

Ex. 7.2.4 Multinomial model

Sample (s_1, s_2, \dots, s_n) ; Dirichlet prior $(\alpha_1, \alpha_2, \dots, \alpha_k)$ on $(\theta_1, \dots, \theta_{k-1}) \Rightarrow$ the posterior distribution is

Dirichlet $(x_1 + \alpha_1, x_2 + \alpha_2, \dots, x_k + \alpha_k)$.

Can show mathematically that, if $(\theta_1, \dots, \theta_{k-1}) \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_k)$ then $\theta_i \sim \text{Dirichlet}(\alpha_i, \alpha_{-i}) = \text{Beta}(\alpha_i, \alpha_{-i})$,

where $\alpha_{-i} = \alpha_1 + \alpha_2 + \dots + \alpha_k - \alpha_i \Rightarrow$ the marginal posterior of θ_1 is $\text{Beta}(x_1 + \alpha_1, x_2 + \alpha_2 + \dots + x_k + \alpha_k)$

Assuming that each $\alpha_i \geq 1$, & keeping in mind that $x_1 + \alpha_1 + \dots + x_k + \alpha_k = n$, the posterior mode is $\hat{\theta}_1 = \frac{x_1 + \alpha_1}{n - 2 + \alpha_1 + \dots + \alpha_k}$

For the uniform prior: $\alpha_1 = \dots = \alpha_k = 1$

$\hat{\theta}_1 = \frac{x_1}{n + k - 2}$

The posterior expectation is $E(\theta_1 | x) = \frac{x_1 + \alpha_1}{n + \alpha_1 + \dots + \alpha_k}$

The plug-in MLE of θ_1 is $\frac{x_1}{n}$ which is close to $\hat{\theta}_1 + 1$

Bayesian estimator for large n .

How to assess the accuracy of Bayesian estimates?

For posterior mean, a possible approach is to compute

Posterior variance. As an example, in Ex. 7.2.3, the posterior variance is $(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2})^{-1}$

and $\rightarrow \frac{\sigma_0^2}{n} = \text{Var}(\bar{x})$ as $\tau_0^2 \rightarrow \infty$

7.2.2. Credible intervals

A credible interval for $\psi(\theta)$ is an interval $C(s) = [l(s), u(s)]$ that we believe will contain a true value of ψ . Typically, we specify $\delta \in (0, 1)$ first. Then, we find an interval $C(s)$ such that $\mathbb{P}(\psi(\theta) \in C(s) | s) = \mathbb{P}(\{\theta : l(s) \leq \psi(\theta) \leq u(s)\} | s) \geq \delta$. Then, $C(s)$ is a δ -credible interval for ψ .

Usually, the goal is to find a δ -credible interval $C(s)$ such that 1) $\mathbb{P}(\psi(\theta) \in C(s) | s)$ is as close to δ as possible and 2) $C(s)$ is as narrow as possible. One option to view

the highest posterior density (HPD) intervals that are $C(s) = \{\psi : w(\psi | s) \geq c\}$

where c is chosen as large as possible to satisfy (1)

Of course, $C(s)$ contains the mode whenever

$c \leq \max_{\psi} w(\psi | s)$. The width of such an interval is the measure of the accuracy of the mode of $w(\cdot | s)$ as an estimator of $\psi(\theta)$.

Ex 7.2.7 Location normal model.

$$x_1, \dots, x_n \sim N(\mu, \sigma_0^2); \mu \sim N(\mu_0, \tau_0^2)$$

Here the posterior mean is $\hat{\mu}$ and the posterior variance is

$$\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1}, \text{ the distribution is symmetric around its mode, } \hat{\mu}$$

$\hat{\mu}$, the shortest δ -HPD interval is

$$\hat{\mu} \pm \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1/2} c$$

where c is chosen in such a way that

$$\delta = \mathbb{P}\left(\mu \in \left[\hat{\mu} \pm \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1/2} c\right] \mid x_1, \dots, x_n\right) = \mathbb{P}\left(-c \leq \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{1/2} (\mu - \hat{\mu}) \leq c \mid x_1, \dots, x_n\right)$$

Because $\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{1/2} (\mu - \hat{\mu}) | x_1, \dots, x_n \sim N(0, 1)$,
 we find that $f = \Phi(c) - \Phi(-c)$. Therefore,

$c = \frac{z_{1-f}}{2}$, and the interval is

$$\hat{\mu} \pm \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1/2} \frac{z_{1-f}}{2}$$

Note that as $\tau_0^2 \rightarrow \infty$, this interval becomes $\bar{x} \pm \frac{\sigma_0}{\sqrt{n}} \frac{z_{1-f}}{2}$

There are other possible types of credible intervals. One common method is to take $[y_l, y_r]$ where y_l is a $\frac{1-f}{2}$ quantile of the posterior distribution and y_r is a $\frac{1+f}{2}$ quantile for it. These tend to avoid some expensive computation that is needed for HPD intervals.

Additional example. The same farmer wants to construct a 95% credible set for the average weight of the cows μ . Recall that the posterior mean $E(\mu | x_1, \dots, x_n) = 412.9$ kg while the posterior variance is $\text{Var}(\mu | x_1, \dots, x_n) = \left(\frac{10}{400} + \frac{1}{100}\right)^{-1} = 28.6$ kg². Thus, the posterior distribution is $\mu | x_1, \dots, x_n \sim N(412.9, 28.6)$ and the 95% credible interval is

$$412.9 \pm z_{0.025} \sqrt{28.6} = (402.4, 423.4)$$

7.2.3. Hypothesis Testing and Bayesian Factors

Let $H_0: \psi(\theta) = \psi_0$. Can compute the posterior probability $\pi(\psi(\theta) = \psi_0 | s)$ and if it is small, we have the evidence against H_0 .

Ex. 7.2.9 - Let $H_0: \theta \in A$. Let $\psi = \mathbb{1}_A$, we can rewrite $H_0: \psi(\theta) = 1$ and $\pi(\psi(\theta) = \mathbb{1}_A | s) = \pi(A | s)$. The approach (*) has its problems. If the posterior distribution is a continuous one, then $\pi(\psi(\theta) = \psi_0 | s) = 0$ for any s .

To avoid this problem, another approach is often used. Realize that a region of low posterior probability occurs where the posterior density $w(\cdot | s)$ is relatively low. This suggests computing a Bayesian P-value

$$\pi(\{\theta: w(\psi(\theta) | s) \leq w(\psi_0 | s)\} | s) \quad (**)$$

If $w(\cdot | s)$ is unimodal, this is a tail probability. If it is small, then ψ_0 is surprising according to our posterior beliefs. If we decide to reject H_0 whenever P-val. $< 1 - \alpha$, this is equivalent to computing a $1 - \alpha$ HPD region for ψ and rejecting H_0 whenever ψ_0 is not in the region.

Ex. 7.2.10 Clearly, $\psi(\theta) = \mathbb{1}_{A^c}(\theta)$ has the posterior $\text{Ber}(\pi(A) | s)$

Thus, $w(\cdot | s)$ is defined as $w(0 | s) = 1 - \pi(A | s) = \pi(A^c | s)$
 $w(1 | s) = \pi(A | s)$

$$\begin{aligned} \text{For } \psi_0 = 1, \text{ so } \{\theta: w(\psi(\theta) | s) \leq w(1 | s)\} &= \{\theta: w(\mathbb{1}_{A^c}(\theta) | s) \leq \pi(A | s)\} = \\ &= \begin{cases} \Omega, & \pi(A | s) \geq \pi(A^c | s) \\ A, & \pi(A | s) < \pi(A^c | s) \end{cases} \end{aligned}$$

Therefore, (**) becomes

$$\pi(\{\theta: w(\psi(\theta) | s) \leq w(1 | s)\} | s) = \begin{cases} 1, & \pi(A | s) \geq \pi(A^c | s) \\ \pi(A | s), & \pi(A | s) < \pi(A^c | s) \end{cases}$$

- again we have evidence against H_0 whenever $\pi(A | s)$ is small -

Thus, we find that using $(**)$ is equivalent to using $(*)$ whenever the parameter $\psi(\theta)$ takes only two values. When the prior distribution is continuous, $(*)$ doesn't work. But $(**)$ has its own issues in this situation.

Ex. 7.2.11. Let the posterior dist. of θ be $w(\theta|s) = 2\theta$ when $0 \leq \theta \leq 1$; $H_0: \theta = \frac{3}{4}$. Then, $w(\theta|s) \leq w(\frac{3}{4}|s)$ iff $\theta \leq \frac{3}{4}$; $(**)$ becomes $\int_0^{\frac{3}{4}} 2\theta d\theta = \frac{9}{16}$.

On the other hand, 1-1 transformation $p = \theta^2$ brings $H_0: p = \frac{9}{16}$.

The posterior distr. of p is $w(p|s) = 1, 0 \leq p \leq 1$; thus it is trivially true that $w(p|s) \leq w(\frac{9}{16}|s) \Rightarrow (*)$ is equal to 1, no evidence against H_0 can be found using this parameterization. Thus, the dependence on parameterization seems inappropriate.

To avoid all of these difficulties, it seems advisable to assign a positive prior probability to the hypothesis H_0 . Then, Ex. 7.2.10 demonstrates that this is the same as using $(*)$ to assess H_0 . Usually, one assigns the prior

$$\Pi = p\Pi_1 + (1-p)\Pi_2,$$

where $\Pi_1(\psi(\theta) = \psi_0) = 1$, $\Pi_2(\psi(\theta) = \psi_0) = 0$, so that Π_1 is degenerate at ψ_0 and Π_2 is continuous at ψ_0 .

Then, $\Pi(\psi(\theta) = \psi_0) = p\Pi_1(\psi(\theta) = \psi_0) + (1-p)\Pi_2(\psi(\theta) = \psi_0) = p > 0$ - this is the prior probability

that H_0 is true.

The prior predictive (marginal) for the data s is

$$m(s) = pm_1(s) + (1-p)m_2(s)$$

where m_i is the prior predictive obtained via prior Π_i .

This implies that the posterior probability measure for A is

$$\Pi(A|s) = \frac{pm_1(s)}{pm_1(s) + (1-p)m_2(s)} \Pi_1(A|s) + \frac{(1-p)m_2(s)}{pm_1(s) + (1-p)m_2(s)} \Pi_2(A|s)$$

where $\Pi(A|s)$ is the posterior measure obtained via the prior Π .

$$\pi(\psi|\theta) = \psi_0(s) = \frac{p m_1(s)}{p m_1(s) + (1-p) m_2(s)} \text{ and}$$

we use this probability to assess H_0 .

Ex. 7.2.12. Location Normal Model

Let $(x_1, \dots, x_n) \sim N(\mu, \sigma_0^2)$; $H_0: \mu = \mu_0$.

Prior for μ is $N(\mu_0, \tau_0^2)$ this will be our $\pi_2 \sim N(\mu_0, \tau_0^2)$

Consider now $\pi = p\pi_1 + (1-p)\pi_2$; take $\pi_1(\{\mu_0\}) = 1$,

and so $\pi(\{\mu_0\}) = p$. Recall that under π_2 the posterior

distribution of μ is $N\left(\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1} \left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma_0^2} \bar{x}\right), \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1}\right)$

and under π_1 the posterior is the distribution degenerate at μ_0 .

To be continued.

Remark. The prior predictive probability measure for the data s with a mixture of π_1 and π_2 prior distribution is

$$\begin{aligned} m(s) &= \mathbb{E}_{\pi} f_0(s) = \sum_{\theta} f_0(s) \pi(\{\theta\}) = \sum_{\theta} f_0(s) (p\pi_1(\{\theta\}) + (1-p)\pi_2(\{\theta\})) \\ &= p \sum_{\theta} f_0(s) \pi_1(\{\theta\}) + (1-p) \sum_{\theta} f_0(s) \pi_2(\{\theta\}) \\ &= p f_0(s) + (1-p) \sum_{\theta} f_0(s) \pi_2(\{\theta\}) = p m_1(s) + (1-p) m_2(s). \end{aligned}$$

The posterior ~~predictive~~ probability is given by

$$\begin{aligned} \pi(A|s) &= \sum_{\theta \in A} \frac{f_0(s) \pi(\{\theta\})}{m(s)} = \sum_{\theta \in A} \frac{f_0(s) [p\pi_1(\{\theta\}) + (1-p)\pi_2(\{\theta\})]}{p m_1(s) + (1-p) m_2(s)} \\ &= \frac{p m_1(s)}{p m_1(s) + (1-p) m_2(s)} \sum_{\theta \in A} \frac{f_0(s) \pi_1(\{\theta\})}{m_1(s)} + \frac{(1-p) m_2(s)}{p m_1(s) + (1-p) m_2(s)} \sum_{\theta \in A} \frac{f_0(s) \pi_2(\{\theta\})}{m_2(s)} \\ &= \frac{p m_1(s)}{p m_1(s) + (1-p) m_2(s)} \pi_1(A|s) + \frac{(1-p) m_2(s)}{p m_1(s) + (1-p) m_2(s)} \pi_2(A|s) \end{aligned}$$

Bayes factors - another method of hypothesis assessment.

Def. 7.2.1. Prob. model with sample space S and probability measure P , the odds in favor of event $A \subset S$ is $\frac{P(A)}{P(A^c)}$.

Bayes factors compare posterior odds with prior odds.

Def. 7.2.2. The Bayes factor BF_{H_0} in favor of

$H_0: \psi(\theta) = \psi_0$ is defined as

$$BF_{H_0} = \left\{ \frac{\pi(\psi(\theta) = \psi_0 | s)}{1 - \pi(\psi(\theta) = \psi_0 | s)} \right\} \bigg/ \left\{ \frac{\pi(\psi(\theta) = \psi_0)}{1 - \pi(\psi(\theta) = \psi_0)} \right\}$$

~~Small relationship~~ BF_{H_0} small provides evidence against H_0 .

Note also that $\pi(\psi(\theta) = \psi_0 | s) = \frac{\Gamma BF_{H_0}}{1 + \Gamma BF_{H_0}}$ where

$\Gamma = \frac{\pi(\psi(\theta) = \psi_0)}{1 - \pi(\psi(\theta) = \psi_0)}$ is the prior odds in favor of H_0 . So, when BF_{H_0} is small, then $\pi(\psi(\theta) = \psi_0 | s)$ is small and conversely.

The following result establishes connections with likelihood ratios.

Th. 7.2.1. If the prior π is a mixture $\pi = p\pi_1 + (1-p)\pi_2$,

where $\pi_1(A) = 1$, $\pi_2(A^c) = 1$, and $H_0: \theta \in A$,

$BF_{H_0} = \frac{m_1(s)}{m_2(s)}$ where m_i is the prior predictive of the data under π_i .

Proof. Recall that, a prior that concentrates all prob. on the set, results in the same posterior.

Then, $BF_{H_0} = \frac{\pi(A|s)}{1 - \pi(A|s)} \bigg/ \frac{\pi(A)}{1 - \pi(A)} = \frac{p m_1(s)}{(1-p) m_2(s)} \bigg/ \frac{p}{1-p} = \frac{m_1(s)}{m_2(s)}$

Note that the Bayes factor is independent of p .

Remark 1 Note that BF_{H_0} does not depend on p .

However, it is not immediately clear how to calibrate it.
 We can use e.g. $\Pi(\psi(\theta) = \psi_0 | S) = \frac{r BF_{H_0}(\psi_0)}{1 + r BF_{H_0}(\psi_0)}$ because

the posterior probability is directly interpretable. This, however, requires the knowledge of p .

Ex. 7.2.13. Location - normal

Continue: compute the prior predictor under Π_2 .

Recall that the joint density of (x_1, \dots, x_n) given μ is

$$m_2(x_1, \dots, x_n) = \int_{-\infty}^{\infty} \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} \right)^n \exp\left(-\frac{n-1}{2\sigma_0^2} s^2\right) \exp\left(-\frac{n}{2\sigma_0^2} (\bar{x} - \mu)^2\right) \times \left(\frac{1}{\sqrt{2\pi\tau_0^2}} \right) \exp\left(-\frac{1}{2\tau_0^2} (\mu - \mu_0)^2\right) d\mu$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} \right)^{n-1} \exp\left(-\frac{n-1}{2\sigma_0^2} s^2\right) \times \frac{1}{\sqrt{2\pi\tau_0^2}} \int_{-\infty}^{\infty} \exp\left(-\frac{n}{2\sigma_0^2} (\bar{x} - \mu)^2\right) \exp\left(-\frac{1}{2\tau_0^2} (\mu - \mu_0)^2\right) d\mu$$

The part after (x) can be shown to be

$$\frac{1}{\sqrt{2\pi\tau_0^2}} \exp\left(-\frac{1}{2} \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2} \right) \left(\frac{\mu_0}{\tau_0^2} + \frac{n\bar{x}}{\sigma_0^2} \right)^2\right) \times \exp\left(-\frac{1}{2} \left(\frac{\mu_0^2}{\tau_0^2} + \frac{n\bar{x}^2}{\sigma_0^2} \right) \left(\frac{1}{\tau_0^2} + \frac{1}{\sigma_0^2} \right)^{-1}\right)$$

Because Π_1 is degenerate at μ_0 we have the prior predictive

under Π_1 as $m_1(x_1, \dots, x_n) = \left(\frac{1}{\sqrt{2\pi\sigma_0^2}} \right)^n \exp\left(-\frac{n-1}{2\sigma_0^2} s^2\right) \exp\left(-\frac{n}{2\sigma_0^2} (\bar{x} - \mu_0)^2\right)$

Therefore, BF_{H_0} is equal to $\exp\left(-\frac{n}{2\sigma_0^2} (\bar{x} - \mu_0)^2\right)$ divided by $(**)$

For example, suppose that $\mu_0 = 0, \tau_0^2 = 2, \sigma_0^2 = 1, n = 10,$

$\bar{x} = 0.2$. Then $\exp\left(-\frac{n}{2\sigma_0^2} (\bar{x} - \mu_0)^2\right) = \exp\left(-\frac{10}{2} (0.2)^2\right) = 0.81873$

and $(**)$ is equal to 0.21615:

$\therefore BF_{H_0} \text{ is } = \frac{0.81873}{0.21615} = 3.7878$ - some evidence in favor of

$H_0: \mu = \mu_0$. If we take $p = \frac{1}{2}$ (complete indifference between H_0 being true and not being true), then $r = 1$ and we find that from $(***)$

$\Pi(\mu = \mu_0 | x_1, \dots, x_n) = \frac{3.7878}{1 + 3.7878} = 0.79114$ - strongly supportive

Ex Trying to guess whether a man is from the UK or

France. You offer him beer, brandy/cognac, whiskey, or wine and try to figure out who he is by the choice of the drink.

We have the following table:

State	beer	brandy	whisky	wine
France	10%	20%	10%	60%
UK	50%	10%	20%	20%

H_0 : France H_1 : UK

Assume some prior probabilities $\delta_0 = P(\text{France})$ and $\delta_1 = P(\text{UK})$.

- perhaps an ^{non-}informative prior $\delta_0 = \delta_1 = 0.5$

We can apply the simplest Bayes rule for this case and reject the null if the likelihood ratio $k_Y = \frac{P_{UK}(Y)}{P_{France}(Y)} > 1$

We decide that the man is from UK if beer or whiskey and France if wine or brandy. The values of λ are: $5 \frac{1}{2}$ $2 \frac{1}{3}$

7.2.4. Prediction.

We have an unobserved response value t in a sample space T and observed response $s \in S$. The statistical model is $\{P_\theta: \theta \in \Omega\}$ for s and the conditional statistical model $\{P_\theta(\cdot|s): \theta \in \Omega\}$ for t given s . Assume that both models have the same true value of $\theta \in \Omega$. The objective is to construct a prediction $\hat{t}(s) \in T$ of the unobserved t based on s . The value of t is unknown because it is e.g. a future outcome.

Use $q_\theta(\cdot|s)$ to obtain the joint distribution of (θ, s, t) : $q_\theta(\cdot|s) f_\theta(s) \pi(\theta)$.

Once we have observed s , the conditional density of (t, θ) given s is

$$\frac{q_\theta(t|s) f_\theta(s) \pi(\theta)}{\int_{\Omega} \int q_\theta(t|s) f_\theta(s) \pi(\theta) dt d\theta} = \frac{q_\theta(t|s) f_\theta(s) \pi(\theta)}{\int_{\Omega} f_\theta(s) \pi(\theta) d\theta} = \frac{q_\theta(t|s) f_\theta(s) \pi(\theta)}{m(s)}$$

The marginal posterior distribution of t , known as the posterior predictive of t , is

$$q(t|s) = \int_{\Omega} \frac{q_0(t|s) f_0(s) d\theta}{m(s)} d\theta = \int_{\Omega} q_0(t|s) \tilde{q}_0(s) d\theta.$$

Remark. The posterior predictive of t is obtained by averaging $q_0(t|s)$ w.r.t. respect to the posterior distribution of θ .

Ex. 7-2.14 Given $(x_1, \dots, x_n) \sim \text{Ber}(\theta)$, $\theta \sim \text{Beta}(\alpha, \beta)$, predict X_{n+1} . The posterior probability of n of X_{n+1} is

$$\begin{aligned}
 q(t|x_1, \dots, x_n) &= \int_0^1 \theta^t (1-\theta)^{n-t} \frac{\Gamma(n+\alpha+\beta)}{\Gamma(n\bar{x}+\alpha)\Gamma(n(1-\bar{x})+\beta)} \theta^{n\bar{x}+\alpha-1} (1-\theta)^{n(1-\bar{x})+\beta-1} d\theta \\
 &= \frac{\Gamma(n+\alpha+\beta)}{\Gamma(n\bar{x}+\alpha)\Gamma(n(1-\bar{x})+\beta)} \int_0^1 \theta^{n\bar{x}+\alpha+t-1} (1-\theta)^{n(1-\bar{x})+\beta+(1-t)-1} d\theta \\
 &= \frac{\Gamma(n+\alpha+\beta)}{\Gamma(n\bar{x}+\alpha)\Gamma(n(1-\bar{x})+\beta)} \frac{\Gamma(n\bar{x}+\alpha+t)\Gamma(n(1-\bar{x})+\beta+1-t)}{\Gamma(n+\alpha+\beta+t)} \\
 &= \begin{cases} \frac{n\bar{x}+\alpha}{n+\alpha+\beta}, & t=1 \\ \frac{n(1-\bar{x})+\beta}{n+\alpha+\beta}, & t=0 \end{cases} \sim \text{Ber}\left(\frac{n\bar{x}+\alpha}{n+\alpha+\beta}\right)
 \end{aligned}$$

(Keep in mind that X_{n+1} is independent of x_1, \dots, x_n)

Using the posterior mode as the predictor we maximize $q(t|x_1, \dots, x_n)$ w.r.t. t ; $\Rightarrow \hat{t} = \begin{cases} 1, & \text{if } \frac{n\bar{x}+\alpha}{n+\alpha+\beta} \geq \frac{n(1-\bar{x})+\beta}{n+\alpha+\beta} \\ 0 & \text{otherwise} \end{cases}$

At the same time, $E[t|x_1, \dots, x_n] = \frac{n\bar{x}+\alpha}{n+\alpha+\beta}$ - can be anywhere in $[0, 1]$