

Jeffreys priors

Lecturer: Michael I. Jordan

Scribe: Timothy Hunter

1 Priors for the multivariate Gaussian

Consider a multivariate Gaussian variable X of size p . Its probability density function can be parametrized by a mean vector $\mu \in \mathbb{R}^p$ and a covariance matrix $\Sigma \in \mathcal{S}_p^+$:

$$p(X|\mu, \Sigma) = \frac{1}{(2\pi)^{p/2}} \frac{1}{\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(X - \mu)^T \Sigma^{-1} (X - \mu)\right) \quad (1)$$

We will consider three cases of conjugate priors: the case when the covariance is fixed, the case when the mean is fixed and the general case.

1.1 The case of fixed variance

The conjugate prior is a multivariate Gaussian of mean μ_0 and covariance matrix Σ_0 . The derivations are the same as in the univariate case.

1.2 The case of fixed mean

The conjugate prior is the *inverse Wishart distribution*. One can see this using the trace trick:

$$(X - \mu)^T \Sigma^{-1} (X - \mu) = \text{Tr}\left(\Sigma^{-1} (X - \mu) (X - \mu)^T\right) \quad (2)$$

This is an inner product in the matrix space between the matrices Σ^{-1} and $(X - \mu) (X - \mu)^T$. This particular inner product corresponds to summing all the pairwise products of elements of each matrix. We will derive here the Wishart distribution, and you can derive the inverse Wishart by a change of variable. The computations are similar for the inverse Wishart distribution.

The distribution over the precision matrix is the Wishart distribution with one degree of freedom is:

$$p(M|V) = \frac{|M|^{p/2} e^{-\frac{1}{2}\text{Tr}(MV^{-1})}}{2^{p/2} |V|^{1/2} \Gamma_p\left(\frac{1}{2}\right)} \quad (3)$$

where Γ_p is the generalized (multivariate) Gamma function, and M, V are positive definite matrices. When we forget the normalization constants, we recognize the product of the determinant (up to some exponent) with the exponential of the trace, which is the familiar form of the multivariate Gaussian:

$$p(M|V) \propto |M|^{p/2} e^{-\frac{1}{2}\text{Tr}(MV^{-1})} \quad (4)$$

This is the multivariate analog of a Gamma distribution. In order to work with the covariance matrix and get the inverse Wishart distribution, one has to apply the change of variable $\Sigma = P^{-1}$. This shape of the inverse Wishart looks very close to that of the inverse gamma:

$$p(\Sigma|V) \propto |\Sigma|^{-(p+1)/2} e^{-\frac{1}{2}\text{Tr}(\Sigma^{-1}V^{-1})} \quad (5)$$

If we want to get a predictive distribution, we integrate the inverse Wishart against the multivariate Gaussian which gives the multivariate Student distribution:

$$T \propto \frac{1}{\left(1 + (X - \mu)^T \Sigma^{-1} (X - \mu)\right)^{p/2}} \quad (6)$$

with a complicated with a heavy tail.

1.3 The general case

The computations are the same as before with an inverse Wishart for the covariance and a scaled Gaussian (scaled by the Wishart).

1.4 Sampling from the Wishart distribution: the Bartlett decomposition

If one needs to sample from the Wishart, there is a nice way to sample it called the Bartlett decomposition. Consider the Cholesky decomposition of the parameter:

$$V = LL^T \quad (7)$$

Samples of Σ are obtained by sampling

$$\Sigma = LZZ^T L^T \quad (8)$$

where

$$Z = \begin{pmatrix} \sqrt{c_1} & & & & \\ z_{21} & \sqrt{c_2} & & & \\ z_{31} & & \sqrt{c_3} & & \\ \vdots & & & \ddots & \\ z_{p1} & z_{p2} & & \dots & \sqrt{c_p} \end{pmatrix} \quad (9)$$

in which the diagonal coefficients are from the χ^2 distribution with p degrees of freedom and the z_{ij} are from the univariate Gaussian distribution $\mathcal{N}(0, 1)$.

There is a similar way to sample from the multivariate Gaussian distribution. Consider the multivariate Gaussian with identity matrix $X \sim \mathcal{N}(0, I_p)$. This is easy to sample from: each coefficient can be sampled independently by a univariate Gaussian. We use the Cholesky (or the square root) decomposition of the covariance matrix

$$\Sigma = LL^T \quad (10)$$

We then define a new random variable $W = LX$, with 0 mean and covariance $\text{Var}(W) = \mathbb{E}[WW^T] - 0 = LL^T = \Sigma$. Therefore W as defined this way can be described a 0-mean Gaussian with covariance Σ . Getting a different mean is simply a matter of translation W .

This concludes our introduction to conjugate priors. Conjugate priors are a matter of convenience, easy to implement and as such widely used in software implementations. They have some nice properties, in particular they are optimal asymptotically. They are often used in applications, when one lacks prior knowledge. Using conjugate priors, only needs to assess the prior parameters.

2 Jeffreys priors

Though conjugate priors are computationally nice, objective Bayesians instead prefer priors which do not strongly influence the posterior distribution. Such a prior is called an *uninformative prior*.

This is a hard problem, and a number of things we might try are not appropriate. The historical approach, followed by Laplace and Bayes, was to assign flat priors. This prior seems reasonably uninformative. We do not know where the actual value lies in the parameter space, so we might as well consider all values equiprobable. This prior however is not invariant. Consider for example a binomial distribution $X \sim \text{Binom}(n, \theta)$ in which we want to put a prior on θ . We know that θ lies between 0 and 1. The flat prior on θ is the uniform distribution: $\pi(\theta) = 1$. Since θ lies between 0 and 1, we can use a new parametrization using the log-odds ratio: $\rho = \log \frac{\theta}{1-\theta}$. This is a perfectly valid parametrization, and a natural one if we want to map θ to the full scale of the reals. Under this parametrization the prior distribution $\pi(\rho)$ is not flat anymore. This example shows a prior that is uninformative in one parametrization, but becomes informative through a change of variables.

This becomes more problematic in higher dimensions: the uniform prior in large dimension does not integrate anymore. In addition, the flat prior becomes very informative: it tells that most of the probability mass lies at $+\infty$, far from the origin. If instead one considers a high-dimensional Gaussian distribution $X \sim \mathcal{N}(0, 1)$, most of the mass is concentrated in a (high dimensional) unit sphere centered at the origin.

Faced with these issues, we see that flat priors and uninformative priors raise mathematical and philosophical problems. These examples show that finding prior distributions that have a minimal impact as possible on the data raises deep practical issues.

We first consider some special cases in one dimension, then consider the general case.

2.1 Examples

2.1.1 The example of an uninformative location prior

Consider the case where we have a location parameter: a probability distribution over a variable X of density $f(X - \theta)$ where θ is a *location parameter* that we endow with a prior. A candidate for a prior would be $\pi(\theta) \propto 1$. If θ lies in an interval, we can consider the uniform distribution as a prior estimate. If θ can take any value in \mathbb{R} , the flat prior is not a probability density because it does not integrate. Such a prior is called an *improper prior*. It expresses our state of ignorance (hence the flat prior) and can be defined as the limit of a proper prior.

2.1.2 The example of an uninformative scaling prior

Consider a density factor θ :

$$f_{\theta}(x) = \frac{1}{\theta} f\left(\frac{x}{\theta}\right) \quad (11)$$

The $\frac{1}{\theta}$ in the front ensures that f_θ still integrates to 1. The prior on θ should be invariant to rescaling by any arbitrary positive constant, i.e.:

$$\pi(\theta) = \frac{1}{c} \pi\left(\frac{\theta}{c}\right) \quad (12)$$

for all $c > 0$. This means if we rescale our variable, the prior will not change, that is, the prior does not give any information when we rescale the variable. The previous relation is a functional equation that admits a single solution for π (up to a scaling factor):

$$\pi(\theta) \propto \frac{1}{\theta} \quad (13)$$

Note how it is an improper prior because it does not have a finite integral. This is the uninformative scale prior.

We now look for a space in which this prior transforms into a flat prior. Consider the change of variable $\rho = \log \theta$ (or equivalently $\theta = e^\rho$). Then the probability density function in the new parametrization is:

$$\begin{aligned} \pi(\rho) &= \pi(\theta) \left| \frac{d\theta}{d\rho} \right| \\ &\propto e^{-\rho} e^\rho = 1 \end{aligned}$$

Thus our scale invariant prior is actually a flat prior in the log scale. It treats equally any order of magnitude.

This prior has another derivation based on the (proper) conjugate prior of the variance of the Gaussian. We saw that the conjugate prior for the variance of the Gaussian is the inverse gamma:

$$p(\sigma^2 | \alpha, \beta) \propto (\sigma^2)^{-(\alpha+1)} e^{-\beta/\sigma^2} \quad (14)$$

which is parametrized by two parameters α and β . The parameter α can be interpreted as the number of observations and β is some inverse concentration of the parameter in a certain region. It then makes sense to take α and β to the 0 limit, as if we had no prior information. When we do that we obtain:

$$p(\sigma^2 | \alpha, \beta) \propto \frac{1}{\sigma^2} \quad (15)$$

which is the same prior as the one we derived.

2.2 Jeffreys priors

Uninformative priors we have seen so far are appealing because they are flat priors in some meaningful parametrization. Jeffreys priors are a generalization of these ideas, and can deliver a broad range of priors that incorporates these special cases. They are quite reasonable in one dimension. They are based on a principle of invariance: one should be able to apply these priors to certain situations, apply a change of variable, and still get the same answer. Suppose we are provided with some model and some data, i.e. with a likelihood function $p(x|\theta)$. One should be able to manipulate the likelihood and get a prior on θ , from the likelihood only. Note how this approach goes contrary to the subjective Bayesian frame of mind, in which one first chooses a prior on then θ and then applies it to the likelihood to derive the posterior.

The answer in the one-dimensional case is:

$$\pi_J(\theta) \propto \mathbf{I}(\theta)^{1/2} \quad (16)$$

in which \mathbf{I} is the *Fisher information*, defined when θ is unidimensional by the second derivative of the log likelihood:

$$\mathbf{I}(\theta) = -\mathbb{E}_\theta \left[\frac{d^2 \log p(X|\theta)}{d\theta^2} \right] \quad (17)$$

This is an integral over the values of X with keeping θ fixed (the expectation in the frequentist sense). Using the maximum likelihood principle, the best parameter θ cancels the first derivative of the log likelihood, and the second derivative gives the curvature of the likelihood around the MLE.

$\frac{d^2 \log p(X|\theta)}{d\theta^2}$ is a random variable over X and \mathbb{E}_θ denotes the fact that we are integrating with respect to the distribution f_θ indexed by the (fixed) variable θ : $f_\theta(X) = p(X|\theta)$. The Fisher information is locally concave around the MLE, globally concave for the exponential family, but not globally concave for all distributions. This is not the case for example for mixture models.

We check that it works for a case we already saw: the Gaussian with fixed variance $X \sim \mathcal{N}(\mu, \sigma^2)$ for which we want to get prior on the location parameter μ . The likelihood is:

$$p(X|\mu) \propto \exp\left(-\frac{1}{2\sigma^2}(X - \mu)^2\right) \tag{18}$$

thus when we take the second derivative with respect to μ :

$$\frac{d^2 \log p(X|\mu)}{d\mu^2} = -\frac{1}{\sigma^2} \tag{19}$$

which is a constant with respect to the random variable X and μ , so when we take the expectation, we get the flat prior we obtained before:

$$\mathbf{I}(\mu) \propto 1 \tag{20}$$

Now that we saw the answer, here are some explanations as to where the result comes from. Let us define a new parameter $\phi = h(\theta)$ as a reparametrization. If we calculate π_J with respect to the variable θ and then transform variables, this will give a prior π on ϕ by the change of variable formula. The question is thus to check if this prior $\pi(\phi)$ is indeed the Jeffreys prior $\pi_J(\phi)$ that we would have computed in the first place by using the variable ϕ . We apply Jeffreys' principle in the ϕ space by using the chain rule and reexpress in terms of θ :

$$\begin{aligned} \mathbf{I}(\phi) &= -\mathbb{E}\left[\frac{d^2 \log p(X|\phi)}{d\phi^2}\right] \\ &= -\mathbb{E}\left[\frac{d^2 \log p(X|\theta)}{d\theta^2} \left(\frac{d\theta}{d\phi}\right)^2 + \frac{d \log p(X|\theta)}{d\theta} \frac{d^2 \theta}{d\phi^2}\right] \\ &= -\mathbb{E}\left[\frac{d^2 \log p(X|\theta)}{d\theta^2}\right] \left(\frac{d\theta}{d\phi}\right)^2 + \mathbb{E}\left[\frac{d \log p(X|\theta)}{d\theta}\right] \frac{d^2 \theta}{d\phi^2} \end{aligned}$$

The previous formulas are simply in application of the chain rule. We know:

$$\mathbb{E}\left[\frac{d \log p(X|\theta)}{d\theta}\right] = 0 \tag{21}$$

One way to see this fact is to use the total probability:

$$\forall \theta, \int p(X|\theta) dX = 1 \tag{22}$$

Assuming sufficient regularity, when we take the derivative with respect to θ :

$$\begin{aligned}
 0 &= \frac{d}{d\theta} \int p(X|\theta) dX \\
 &= \int \frac{dp(X|\theta)}{d\theta} \frac{p(X|\theta)}{p(X|\theta)} dX \\
 &= \int \left[\frac{dp(X|\theta)}{d\theta} \frac{1}{p(X|\theta)} \right] p(X|\theta) dX \\
 &= \int \left[\frac{d \log p(X|\theta)}{d\theta} \right] p(X|\theta) dX \\
 &= \mathbb{E} \left[\frac{d \log p(X|\theta)}{d\theta} \right]
 \end{aligned}$$

This is formally proved using the dominated convergence theorem on distributions. Using this result, taking the expectation over X with θ fixed is equivalent to take the expectation with ϕ fixed, so we get:

$$\mathbf{I}(\phi) = \mathbf{I}(\theta) \left(\frac{d\theta}{d\phi} \right)^2 \tag{23}$$

and by taking the square root:

$$\sqrt{\mathbf{I}(\phi)} = \sqrt{\mathbf{I}(\theta)} \left| \frac{d\theta}{d\phi} \right| \tag{24}$$

By the change of variable formula, this shows that the Jeffreys prior $\pi_J(\theta) = \sqrt{\mathbf{I}(\theta)}$ is invariant to a change of variable.