

7.4. Choosing priors

Important issue - how to reflect the degree of prior belief? How subjective is the choice?

Notation: family of priors $\{f_\lambda : \lambda \in \Lambda\}$, selecting a prior then means selecting $\lambda \in \Lambda$. λ is a hyperparameter.

7.4.1 Conjugate priors

Def. The family of priors $\{f_\lambda : \lambda \in \Lambda\}$ for the parameter θ of the model $\{f_\theta : \theta \in \Theta\}$ is conjugate, if for all $s \in S$ and $\lambda \in \Lambda$ the posterior $\{f_{\lambda+s} : s \in S\} \in \{f_\lambda : \lambda \in \Lambda\}$.

Example: conjugate families we described earlier.

7.4.2 Elicitation:

- means using the statistician's beliefs about θ to specify the prior in $\{f_\lambda : \lambda \in \Lambda\}$ that reflects these beliefs.

Ex. 7.4.2. Location Normal

Data - $N(\mu, \sigma^2)$; restrict to $N(\mu_0, \sigma_0^2)$: $\mu_0 \in \mathbb{R}^1, \sigma_0^2 > 0$

the priors for μ . So, $\lambda = (\mu_0, \sigma_0^2)$, 2 df.

E.g. We can ask an expert to specify two quantiles of the prior distr. for $\mu \Rightarrow$ defines the prior. We may want to specify μ_0 such that f_λ is as likely to greater than as less than μ_0 , so that μ_0 is the median of the prior. We may also want to specify V_0 s.t. there is 99% certainty that the true $\mu < V_0$. - that is 0.99 quantile of the prior.

7.4.2 Empirical Bayes:

Here, the main idea is to base the choice of λ on the data s . Note this approach seems to violate the principle of conditional probability. Indeed, in this case the posterior distribution of θ given s is not, in general, the conditional distribution of θ given data

One possibility is to compute the prior predictive $m_\lambda(s)$ for the data s , and then base the choice of λ on these values.

Note that the prior predictive is like a likelihood function for λ so that it is logical to select λ that maximizes $m_\lambda(s)$.

Ex. 7.4.3. Bernoulli model

$(X_1, \dots, X_n) \sim \text{Ber}(\theta)$, $\theta \sim \text{Beta}(\lambda, \lambda)$ for some $\lambda > 0$

The prior is symmetric around $\frac{1}{2}$. Check that the prior mean is $\frac{1}{2}$ and the prior variance is $\frac{\lambda^2}{(2\lambda+1)(2\lambda)^2} = \frac{1}{4(2\lambda+1)} \rightarrow 0$

as $\lambda \rightarrow \infty$. Choosing large λ makes for a very precise prior.

$$\text{Then, } m_\lambda(X_1, \dots, X_n) = \frac{\Gamma(2\lambda)}{\Gamma^2(\lambda)} \int_0^{n\bar{x}+\lambda-1} \theta^{n\bar{x}+\lambda-1} (1-\theta)^{\lambda-1} d\theta =$$

$$= \frac{\Gamma(2\lambda)}{\Gamma^2(\lambda)} \frac{\Gamma(n\bar{x}+\lambda) \Gamma(n(1-\bar{x})+\lambda)}{\Gamma(n+2\lambda)} \quad \text{Maximisation}$$

w.r.t λ has to be numerical.

E.g. $n=20$, $n\bar{x}=2.3$ - the number of 1's out of 20.

See Fig. 7.4.1 - max is approx. $\lambda=2.3$

7.4.4. Hierarchical Bayes.

Here, yet another (hyper) prior is put on λ .

Then, the prior for θ becomes $\pi_\lambda(\theta) = \int \pi_\lambda(\theta) w(\lambda) d\lambda$

Typically, default choices are made for w .

$$\pi_\lambda(s) = \frac{f_\theta(s) \int \pi_\lambda(\theta) w(\lambda) d\lambda}{m(s)} = \int \frac{f_\theta(s) \pi_\lambda(\theta) m_\lambda(s) w(\lambda)}{m_\lambda(s) m(s)} d\lambda$$

$$\text{Where } m(s) = \int \int f_\theta(s) \pi_\lambda(\theta) w(\lambda) d\theta d\lambda = \int m_\lambda(s) w(\lambda) d\lambda \text{ and}$$

$$m_\lambda(s) = \int f_\theta(s) \pi_\lambda(\theta) d\theta \quad [\text{if } \lambda \text{ is continuous with a prior density given by } w]$$

Note that the posterior density of λ is $m_\lambda(s)w(\lambda)$
Write $\frac{f_\theta(s)\delta_\lambda(\theta)}{m_\lambda(s)}$ is the posterior density of $M(s)$

& given λ . Thus, we can use $\delta_\theta(s)$ for inferences
about θ and $\frac{m_\lambda(s)w(\lambda)}{m(s)}$ for inferences about λ .

Ex. 7.4.2 Location-scale normal -

Data $X \sim N(\mu, \sigma^2)$ where $\mu \sim N(\mu_0, \sigma_0^2)$

$\sigma^2 \sim \text{Gamma}(\alpha_0, \beta_0)$

Just need to have one more prior $w(\mu_0, \sigma_0^2, \alpha_0, \beta_0)$ -

7.4.5 Improper Priors and non Informativity.

Where to stop the chain of priors in hierarchical Bayes? One possibility is to stop at a noninformative prior (also, default prior or reference prior). Just how to express that ignorance is often difficult to say. In some cases, the default prior is improper.

$\int \pi(\theta) d\theta = +\infty$. How to interpret this, is unclear.
It also implies that (s, θ) no longer has a joint prob distribution; the use of the principle of conditional probability to justify our inference based on posterior also doesn't work. Need to make sure that at least the posterior is proper.

Ex. 7.4.5 Location-Normal model.

Let $(x_1, x_2, \dots, x_n) \sim N(\mu, \sigma_0^2)$; choose $w(\mu) \propto 1$

\Rightarrow the posterior density of μ is $\propto e^{-\frac{n}{2\sigma_0^2}(x-\mu)^2}$

$\Rightarrow \mu | s \sim N\left(\bar{x}, \frac{\sigma_0^2}{n}\right)$. This is the same as the limiting posterior in Ex. 7.1.2 as $\sigma_0 \rightarrow \infty$.

One common method of default prior selection is
to use Jeffreys' prior $\propto \frac{1}{\theta^2} T^{1/2}(\theta)$ for $\theta \in \mathbb{R}^+$.
It is dependent on the model. It is often improper

Ex. 7.4.6. Location normal model -

Rec $p(\theta) = \frac{\sqrt{n}}{\theta}$ - the same posterior is obtained as
in Ex. 7.4.5. This is because the Jeffreys' prior
is a constant w.r.t. θ .