

§2 Bayesian Inference

Outline

1. Bayesian inference for the normal distribution
2. Summarisation of posterior distributions

1. Bayesian inference for the normal distribution

A) Unknown mean, known precision

Suppose we have a sample y of n independent observations

$$Y_i | \theta \sim \text{Normal}(\theta, \sigma^2) , \quad i = 1, \dots, n,$$

where σ^2 is known and θ is unknown; ie

$$\begin{aligned} p(y_i | \theta) &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2\sigma^2} (y_i - \theta)^2 \right] \\ &\propto \exp \left[-\frac{1}{2} \left\{ \frac{1}{\sigma^2} (y_i - \theta)^2 \right\} \right] . \end{aligned}$$

Sometimes, it is more convenient to work with the *precision* parameter $\tau = \frac{1}{\sigma^2}$. So,

$$Y_i | \theta \sim \text{Normal}(\theta, \tau^{-1}) , \quad i = 1, \dots, n,$$

and thus

$$p(y_i | \theta) \propto \exp \left[-\frac{1}{2} \left\{ \tau (y_i - \theta)^2 \right\} \right] .$$

Classical inference

Maximum likelihood estimation: $\hat{\theta} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\text{Var}(\hat{\theta}) = \sigma^2/n = (n\tau)^{-1}$.

Bayesian inference

Suppose we have a normal prior distribution for the unknown parameter θ :

$$\theta \sim \text{Normal}(\mu_0, \phi_0^{-1}),$$

where μ_0 (prior mean) and ϕ_0 (prior precision) are fixed. Hence,

$$p(\theta) \propto \exp \left[-\frac{1}{2} \{ \phi_0 (\theta - \mu_0)^2 \} \right] .$$

Then the posterior distribution for θ is

$$\begin{aligned} p(\theta | \mathbf{y}) &\propto p(\theta) p(\mathbf{y} | \theta) = p(\theta) \prod_{i=1}^n p(y_i | \theta) \\ &\propto \exp \left[-\frac{1}{2} \left\{ \phi_0 (\theta - \mu_0)^2 + \tau \sum_{i=1}^n (y_i - \theta)^2 \right\} \right] \end{aligned}$$

$$p(\theta|\mathbf{y}) \propto \exp \left[-\frac{1}{2} \left\{ \phi_0(\theta - \mu_0)^2 + \tau \sum_{i=1}^n (y_i - \theta)^2 \right\} \right]$$

The part inside $\{\dots\}$

$$\begin{aligned} &= \theta^2 (\phi_0 + n\tau) - 2\theta (\mu_0\phi_0 + n\bar{y}\tau) + \text{const}_1 \\ &= (\phi_0 + n\tau) \left(\theta - \frac{\mu_0\phi_0 + n\bar{y}\tau}{\phi_0 + n\tau} \right)^2 + \text{const}_2 \end{aligned}$$

Therefore,

$$p(\theta|\mathbf{y}) \propto \exp \left[-\frac{1}{2} \left\{ (\phi_0 + n\tau) \left(\theta - \frac{\mu_0\phi_0 + n\bar{y}\tau}{\phi_0 + n\tau} \right)^2 \right\} \right] .$$

Hence $\theta | \mathbf{y} \sim \text{Normal}(\mu_1, \phi_1^{-1})$,

where $\phi_1 = \phi_0 + n\tau$,

$$\mu_1 = \frac{\phi_0\mu_0 + n\tau\bar{y}}{\phi_0 + n\tau} = w\mu_0 + (1 - w)\bar{y} ,$$

$$w = \frac{\phi_0}{\phi_0 + n\tau}, \quad 1 - w = \frac{n\tau}{\phi_0 + n\tau} .$$

Three interesting notes:

1. The posterior precision ϕ_1 is the sum of the prior precision ϕ_0 and the sample precision $n\tau$.
2. The posterior mean μ_1 is a weighted average of the prior mean μ_0 and the sample mean \bar{y} , weighted by their relative precisions.
3. Both θ and $\theta|\mathbf{y}$ follow normal distributions.

$$\begin{aligned}\phi_1 &= \phi_0 + n\tau, \quad \mu_1 = w\mu_0 + (1-w)\bar{y}, \\ w &= \frac{\phi_0}{\phi_0 + n\tau}.\end{aligned}$$

Three important comments

1. If $n \rightarrow \infty$ with ϕ_0 fixed, then $w \rightarrow 0$ and $1 - w \rightarrow 1$. So, for large enough n ,

$$\theta \mid \mathbf{y} \sim \text{Normal}(\bar{y}, (n\tau)^{-1})$$

approximately, ie posterior does not depend on the prior.

2. If $\phi_0 \rightarrow 0$ with n fixed, then $w \rightarrow 0$ and $1 - w \rightarrow 1$. So, for very diffused prior beliefs (ie non-informative priors),

$$\theta \mid \mathbf{y} \sim \text{Normal}(\bar{y}, (n\tau)^{-1})$$

approximately, ie posterior does not depend on the prior.

3. If we write $\phi_0 = \kappa_0\tau$, so that prior is $\text{Normal}(\mu_0, (\kappa_0\tau)^{-1})$, then

$$\theta \mid \mathbf{y} \sim \text{Normal}\left(\frac{n}{n + \kappa_0}\bar{y} + \frac{\kappa_0}{n + \kappa_0}\mu_0, [(n + \kappa_0)\tau]^{-1}\right).$$

Hence κ_0 may be viewed as a ‘prior sample size’.

Example 2.1

The first reliable geochemical datings for the age of the Ennerdale granophyre rock strata were obtained in the 1960s using the K/Ar method (based on the relative proportions of potassium 40 and argon 40 in the rock). **The resulting estimate was 370 (SE 20) million years.** In the 1970s a more accurate method based on the relative proportions of rubidium 87 and strontium 87 became available. **We shall assume that for the latter method the standard deviation of a measurement is known to be 8 million years.**

Rb/Sr measurements were made on **a sample of 5 observations** taken from the rocks and the estimated age calculated from each measurement. **The mean of these estimated ages was 421 million years.**

Using the information from the K/Ar method as a prior, and assuming normally distributed measurement errors, obtain a posterior distribution for θ , the true age of the Ennerdale granophyre.

Prior: $\theta \sim \text{Normal}(370, 20^2)$

Likelihood: $Y_i \mid \theta \sim \text{Normal}(\theta, 8^2)$, $i = 1, \dots, 5$, with $\bar{y} = 421$

\Rightarrow Posterior: $\theta \mid \mathbf{y} \sim \text{Normal}(\mu_1, \phi_1^{-1})$, where

$$\phi_1 = \frac{1}{20^2} + \frac{5}{8^2} = 0.081 ,$$

$$w = \frac{1/20^2}{1/20^2 + 5/8^2} = 0.031 ,$$

$$\mu_1 = 0.031 \times 370 + (1 - 0.031) \times 421 = 419.42 .$$

B) Known mean, unknown precision

Suppose we have an independent sample y :

$$Y_i \mid \tau \sim \text{Normal}(\theta, \tau^{-1}) , \quad i = 1, \dots, n,$$

where this time τ is unknown and θ is known;
ie

$$p(y_i \mid \tau) \propto \tau^{1/2} \exp \left[-\frac{1}{2} \{ \tau (y_i - \theta)^2 \} \right] .$$

For reasons to be discussed later ($\tau > 0$; a conjugate prior for this likelihood) it is convenient to choose a gamma distribution as a prior for τ , i.e. $\tau \sim \text{Gamma}(\alpha, \beta)$:

$$p(\tau) = \frac{\beta^\alpha}{\Gamma(\alpha)} \tau^{\alpha-1} e^{-\beta\tau} \quad (\tau > 0)$$

where α (the shape parameter, > 0) and β (the inverse scale parameter, > 0) are fixed.

Note that if $\tau \sim \text{Gamma}(\alpha, \beta)$ (precision)
then $\tau^{-1} \sim \text{InvGamma}(\alpha, \beta)$ (variance).

The posterior distribution for τ is

$$\begin{aligned}
p(\tau \mid \mathbf{y}) &\propto p(\tau) p(\mathbf{y} \mid \tau) = p(\tau) \prod_{i=1}^n p(y_i \mid \tau) \\
&\propto p(\tau) \prod_{i=1}^n \left\{ \tau^{1/2} \exp \left[-\frac{\tau}{2} (y_i - \theta)^2 \right] \right\} \\
&\propto \tau^{\alpha-1} \exp[-\beta\tau] \times \tau^{n/2} \exp \left[-\frac{\tau}{2} n s_{(n)}^2 \right] \\
&= \tau^{\frac{n}{2} + \alpha - 1} \exp \left[- \left(\frac{n s_{(n)}^2}{2} + \beta \right) \tau \right] \\
\tau \mid \mathbf{y} &\sim \text{Gamma} \left(\frac{n}{2} + \alpha, \frac{n s_{(n)}^2}{2} + \beta \right)
\end{aligned}$$

where $s_{(n)}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2$

- Both τ and $\tau \mid \mathbf{y}$ follow Gamma distributions.

$$\tau \mid \mathbf{y} \sim \text{Gamma} \left(\frac{n}{2} + \alpha, \frac{ns_{(n)}^2}{2} + \beta \right)$$

Comments

1. shape and inverse scale parameters of posterior distribution are sums of the prior parameters and statistics of the data.
2. posterior shape is $\frac{n}{2} + \alpha = \frac{n+2\alpha}{2}$;
posterior inverse scale is $\frac{ns_{(n)}^2}{2} + \beta = \frac{ns_{(n)}^2 + 2\beta}{2}$.

So, prior $\text{Gamma}(\alpha, \beta)$ can be thought of as providing information equivalent to 2α observations with total sum of squares (ie $\sum_{i=1}^{2\alpha} (y_i - \theta)^2$) equal to 2β .

3. Mean of the posterior distribution

$$E[\tau \mid \mathbf{y}] = \frac{n/2 + \alpha}{ns_{(n)}^2/2 + \beta}$$

- If $s_{(n)}^2$ is large (for fixed n), $E[\tau \mid \mathbf{y}]$ is small; if $s_{(n)}^2$ is small, $E[\tau \mid \mathbf{y}]$ is large. Make sense?
- When $\alpha \rightarrow 0$ and $\beta \rightarrow 0$, $E[\tau \mid \mathbf{y}] \rightarrow 1/s_{(n)}^2$. Make sense?

C) Unknown mean, unknown precision

Suppose we have an independent sample y :

$$Y_i | \theta, \tau \sim \text{Normal}(\theta, \tau^{-1}) , \quad i = 1, \dots, n,$$

where both τ and θ are unknown.

We now need to specify a *joint* prior distribution $p(\theta, \tau)$. One way to do this is as follows:

$$\begin{aligned}\tau &\sim \text{Gamma}(\alpha, \beta) \\ \theta | \tau &\sim \text{Normal}(\mu_0, (\kappa_0 \tau)^{-1}) \\ p(\theta, \tau) &= p(\theta | \tau) p(\tau)\end{aligned}$$

However, for now we shall simplify the maths by instead assuming a *non-informative* prior for θ and τ :

$$p(\theta, \tau) \propto \tau^{-1}$$

(This is the product of a uniform prior for θ and the Jeffreys' prior for τ . See Lee p84. More on non-informative priors later.)

Joint posterior: $p(\theta, \tau \mid \mathbf{y})$

Multiplying this prior by the normal likelihood $\prod_{i=1}^n p(y_i \mid \theta, \tau)$ gives the joint posterior

$$\begin{aligned} p(\theta, \tau \mid \mathbf{y}) &\propto \tau^{-1} \tau^{n/2} \exp \left[-\frac{\tau}{2} \sum_{i=1}^n (y_i - \theta)^2 \right] \\ &= \tau^{n/2-1} \exp \left[-\frac{\tau}{2} \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - \theta)^2 \right] \\ &= \tau^{n/2-1} \times \\ &\quad \exp \left[-\frac{\tau}{2} \left\{ \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \theta)^2 \right\} \right] \\ &= \tau^{n/2-1} \times \\ &\quad \exp \left[-\frac{\tau}{2} \left\{ (n-1)s^2 + n(\bar{y} - \theta)^2 \right\} \right] \end{aligned}$$

where $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$.

Marginal posterior for τ : $p(\tau \mid \mathbf{y})$

$$\begin{aligned}
 p(\tau \mid \mathbf{y}) &= \int_{-\infty}^{\infty} p(\theta, \tau \mid \mathbf{y}) d\theta \\
 &\propto \int_{-\infty}^{\infty} \tau^{\frac{n}{2}-1} \exp \left[-\frac{\tau}{2} ((n-1)s^2 + n(\bar{y} - \theta)^2) \right] d\theta \\
 &= \tau^{\frac{n}{2}-1} \exp \left[-\frac{\tau}{2} (n-1)s^2 \right] \int_{-\infty}^{\infty} \exp \left[-\frac{\tau n}{2} (\bar{y} - \theta)^2 \right] d\theta \\
 &\quad \text{Noting that } \int_{-\infty}^{\infty} \exp \left[-\frac{\tau n}{2} (\theta - \bar{y})^2 \right] d\theta \text{ is an} \\
 &\quad \text{un-normalised normal integral gives:} \\
 &= \tau^{\frac{n}{2}-1} \exp \left[-\frac{\tau}{2} (n-1)s^2 \right] \sqrt{\frac{2\pi}{n\tau}} \\
 &\propto \tau^{\frac{n-1}{2}-1} \exp \left[-\frac{(n-1)s^2}{2} \tau \right] .
 \end{aligned}$$

Therefore,

$$\tau \mid y \sim \text{Gamma} \left(\frac{n-1}{2}, \frac{(n-1)s^2}{2} \right) .$$

Note: $E[\tau \mid \mathbf{y}] = \frac{n-1}{2} / \frac{(n-1)s^2}{2} = \frac{1}{s^2}$, where s^2 is the usual unbiased classical estimator of the variance.

Marginal posterior for θ : $p(\theta \mid \mathbf{y})$

$$\begin{aligned} p(\theta \mid \mathbf{y}) &= \int_0^\infty p(\theta, \tau \mid \mathbf{y}) d\tau \\ &\propto \Gamma\left(\frac{n}{2}\right) \left[\frac{(n-1)s^2 + n(\bar{y} - \theta)^2}{2} \right]^{-n/2} \\ &\propto \left[1 + \frac{n(\theta - \bar{y})^2}{(n-1)s^2} \right]^{-n/2} . \end{aligned}$$

Using transformation of variable: $\phi = \frac{\theta - \bar{y}}{\sqrt{s^2/n}}$,
we get

$$\frac{\theta - \bar{y}}{\sqrt{s^2/n}} \mid \mathbf{y} \sim t_{n-1} .$$

Conditional posterior for $\theta \mid \tau, \mathbf{y}$

$$\theta \mid \tau, \mathbf{y} \sim \text{Normal} \left(\bar{y}, (n\tau)^{-1} \right) .$$

2. Summarisation of posterior distributions

Bayesian inference is based on the posterior distribution $p(\theta \mid y)$. The posterior encapsulates everything that is known about θ following observation of the data y .

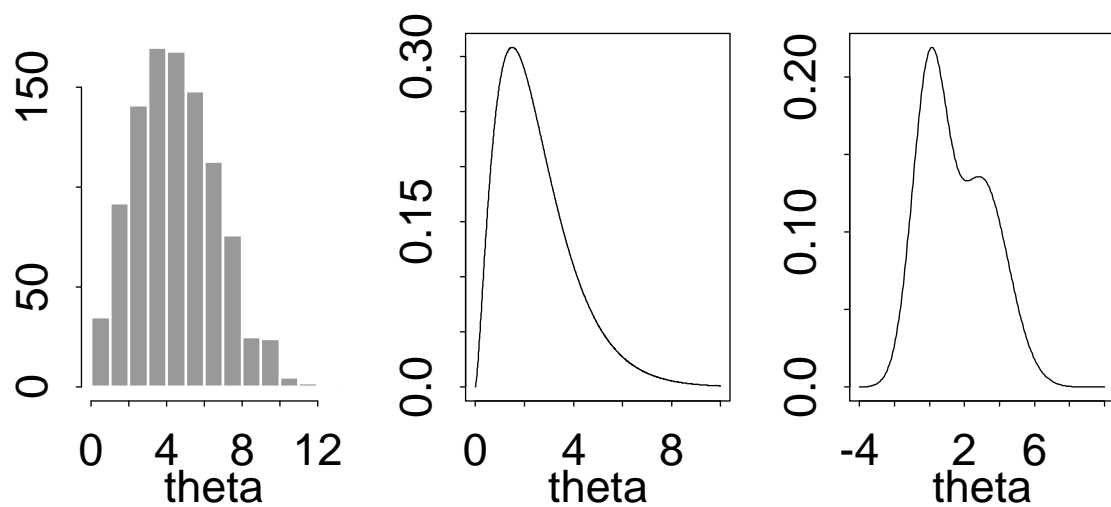
To help understand the mathematical formula for $p(\theta \mid y)$ and to identify its interesting features in a clear and concise way, we could *summarise* the posterior.

Three main ways to summarise $p(\theta \mid y)$:

1. Graphical summaries - e.g. plot the shape of the posterior
2. Quantitative summaries - e.g. measures of location and dispersion
3. Summaries relating to specific hypotheses - e.g. $P(\theta \in \Theta)$

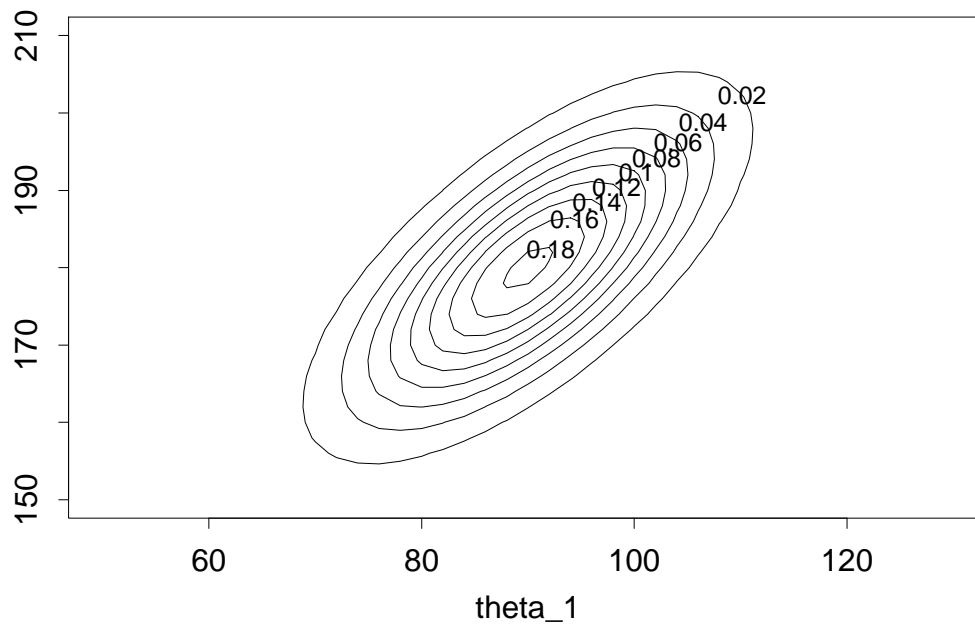
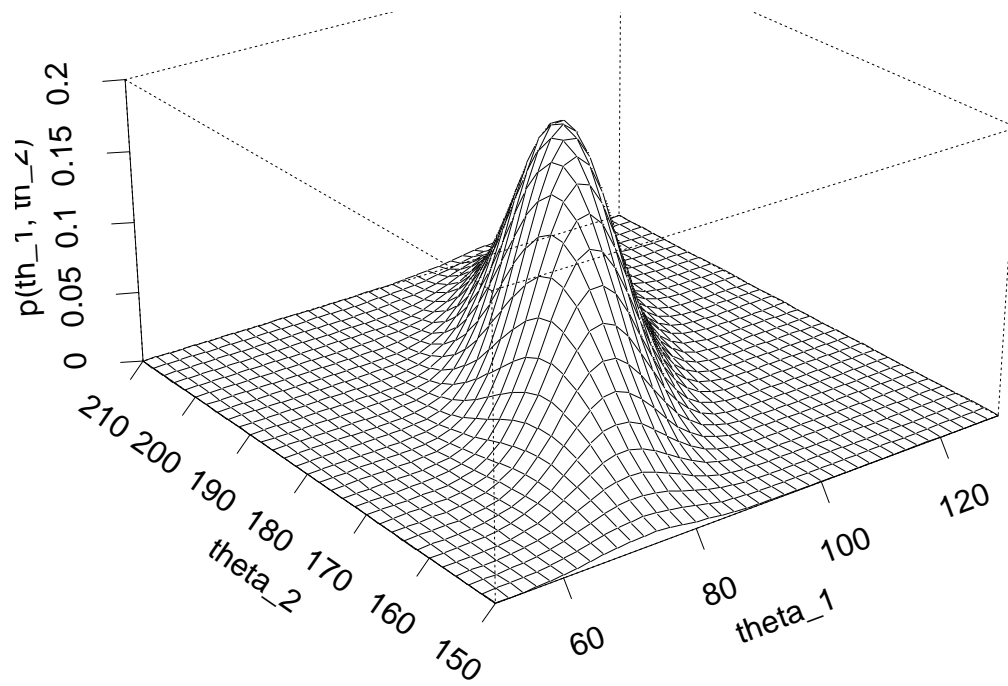
A) Graphical summaries

Univariate posteriors, $p(\theta \mid y)$



- Summarises shape of posterior
- Shows range of θ values that are with highest probabilities

Bivariate posteriors, $p(\theta_1, \theta_2 \mid y)$



Multivariate posteriors, $p(\boldsymbol{\theta} \mid y)$

- Most real-life problems involve many parameters $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$
- Difficult to visualise more than 2 dimensions graphically
- We can reduce the dimension by integrating to obtain one-dimensional marginal posteriors, $p(\theta_i \mid y) = \int p(\boldsymbol{\theta} \mid y) d\boldsymbol{\theta}_{-i}$, and conditional posteriors, $p(\theta_i \mid \boldsymbol{\theta}_{-i}, y)$, for different values of $\boldsymbol{\theta}_{-i}$ (see e.g. pp12-13). However, in many real-life problems these integrals are analytically intractable. \Rightarrow need simulation methods (e.g. MCMC) to evaluate them.

B) Quantitative summaries

Quantitative summaries of Bayesian posteriors usually take the form of *point* and *interval estimates* for the parameters of interest.

- In principle, a Bayesian is free to choose *any* (sensible) numerical summary of the ‘location’ of θ (ie that provides an indication of the ‘typical’ value of θ) as a point estimate.

Common choices include the *mean*, *median* and *mode* of the posterior.

- Likewise, any (sensible) numerical summary of the ‘dispersion’ of the posterior may be used.

Common choices include the standard deviation of the posterior distribution, quantile ranges/credible intervals, and highest posterior density (HPD) regions.

Formal ‘rules’ for choosing optimal point or interval estimators can be obtained by appealing to *Statistical Decision Theory*.

Bayesian point estimates

- *Posterior mean*: $E(\theta) = \int \theta p(\theta | y) d\theta$
- *Posterior mode*: value θ_{mode} of θ for which $p(\theta | y)$ equals its maximum (there may be more than one modes)
- *Posterior median*: value θ_{med} of θ for which $P(\theta \geq \theta_{\text{med}} | y) = P(\theta \leq \theta_{\text{med}} | y) = 0.5$

Recall Example 2.1: $\theta | \mathbf{y} \sim \text{Normal}(419.42, 12.4)$.
A point estimate for the age of the rocks is

$$E(\theta | \mathbf{y}) = \theta_{\text{mode}} = \theta_{\text{med}} = 419.42 \text{ million years}$$

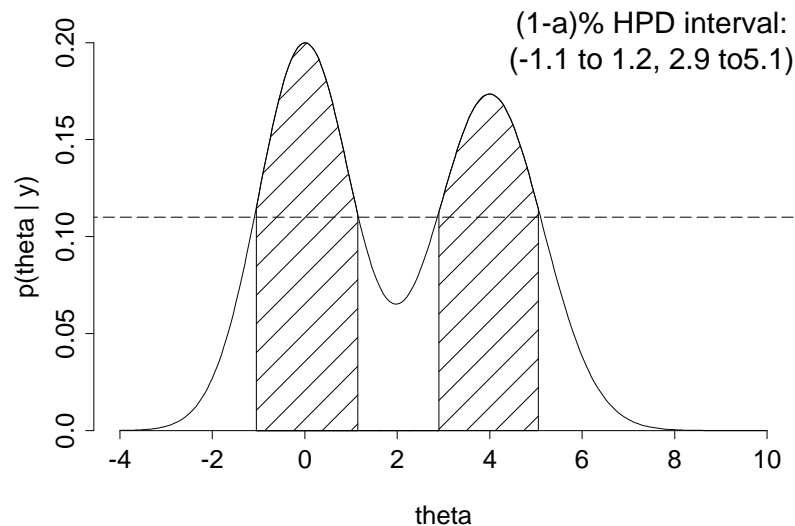
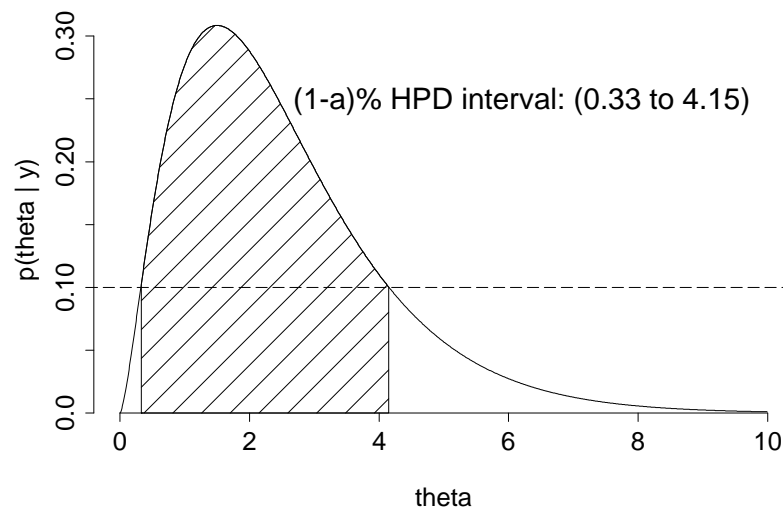
Bayesian interval estimates

- *Credible interval*: An interval $[c_1, c_2]$ is a $100(1-\alpha)\%$ credible interval for θ if

$$P(\theta \in [c_1, c_2] \mid y) = 1 - \alpha .$$

- E.g.: Let θ_q be the $(100 \times q)$ -th percentile of $p(\theta \mid y)$, ie $P(\theta \leq \theta_q \mid y) = q$. Then $[\theta_{.025}, \theta_{.975}]$ is a 95% (central) credible interval for θ .
 - If $p(\theta \mid y)$ is (approximately) $\text{Normal}(\mu_\theta, s^2)$, then the interval $[\mu_\theta - z_{\frac{\alpha}{2}}s, \mu_\theta + z_{\frac{\alpha}{2}}s]$ is a $100(1-\alpha)\%$ credible interval for θ , where $z_{\frac{\alpha}{2}}$ is the value such that $P(Z \leq -z_{\frac{\alpha}{2}}) = P(Z \geq z_{\frac{\alpha}{2}}) = \frac{\alpha}{2}$ for a standard Normal random variable Z .
- *Credible region*: For multivariate θ , a region \mathbf{C} is a $100(1-\alpha)\%$ (simultaneous) credible region for θ if $P(\theta \in \mathbf{C} \mid y) = 1 - \alpha$.

- *Highest posterior density (HPD) interval*:
An interval such that the posterior density at any point inside the interval is greater than that at any point outside.



- A $100(1-\alpha)\%$ HPD interval is the $100(1-\alpha)\%$ credible interval with the shortest width.
- Calculation of HPD intervals usually requires the aid of a computer.

C) Summaries for specific hypotheses

One of the most powerful features of Bayesian inference is that it can answer *any* question of interest by appropriate summarisation of the posterior.

E.g. recall Example 2.1. Suppose the hypotheses are $H_0 : \theta < \theta_0$ and $H_1 : \theta \geq \theta_0$ with equal losses for type I and type II errors. Then we reject H_0 if

$$P(\theta < \theta_0 \mid \mathbf{y}) < P(\theta > \theta_0 \mid \mathbf{y})$$

$$\begin{aligned} P(\theta < 420 \mid \mathbf{y}) &= \int_{-\infty}^{420} p(\theta \mid \mathbf{y}) d\theta \\ &= \Phi\left(\frac{420 - 419.42}{\sqrt{12.4}}\right) \\ &= 0.57 \end{aligned}$$

Suppose the losses for type I and II errors are c_α and c_β , respectively. Then we reject H_0 if

$$c_\alpha P(\theta < \theta_0 \mid \mathbf{y}) < c_\beta P(\theta > \theta_0 \mid \mathbf{y}) .$$

Comments

Consider a classical hypothesis test of

$$H_0 : \theta < \theta_0 \text{ vs } H_1 : \theta \geq \theta_0,$$

and suppose that we got a p -value of 0.02.

- Does this mean that $P(\theta < \theta_0) = 0.02$?
- This means that: if H_0 is true (ie $\theta < \theta_0$), then the probability of observing data at least as extreme as the data actually observed is at most 0.02.

Classical inference cannot provide probability statements about parameters!

Brief summary of Bayesian inference:

Bayesian inference involves calculating posterior distributions and summarising them using appropriate graphical and numerical summaries.

*A Bayesian and a Frequentist were to be executed. The judge asked them what were their last wishes. The Bayesian replied that he would like to give the Frequentist one more lecture. The judge granted the Bayesian's wish and then turned to the Frequentist for his last wish. The Frequentist quickly responded that he wished to hear the lecture again and again and again and again
— X.-L. Meng*

Outline revisited

1. Bayesian inference for the normal distribution
2. Summarisation of posterior distributions

Next week: Prior Distributions