete data (that is, assigning probability mass to a countable set of outcomes), the terms in the likelihood are probabilities and thus the likelihood becomes smaller as we gather more observations (since we multiply terms between zero and one). The marginal likelihood term becomes smaller and smaller, so its reciprocal is big and this can lead to arithmetic underflow.

1.2 Posterior distribution

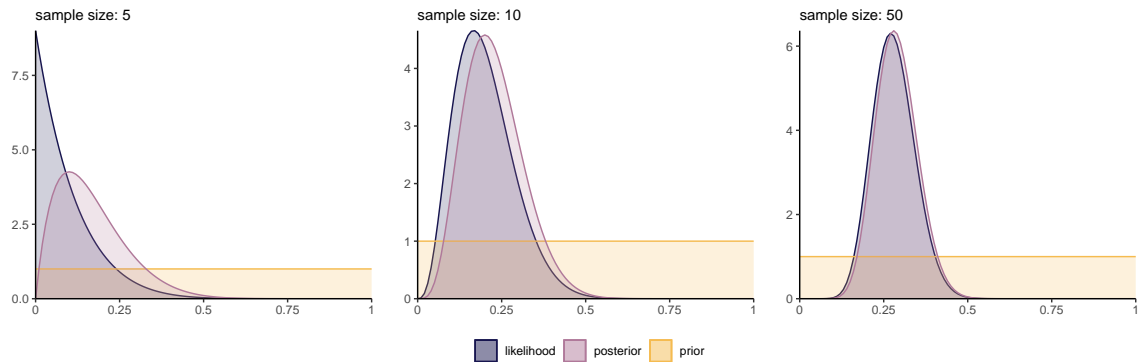


Figure 1.3: Beta posterior and binomial likelihood with a uniform prior for increasing number of observations (from left to right) out of a total of 100 trials.

```
k <- 6L # number of successes
n <- 14L # number of trials
alpha <- beta <- 1.5 # prior parameters
unnormalized_posterior <- function(p){
  p^(k+alpha-1) * (1-p)^(n-k + beta - 1)
}
integrate(f = unnormalized_posterior,
          lower = 0,
          upper = 1)
```

1.066906e-05 with absolute error < 1e-12

```
# Compare with known constant
beta(k + alpha, n - k + beta)
```

[1] 1.066906e-05

```
# Monte Carlo integration
mean(unnormalized_posterior(runif(1e5)))
```

[1] 1.064067e-05

1 Bayesics

When θ is high-dimensional, the marginal likelihood is intractable. This is one of the main challenges of Bayesian statistics and the popularity and applicability has grown drastically with the development and popularity of numerical algorithms, following the publication of Geman and Geman (1984) and Gelfand and Smith (1990). Markov chain Monte Carlo methods circumvent the calculation of the denominator by drawing approximate samples from the posterior.

1.3 Posterior predictive distribution

Prediction in the Bayesian paradigm is obtained by considering the *posterior predictive distribution*,

$$p(y_{\text{new}} | \mathbf{y}) = \int_{\Theta} p(y_{\text{new}} | \theta) p(\theta | \mathbf{y}) d\theta$$

Given draws from the posterior distribution, say θ_b ($b = 1, \dots, B$), we sample from each a new realization from the distribution appearing in the likelihood $p(y_{\text{new}} | \theta_b)$. This is different from the frequentist setting, which fixes the value of the parameter to some estimate $\hat{\theta}$; by contrast, the posterior predictive, here a beta-binomial distribution $\text{BetaBin}(n, \alpha + k, n - k + \beta)$, carries over the uncertainty so will typically be wider and overdispersed relative to the corresponding binomial model. This can be easily seen from the left-panel of Figure 1.4, which contrasts the binomial mass function evaluated at the maximum likelihood estimator $\hat{p} = 6/14$ with the posterior predictive.

```
npost <- 1e4L
# Sample draws from the posterior distribution
post_samp <- rbeta(n = npost, k + alpha, n - k + beta)
# For each draw, sample new observation
post_pred <- rbinom(n = npost, size = n, prob = post_samp)
```

Example 1.3. Consider an n sample of independent and identically distributed Gaussian, $Y_i \sim \text{No}(0, \tau^{-1})$ ($i = 1, \dots, n$), where we assign a gamma prior on the precision $\tau \sim \text{Ga}(\alpha, \beta)$. The posterior is

$$p(\tau | \mathbf{y}) \propto \prod_{i=1}^n \tau^{n/2} \exp\left(-\tau \frac{\sum_{i=1}^n y_i^2}{2}\right) \times \tau^{\alpha-1} \exp(-\beta\tau)$$

and rearranging the terms to collect powers of τ , etc. we find that the posterior for τ must also be gamma, with shape parameter $\alpha^* = \alpha + n/2$ and rate $\beta^* = \beta + \sum_{i=1}^n y_i^2/2$.

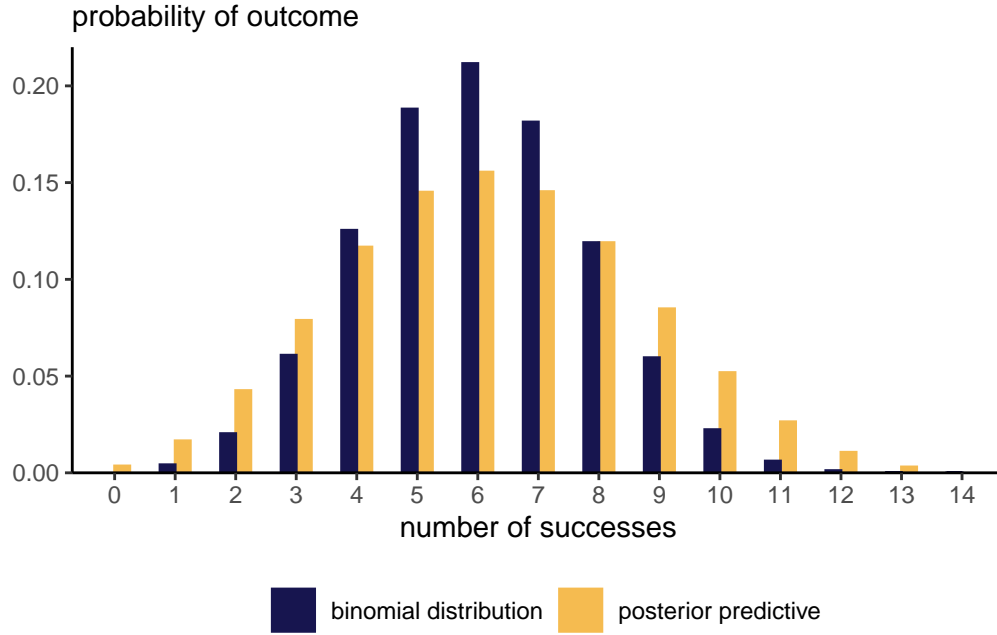


Figure 1.4: Beta-binomial posterior predictive distribution with corresponding binomial mass function evaluated at the maximum likelihood estimator.

The posterior predictive is

$$\begin{aligned}
p(y_{\text{new}} | \mathbf{y}) &= \int_0^\infty \frac{\tau^{1/2}}{(2\pi)^{1/2}} \exp(-\tau y_{\text{new}}^2/2) \frac{\beta^{*\alpha^*}}{\Gamma(\alpha^*)} \tau^{\alpha^*-1} \exp(-\beta^* \tau) d\tau \\
&= (2\pi)^{-1/2} \frac{\beta^{*\alpha^*}}{\Gamma(\alpha^*)} \int_0^\infty \tau^{\alpha^*-1/2} \exp\left\{-\tau(y_{\text{new}}^2/2 + \beta^*)\right\} d\tau \\
&= (2\pi)^{-1/2} \frac{\beta^{*\alpha^*}}{\Gamma(\alpha^*)} \frac{\Gamma(\alpha^* + 1/2)}{(y_{\text{new}}^2/2 + \beta^*)^{\alpha^*+1/2}} \\
&= \frac{\Gamma\left(\frac{2\alpha^*+1}{2}\right)}{\sqrt{2\pi}\Gamma\left(\frac{2\alpha^*}{2}\right) \beta^{*1/2}} \left(1 + \frac{y_{\text{new}}^2}{2\beta^*}\right)^{-\alpha^*-1/2} \\
&= \frac{\Gamma\left(\frac{2\alpha^*+1}{2}\right)}{\sqrt{\pi}\sqrt{2\alpha^*}\Gamma\left(\frac{2\alpha^*}{2}\right) (\beta^*/\alpha^*)^{1/2}} \left(1 + \frac{1}{2\alpha^*} \frac{y_{\text{new}}^2}{(\beta^*/\alpha^*)}\right)^{-\alpha^*-1/2}
\end{aligned}$$

which entails that Y_{new} is a scaled Student- t distribution with scale $(\beta^*/\alpha^*)^{1/2}$ and $2\alpha + n$ degrees of freedom. This example also exemplifies the additional variability relative to the distribution generating the data: indeed, the Student- t distribution is more heavy-tailed

than the Gaussian, but since the degrees of freedom increase linearly with n , the distribution converges to a Gaussian as $n \rightarrow \infty$, reflecting the added information as we collect more and more data points and the variance gets better estimated through $\sum_{i=1}^n y_i^2 / n$.

1.4 Summarizing posterior distributions

Most of the field of Bayesian statistics revolves around the creation of algorithms that either circumvent the calculation of the normalizing constant (notably using Monte Carlo and Markov chain Monte Carlo methods) or else provide accurate numerical approximation of the posterior pointwise, including for marginalizing out all but one parameters (integrated nested Laplace approximations, variational inference, etc.) The target of inference is the whole posterior distribution, a potentially high-dimensional object which may be difficult to summarize or visualize. We can thus report only characteristics of the the latter.

The choice of point summary to keep has it's root in decision theory.

Definition 1.1 (Loss function). A loss function $c(\theta, v)$ is a mapping from $\Theta \rightarrow \mathbb{R}^k$ that assigns a weight to each value of θ , corresponding to the regret or loss arising from choosing this value. The corresponding point estimator \hat{v} is the minimizer of the expected loss,

$$\hat{v} = \operatorname{argmin}_v \int_{\Theta} c(\theta, v) p(\theta | \mathbf{y}) d\theta$$

For example, in a univariate setting, the quadratic loss $c(\theta, v) = (\theta - v)^2$ returns the posterior mean, the absolute loss $c(\theta, v) = |\theta - v|$ returns the posterior median and the 0-1 loss $c(\theta, v) = \mathbb{I}(v \neq \theta)$ returns the posterior mode. All of these point estimators are central tendency measures, but some may be more adequate depending on the setting as they can correspond to potentially different values, as shown in the left-panel of Figure 1.5. The choice is application specific: for multimodal distributions, the mode is likely a better choice.

If we know how to evaluate the distribution numerically, we can optimize to find the mode or else return the value for the pointwise evaluation on a grid at which the density achieves it's maximum. The mean and median would have to be evaluated by numerical integration if there is no closed-form expression for the latter.

If we have rather a sample from the posterior with associated posterior density values, then we can obtain the mode as the parameter combination with the highest posterior, the median from the value at rank $\lfloor n/2 \rfloor$ and the mean through the sample mean of posterior draws.

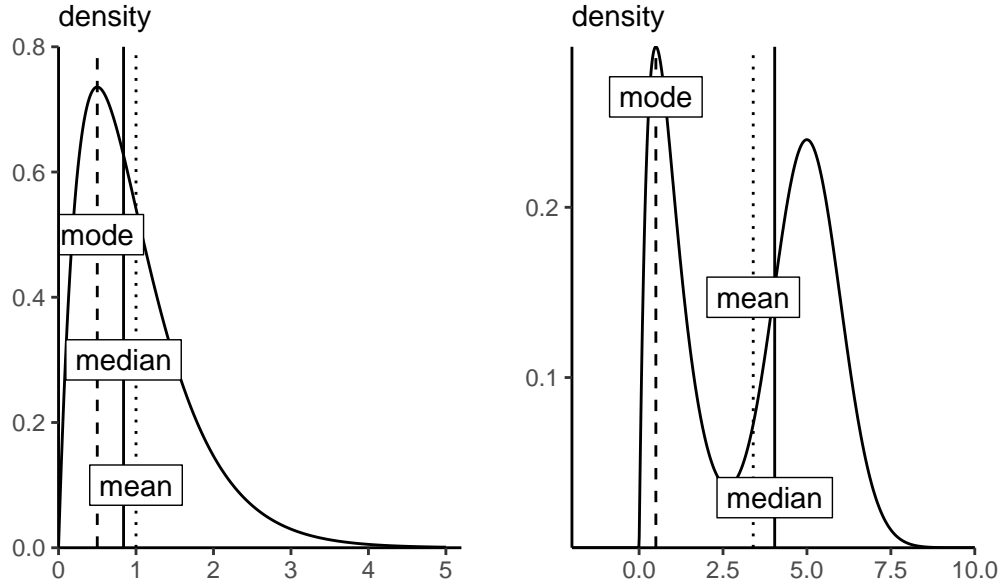


Figure 1.5: Point estimators from a left-skewed distribution (left) and from a multimodal distribution (right).

The loss function is often a functional (meaning a one-dimensional summary) from the posterior. The following example shows how it reduces a three-dimensional problem into a single risk measure.

Example 1.4. In extreme value, we are often interested in assessing the risk of events that are rare enough that they lie beyond the range of observed data. To provide a scientific extrapolation, it often is justified to fit a generalized Pareto distribution to exceedances of $Z = Y - u$, for some user-specified threshold u which is often taken as a large quantile of the distribution of Y . The generalized Pareto distribution function is

$$F(z; \tau, \xi) = 1 - \begin{cases} (1 + \xi/\tau z)_+^{-1/\xi}, & \xi \neq 0 \\ \exp(-z/\tau), & \xi = 0. \end{cases}$$

The shape ξ governs how heavy-tailed the distribution is, while τ is a scale parameter.

Insurance companies provide coverage in exchange for premiums, but need to safeguard themselves against very high claims by buying reinsurance products. These risks are often communicated through the value-at-risk (VaR), a high quantile exceeded with probability p . We model Danish fire insurance claim amounts for inflation-adjusted data collected

1 Bayesics

from January 1980 until December 1990 that are in excess of a million Danish kroner, found in the `evir` package and analyzed in Example 7.23 of McNeil, Frey, and Embrechts (2005). These claims are denoted Y and there are 2167 observations.

We fit a generalized Pareto distribution to exceedances above 10 millions kroner, keeping 109 observations or roughly the largest 5% of the original sample. Preliminary analysis shows that we can treat data as roughly independent and identically distributed and goodness-of-fit diagnostics (not shown) for the generalized Pareto suggest that the fit is adequate for all but the three largest observations, which are (somewhat severely) underestimated by the model.

The generalized Pareto model only describes the n_u exceedances above $u = 10$, so we need to incorporate in the likelihood a binomial contribution for the probability ζ_u of exceeding the threshold u . Provided that the priors for (τ, ξ) are independent of those for ζ_u , the posterior also factorizes as a product, so ζ_u and (τ, ξ) are a posteriori independent.

Suppose for now that we set a $\text{Be}(0.5, 0.5)$ prior for ζ_u and a non-informative prior for the generalized Pareto parameters. The `post_samp` matrix contains exact samples from the posterior distribution of (τ, ξ, ζ_u) , obtained using a Monte Carlo algorithm. Our aim is to evaluate the posterior distribution for the value-at-risk, the α quantile of Y for high values of α and see what point estimator one would obtain depending on our choice of loss function. For any $\alpha > 1 - \zeta_u$, the q_α is

$$\begin{aligned} 1 - \alpha &= \Pr(Y > q_\alpha \mid Y > u) \Pr(Y > u) \\ &= \left(1 + \xi \frac{q_\alpha - u}{\tau}\right)_+^{-1/\xi} \zeta_u \end{aligned}$$

and solving for q_α gives

$$q_\alpha = u + \frac{\tau}{\xi} \left\{ \left(\frac{\zeta_u}{1 - \alpha} \right)^\xi - 1 \right\}.$$

To obtain the posterior distribution of the α quantile, q_α , it suffices to plug in each posterior sample and evaluate the function: the uncertainty is carried over from the simulated values of the parameters to those of the quantile q_α . The left panel of Figure 1.6 shows the posterior density estimate of the $\text{VaR}(0.99)$ along with the maximum a posteriori (mode) of the latter.

Suppose that we prefer to under-estimate the value-at-risk rather than overestimate: this could be captured by the custom loss function

$$c(q, q_0) = \begin{cases} 0.5(0.99q - q_0), & q > q_0 \\ 0.75(q_0 - 1.01q), & q < q_0. \end{cases}$$

1.4 Summarizing posterior distributions

For a given value of the value-at-risk q_0 evaluated on a grid, we thus compute

$$r(q_0) = \int_{\Theta} c(q(\boldsymbol{\theta}), q_0) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}$$

and we seek to minimize the risk, $\hat{q} = \operatorname{argmin}_{q_0 \in \mathbb{R}_+} r(q_0)$. The value returned that minimizes the loss, shown in Figure 1.6, is to the left of the posterior mean for q_α .

```
# Compute value at risk from generalized Pareto distribution quantile fn
VaR_post <- with(post_samp, # data frame of posterior draws
  revdbayes::qgp( # with columns 'probexc', 'scale', 'shape'
    p = 0.01/probexc,
    loc = 10,
    scale = scale,
    shape = shape,
    lower.tail = FALSE))
# Loss function
loss <- function(qhat, q){
  mean(ifelse(q > qhat,
    0.5*(0.99*q-qhat),
    0.75*(qhat-1.01*q)))
}
# Create a grid of values over which to estimate the loss for VaR
nvals <- 101L
VaR_grid <- seq(
  from = quantile(VaR_post, 0.01),
  to = quantile(VaR_post, 0.99),
  length.out = nvals)
# Create a container to store results
risk <- numeric(length = nvals)
for(i in seq_len(nvals)){
  # Compute integral (Monte Carlo average over draws)
  risk[i] <- loss(q = VaR_post, qhat = VaR_grid[i])
}
```

The output of the Bayesian learning problem will be either of:

1. a fully characterized distribution
2. a numerical approximation to the posterior distribution (pointwise)
3. an exact or approximate sample drawn from the posterior distribution

In the first case, we will be able to directly evaluate quantities of interest if there are closed-form expressions for the latter, or else we could draw samples from the distribution and

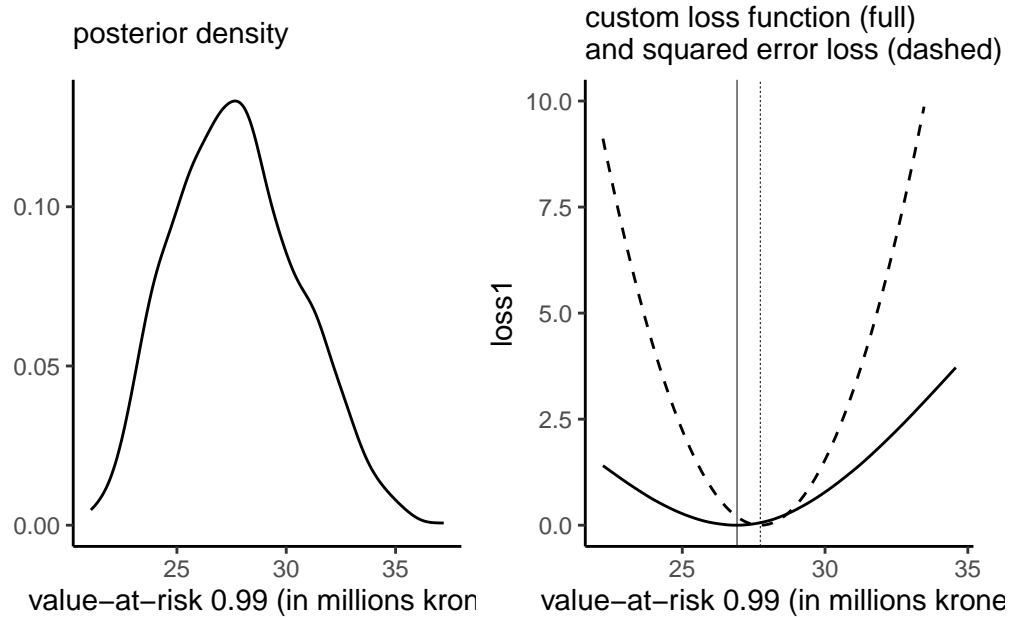


Figure 1.6: Posterior density (left) and losses functions for the 0.99 value-at-risk for the Danish fire insurance data. The vertical lines denote point estimates of the quantiles that minimize the loss functions.

evaluate them via Monte-Carlo. In case of numerical approximations, we will need to resort to numerical integration or otherwise to get our answers.

Often, we will also be interested in the marginal posterior distribution of each component θ_j in turn ($j = 1, \dots, p$). To get these, we carry out additional integration steps,

$$p(\theta_j | \mathbf{y}) = \int p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}_{-j}.$$

With a posterior sample, this is trivial: it suffices to keep the column corresponding to θ_j and discard the others.

To communicate uncertainty, we may resort to credible regions and intervals.

Definition 1.2. A $(1 - \alpha)$ **credible region** (or credible interval in the univariate setting) is a set \mathcal{S}_α such that, with probability level α ,

$$\Pr(\boldsymbol{\theta} \in \mathcal{S}_\alpha | \mathbf{Y} = \mathbf{y}) = 1 - \alpha$$

1.4 Summarizing posterior distributions

These intervals are not unique, as are confidence sets. In the univariate setting, the central or equi-tailed interval are the most popular, and easily obtained by considering the $\alpha/2, 1 - \alpha/2$ quantiles. These are easily obtained from samples by simply taking empirical quantiles. An alternative, highest posterior density credible sets, which may be a set of disjoint intervals obtained by considering the parts of the posterior with the highest density, may be more informative. The top panel Figure 1.7 shows the distinction for a bimodal mixture distribution, and a even more striking difference for 50% credible intervals for a symmetric beta distribution whose mass lie near the endpoints of the distribution, leading to no overlap between the two intervals.

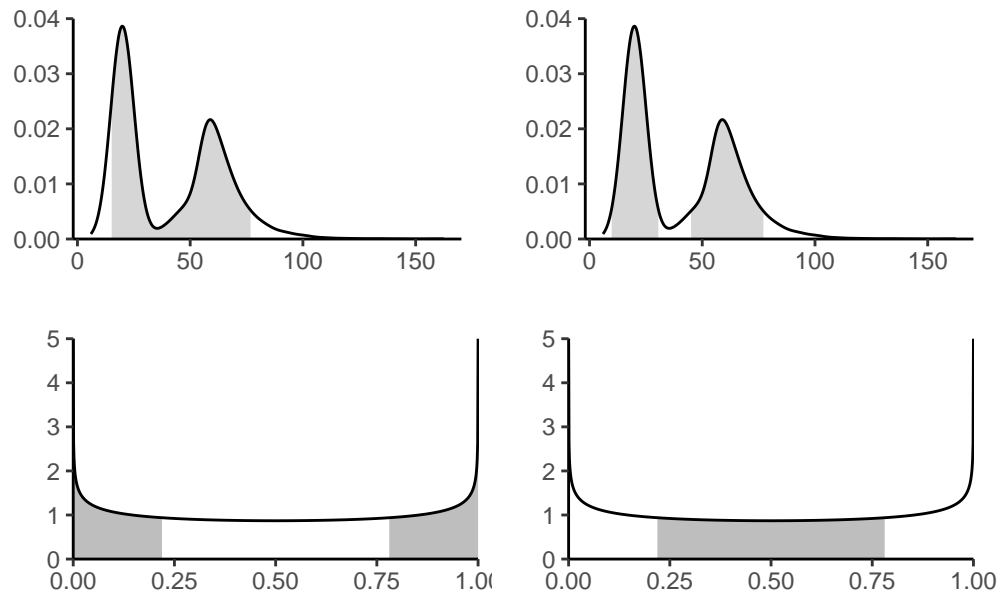


Figure 1.7: Density plots with 89% (top) and 50% (bottom) equitailed or central credible (left) and highest posterior density (right) intervals for two data sets, highlighted in grey.

Another example where these would differ

2 Priors

The posterior distribution combines two ingredients: the likelihood and the prior. If the former is a standard ingredient of any likelihood-based inference, prior specification requires some care. The purpose of this chapter is to consider different standard way of constructing prior functions.

2.1 Conjugate priors

A distribution belongs to an exponential family with parameter vector $\boldsymbol{\theta} \in \mathbb{R}^D$ if it can be written as

$$f(y; \boldsymbol{\theta}) = \exp \left\{ \sum_{k=1}^K Q_k(\boldsymbol{\theta}) t_k(y) + D(\boldsymbol{\theta}) \right\}$$

and in particular, the support does not depend on unknown parameters. If we have an independent and identically distributed sample of observations y_1, \dots, y_n , the log likelihood is thus of the form

$$\ell(\boldsymbol{\theta}) = \sum_{k=1}^K \phi_k(\boldsymbol{\theta}) \sum_{i=1}^n t_k(y_i) + nD(\boldsymbol{\theta}),$$

where the collection $\sum_{i=1}^n t_k(y_i)$ ($k = 1, \dots, K$) are sufficient statistics and $\phi_k(\boldsymbol{\theta})$ are the canonical parameters. The number of sufficient statistics are the same regardless of the sample size. Exponential families play a prominent role in generalized linear models, in which the natural parameters are modelled as linear function of explanatories.

A log prior density that is proportional to

$$\log p(\boldsymbol{\theta}) \propto \eta D(\boldsymbol{\theta}) + \sum_{k=1}^K Q_k(\boldsymbol{\theta}) \nu_k$$

is conjugate.

2 Priors

Example 2.1 (Conjugate priors for the binomial model). The binomial log density with y successes out of n trials is proportional to

$$y \log(p) + (n - y) \log(1 - p) = y \log\left(\frac{p}{1 - p}\right) + n \log(1 - p)$$

with canonical parameter $\text{logit}(p)$, which is the natural link function for Bernoulli, giving rise to logistic regression model. The binomial distribution is thus an exponential family.

Since the density of the binomial is of the form $p^y(1 - p)^{n-y}$, the beta distribution $\text{Be}(\alpha, \beta)$ with density

$$f(x) \propto x^{\alpha-1}(1 - x)^{\beta-1}$$

is the conjugate prior.

Example 2.2 (Conjugate prior for the Poisson model). The Poisson distribution with mean μ has log density proportional to $f(y; \mu) \propto y \log(\mu) - \mu$, so is an exponential family with natural parameter $\log(\mu)$. The gamma density,

$$f(x) \propto \beta^\alpha / \Gamma(\alpha) x^{\alpha-1} \exp(-\beta x)$$

with shape α and rate β is the conjugate prior for the Poisson. For an n -sample of independent observations $\text{Po}(\mu)$ observations with $\mu \sim \text{Ga}(\alpha, \beta)$, the posterior is $\text{Ga}(\sum_{i=1}^n y_i + \alpha, \beta + n)$.

Example 2.3 (Posterior rates for A/B tests using conjugate Poisson model). Upworthy.com, a US media publisher, revolutionized headlines online advertisement by running systematic A/B tests to compare the different wording of headlines, placement and image and what catches attention the most. The Upworthy Research Archive (Matias et al. 2021) contains results for 22743 experiments, with a click through rate of 1.58% on average and a standard deviation of 1.23%. The `clickability_test_id` gives the unique identifier of the experiment, `clicks` the number of conversion out of `impressions`. See Section 8.5 of Alexander (2023) for more details about A/B testing and background information.

Consider an A/B test from November 23st, 2014, that compared four different headlines for a story on Sesame Street workshop with interviews of children whose parents were in jail and visiting them in prisons. The headlines tested were:

1. Some Don't Like It When He Sees His Mom. But To Him? Pure Joy. Why Keep Her From Him?
2. They're Not In Danger. They're Right. See True Compassion From The Children Of The Incarcerated.
3. Kids Have No Place In Jail ... But In This Case, They *Totally* Deserve It.
4. Going To Jail *Should* Be The Worst Part Of Their Life. It's So Not. Not At All.

2.1 Conjugate priors

At first glance, the first and third headlines seem likely to lead to a curiosity gap. The wording of the second is more explicit (and searchable), whereas the first is worded as a question.

We model the conversion rate λ_i for each headline separately using a Poisson distribution and compare the posterior distributions for all four choices. Using a conjugate prior and selecting the parameters by moment matching yields approximately $\alpha = 1.64$ and $\beta = 0.01$ for the hyperparameters.

Table 2.1: Number of views, clicks for different headlines for the Upworthy data.

headline	impressions	clicks
H1	3060	49
H2	2982	20
H3	3112	31
H4	3083	9

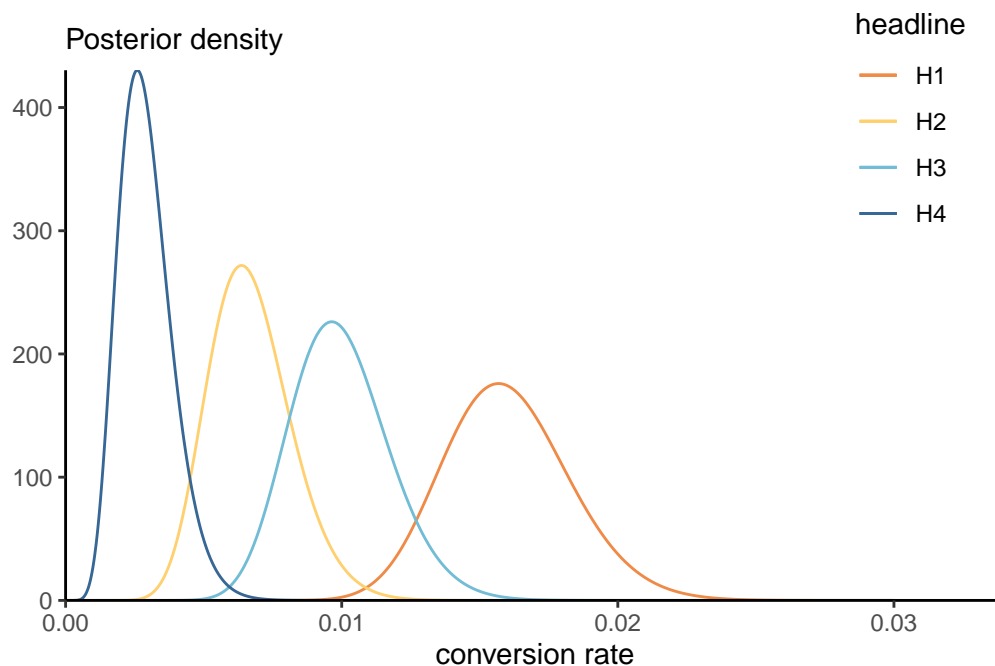


Figure 2.1: Gamma posterior for the Upworthy Sesame street headline.

We can visualize the posterior distributions. In this context, the large sample size lead to the dominance of the likelihood contribution $p(Y_i | \lambda_i) \sim \text{Po}(n_i \lambda_i)$ relative to the prior. We can

2 Priors

see there is virtually no overlap between different rates for headers H1 (preferred) relative to H4 (least favorable). The probability that Headline 3 is better than Headline 1 can be approximated by simulating samples from both posterior and computing the proportion of times one is larger: the probability of superiority is 1.7%, indicating a clear preference for the first headline H1.

Example 2.4 (Conjugate priors in the Bayesian linear model). Consider a linear regression model with observation-specific mean $\mu_i = \mathbf{x}_i \beta$ ($i = 1, \dots, n$) with \mathbf{x}_i the i th row of the $n \times p$ design matrix \mathbf{X} .

Concatenating records, $\mathbf{Y} \sim \text{No}_n(\mathbf{X}\beta, \sigma^2 \mathbf{Q}_y^{-1})$, for a known precision matrix \mathbf{Q}_y , typically \mathbf{I}_n . To construct a conjugate joint prior for $p(\beta, \sigma^2)$, we consider the sequential formulation

$$\beta \mid \sigma^2 \sim \text{No}_p(\nu_\beta, \sigma^2 \mathbf{Q}_\beta^{-1}), \quad \sigma^2 \sim \text{IG}(\alpha, \beta)$$

where IG denotes the inverse gamma distribution¹

The joint posterior is Gaussian-inverse gamma and can be factorized

$$p(\beta, \sigma^2 \mid y) = p(\sigma^2 \mid y) p(\beta \mid \sigma^2, y)$$

where $p(\sigma^2 \mid y) \sim \text{IG}(\alpha^*, \beta^*)$ and $p(\beta \mid \sigma^2, y) \sim \text{No}_p(\mathbf{M}\mathbf{m}, \sigma^2 \mathbf{M})$ with $\alpha^* = \alpha + n/2$, $\beta^* = \beta + 0.5 \nu_\beta^\top \mathbf{Q}_\beta \nu_\beta + \mathbf{y}^\top \mathbf{y} - \mathbf{m}^\top \mathbf{M} \mathbf{m}$, $\mathbf{m} = \mathbf{Q}_\beta \nu_\beta + \mathbf{X}^\top \mathbf{Q}_y \mathbf{y}$ and $\mathbf{M} = (\mathbf{Q}_\beta + \mathbf{X}^\top \mathbf{Q}_y \mathbf{X})^{-1}$; the latter can be evaluated efficiently using Sherman–Morrisson–Woodbury identity.

The exponential family is quite large; Fink (1997) *A Compendium of Conjugate Priors* gives multiple examples of conjugate priors and work out parameter values.

One criticism of the Bayesian approach is the arbitrariness of prior functions. However, the role of the prior is often negligible in large samples (consider for example the posterior of exponential families with conjugate priors). Moreover, the likelihood is also chosen for convenience, and arguably has a bigger influence on the conclusion. Data fitted using a linear regression model seldom follow Gaussian distributions conditionally, in the same way that the linearity is a convenience (and first order approximation).

In general, unless the sample size is small and we want to add expert opinion, we may wish to pick an *uninformative prior*, i.e., one that does not impact much the outcome. For conjugate models, one can often show that the relative weight of prior parameters (relative to the random sample likelihood contribution) becomes negligible by investigating their relative weights.

¹This simply means that the precision σ^{-2} , the reciprocal of the variance, has a gamma distribution with shape α and rate β .

2.2 Uninformative priors

Definition 2.1 (Proper prior). We call a prior *proper* if its integral is finite; such prior function automatically leads to a valid posterior.

The best example of prior priors arise from probability density function. We can still employ this rule for improper priors: for example, taking $\alpha, \beta \rightarrow 0$ in the beta prior leads to a prior proportional to $x^{-1}(1-x)^{-1}$, the integral of which diverges on the unit interval $[0, 1]$. However, as long as the number of success and the number of failures is larger than 1, meaning $k \geq 1, n - k \geq 1$, the posterior distribution would be proper, i.e., integrable. To find the posterior, normalizing constants are also superfluous.

Many uninformative priors are flat, or proportional to a uniform on some subset of the real line and therefore improper. It may be superficially tempting to set a uniform prior on a large range to ensure posterior property, but the major problem is that a flat prior may be informative in a different parametrization, as the following example suggests.

Example 2.5. Consider the parameter $\log(\tau) \in \mathbb{R}$ and the prior $p(\log \tau) \propto 1$. If we reparametrize the model in terms of τ , the new prior (including the Jacobian of the transformation) is τ^{-1}

Some priors are standard and widely used. In location scale families with location ν and scale τ , the density is such that

$$f(x; \nu, \tau) = \frac{1}{\tau} f\left(\frac{x - \nu}{\tau}\right), \quad \nu \in \mathbb{R}, \tau > 0.$$

We thus wish to have a prior so that $p(\tau) = c^{-1}p(\tau/c)$ for any scaling $c > 0$, whence it follows that $p(\tau) \propto \tau^{-1}$, which is uniform on the log scale.

The priors $p(\nu) \propto 1$ and $p(\tau) \propto \tau^{-1}$ are both improper but lead to location and scale invariance, hence that the result is the same regardless of the units of measurement.

Definition 2.2 (Jeffrey's prior). In single parameter models, taking a prior function for θ proportional to the square root of the determinant of the information matrix $p(\theta) \propto \iota(\theta)$ yields a prior that is invariant to parametrization, so that inferences conducted in different parametrizations are equivalent.²

²The Fisher information is linear in the sample size for independent and identically distributed data so we can derive the result for $n = 1$ without loss of generality.

2 Priors

To see this, consider a bijective transformation $\theta \mapsto \vartheta$. Under the reparametrized model and suitable regularity conditions³, the chain rule implies that

$$\begin{aligned} i(\vartheta) &= -\mathbb{E} \left(\frac{\partial^2 \ell(\vartheta)}{\partial^2 \vartheta} \right) \\ &= -\mathbb{E} \left(\frac{\partial^2 \ell(\theta)}{\partial \theta^2} \right) \left(\frac{d\theta}{d\vartheta} \right)^2 + \mathbb{E} \left(\frac{\partial \ell(\theta)}{\partial \theta} \right) \frac{d^2 \theta}{d\vartheta^2} \end{aligned}$$

Since the score has mean zero, $\mathbb{E} \{ \partial \ell(\theta) / \partial \theta \} = 0$, the rightmost term vanishes. We can thus relate the Fisher information in both parametrizations, with

$$i^{1/2}(\vartheta) = i^{1/2}(\theta) \left| \frac{d\theta}{d\vartheta} \right|,$$

implying invariance.

Most of the times, Jeffrey's prior is improper. For the binomial model, it can be viewed as a limiting conjugate beta prior with $\alpha, \beta \rightarrow 0$. Unfortunately, in multiparameter models, the system isn't invariant to reparametrization if we consider the determinant of the Fisher information.

Example 2.6. Consider the binomial distribution $f(y; \theta, n) \propto \theta^y (1 - \theta)^{n-y} \mathbf{1}_{\theta \in [0,1]}$. The negative of the second derivative of the log likelihood with respect to θ is

$$j(\theta) = -\partial^2 \ell(\theta; y) / \partial \theta^2 = y/\theta^2 + (1 - y)/(1 - \theta)^2$$

and since $\mathbb{E}(Y) = n\theta$, the Fisher information is

$$i = \mathbb{E}\{j(\theta)\} = n/\theta + n/(1 - \theta) = n/\{\theta(1 - \theta)\}$$

Jeffrey's prior is thus $p(\theta) \propto \theta^{-1}(1 - \theta)^{-1}$.

Check that for the Gaussian distribution $\text{No}(\mu, \sigma^2)$, the Jeffrey's prior obtained by treating each parameter in turn, fixing the value of the other, are $p(\mu) \propto 1$ and $p(\sigma) \propto 1/\sigma$, which also correspond to the default uninformative priors for location-scale families.

Example 2.7. The Poisson distribution with $\ell(\lambda) \propto -\lambda + y \log \lambda$, with second derivative $-\partial^2 \ell(\lambda) / \partial \lambda^2 = y/\lambda^2$. Since the mean of the Poisson distribution is λ , the Fisher information is $i(\lambda) = \lambda^{-1}$ and Jeffrey's prior is $\lambda^{-1/2}$.

³Using Bartlett's identity; Fisher consistency can be established using the dominated convergence theorem.

2.3 Expert knowledge

The prior distribution may have parameters themselves that need to be specified by experts. One may also wish to add another layer and set an hyperprior distribution on the parameters, resulting in a hierarchical model.

Setting parameters of priors is often done by reparametrizing the latter in terms of moments. Sometimes, it may be easier to set priors in a different scale where subject-matter expertise is most easily elicited.

Example 2.8. The generalized extreme value distribution arises as the limiting distribution for the maximum of m independent observations. The $\text{GEV}(\mu, \sigma, \xi)$ distribution is a location-scale with distribution function

$$F(x) = \exp \left[- \{1 + \xi(x - \mu)/\sigma\}_+^{-1/\xi} \right]$$

where $x_+ = \max\{0, x\}$.

Inverting the distribution function yields the quantile function

$$Q(p)\mu + \sigma \frac{(-\log p)^{-\xi} - 1}{\xi}$$

In environmental data, we often model annual maximum. Engineering designs are often specified in terms of the k -year return levels, defined as the quantile of the annual maximum exceeded with probability $1/k$ in any given year. Using a GEV for annual maximum, Coles and Tawn (1996) proposed modelling annual daily rainfall and specifying a prior on the quantile scale $q_1 < q_2 < q_3$ for tail probabilities $p_1 > p_2 > p_3$. To deal with the ordering constraints, gamma priors are imposed on the differences $q_1 - o \sim \text{Ga}(\alpha_1, \beta_1)$, $q_2 - q_1 \sim \text{Ga}(\alpha_2, \beta_2)$ and $q_3 - q_2 \sim \text{Ga}(\alpha_3, \beta_3)$, where o is the lower bound of the support. The prior is thus of the form

$$p(\mathbf{q}) \propto q_1^{\alpha_1-1} \exp(-\beta_1 q_1) \prod_{i=2}^3 (q_i - q_{i-1}^{\alpha_i-1} \exp\{\beta_i(q_i - q_{i-1})\}).$$

where $0 \leq q_1 \leq q_2 \leq q_3$. We can then relate the prior parameters to moments.

Consider the annual maximum rainfall in Abisko, Sweden.

Example 2.9 (Prior simulation). Are the prior reasonable? One way to see this is to sample values from the priors and generate new observations.

Example 2.10. We can specify gamma back-transform them to location μ , scale σ and shape ξ and simulate observations from the $\text{GEV}(\mu, \sigma, \xi)$ and compare them to observations.

References

- Alexander, Rohan. 2023. *Telling Stories with Data: With Applications in R*. Boca Raton, FL: CRC Press.
- Coles, Stuart G., and Jonathan A. Tawn. 1996. "A Bayesian Analysis of Extreme Rainfall Data." *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 45 (4): 463–78. <https://doi.org/10.2307/2986068>.
- Finetti, Bruno de. 1974. *Theory of Probability: A Critical Introductory Treatment*. Vol. 1. New York: Wiley.
- Gelfand, Alan E., and Adrian F. M. Smith. 1990. "Sampling-Based Approaches to Calculating Marginal Densities." *Journal of the American Statistical Association* 85 (410): 398–409. <https://doi.org/10.1080/01621459.1990.10476213>.
- Geman, Stuart, and Donald Geman. 1984. "Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images." *IEEE Transactions on Pattern Analysis and Machine Intelligence* PAMI-6 (6): 721–41. <https://doi.org/10.1109/TPAMI.1984.4767596>.
- Matias, J. Nathan, Kevin Munger, Marianne Aubin Le Quere, and Charles Ebersole. 2021. "The Upworthy Research Archive, a Time Series of 32,487 Experiments in U.S. Media." *Scientific Data* 8 (195). <https://doi.org/10.1038/s41597-021-00934-7>.
- McNeil, A. J., R. Frey, and P. Embrechts. 2005. *Quantitative Risk Management: Concepts, Techniques, and Tools*. 1st ed. Princeton, NJ: Princeton University Press.

