

Decision Theory and Bayesian Analysis

Dr. Vilda Purutcuoglu

1

¹Edited by Anil A. Aksu based on lecture notes of STAT 565 course by Dr. Vilda Purutcuoglu

Contents

Decision Theory and Bayesian Analysis	1
Lecture 1. Bayesian Paradigm	5
1.1. Bayes theorem for distributions	5
1.2. How Bayesian Statistics Uses Bayes Theorem	6
1.3. Prior to Posterior	8
1.4. Triplot	8
Lecture 2. Some Common Probability Distributions	13
2.1. Posterior	15
2.2. Weak Prior	17
2.3. Sequential Updating	19
2.4. Normal Sample	19
2.5. NIC distributions	19
2.6. Posterior	20
2.7. Weak prior	21
Lecture 3. Inference	23
3.1. Shape	24
3.2. Visualizing multivariate densities	26
3.3. Informal Inferences	27
3.4. Multivariate inference	28
Lecture 4. Formal Inference	29
4.1. Utility and decisions	29
4.2. Formal Hypothesis Testing	30
4.3. Nuisance Parameter	31
4.4. Transformation	32
4.5. The prior distribution	32
4.6. Subjectivity	32
4.7. Noninformative Priors	33
4.8. Informative Priors	34
4.9. Prior Choices	34
Lecture 5. Structuring Prior Information	39

4	DR. VILDA PURUTCUOGLU, BAYESIAN ANALYSIS	
5.1.	Binary Exchangeability	40
5.2.	Exchangeable Parameters	40
5.3.	Hierarchical Models	41
Lecture 6.	Sufficiency and Ancillary	43
6.1.	The Likelihood Principle	43
6.2.	Identifiability	44
6.3.	Asymptotic Theory	44
6.4.	Preposterior Properties	45
6.5.	Conjugate Prior Forms	45
Lecture 7.	Tackling Real Problems	47
7.1.	What is a Markov Chain?	47
7.2.	The Chapman-Kolmogorov Equations	49
7.3.	Marginal Distributions	49
7.4.	General Properties of Markov Chains	49
7.5.	Noniterative Monte Carlo Methods	51
Lecture 8.		55
8.1.	The Gibbs Sampler	55
Lecture 9.	Summary of the properties of Gibbs sampler	61
9.1.	Metropolis Algorithm	61
9.2.	Metropolis Algorithm	62
9.3.	Data Augmentation	64
GIBBS SAMPLING		67
DATA AUGMENTATION		71
R Codes		75
Bibliography		79

LECTURE 1

Bayesian Paradigm

1.1. Bayes theorem for distributions

If A and B are two events,

$$(1.1) \quad P(A | B) = \frac{P(A)P(B | A)}{P(B)}.$$

This is just a direct consequence of the multiplication law of probabilities that says we can express $P(A | B)$ as either $P(A)P(B | A)$ or $P(B)P(A | B)$. For discrete distributions, if Z, Y are discrete random variables

$$(1.2) \quad P(Z = z | Y = y) = \frac{P(Z = z)P(Y = y | Z = z)}{P(Y = y)}.$$

- How many distributions do we deal with here?

We can express the denominator in terms of the distribution in the numerator[1].

$$(1.3) \quad P(Y = y) = \sum_z P(Y = y, Z = z) = \sum_z P(Z = z)P(Y = y | Z = z).$$

- This is sometimes called the law of total probability

In this context, it is just an expression of the fact that as z ranges over the possible values of Z , the probabilities on the left hand-side of equation 1.2 make up the distribution of Z given $Y = y$, and so they must add up to one. The extension to continuous distribution is easy. If Z, Y are continuous random variable,

$$(1.4) \quad f(Z | Y) = \frac{f(Z)f(Y | Z)}{f(Y)}.$$

where the denominator is now expressed as an integral:

$$(1.5) \quad f(Y) = \int f(Z)f(Y | Z)dZ.$$

$$(1.6) \quad f = \begin{cases} \textit{continous name?} \\ \textit{discrete name?} \end{cases}$$

1.2. How Bayesian Statistics Uses Bayes Theorem

Theorem 1.7 (Bayes' theorem).

$$P(A | B) = \frac{P(A)P(B | A)}{P(B)}$$

$P(B)$ = if we are interested in the event B , $P(B)$ is the initial or prior probability of the occurrence of event B . Then we observe event A

$P(B | A)$ = How likely B is when A is known to have occurred is the posterior probability $P(B | A)$.

Bayes' theorem can be understood as a formula for updating from prior to posterior probability, the updating consists of multiplying by the ratio $P(B | A)/P(A)$. It describes how a probability changes as we learn new information. Observing the occurrence of A will increase the probability of B if $P(B | A) > P(A)$. From the law of total probability,

$$(1.8) \quad P(A) = P(A | B)P(B) + P(A | B^c)P(B^c).$$

where $P(B^c) = 1 - P(B)$.

Lemma 1.9.

$$P(A | B) - P(A) = \frac{P(A) - P(A | B^c)P(B^c)}{1 - P(B^c)} - P(A)$$

Proof.

$$P(A | B) - P(A) = \frac{P(A) - P(A | B^c)P(B^c) - P(A) + P(A)P(B^c)}{P(B)}$$

$$P(A | B) - P(A) = \frac{P(B^c)(P(A) - P(A | B^c))}{P(B)}$$

$$P(A | B) - P(A) = P(B^c) \left(\frac{P(B)P(A | B) + P(B^c)P(A | B^c)}{P(B)} - \frac{P(A | B^c)}{P(B)} \right)$$

$$P(A | B) - P(A) = P(B^c) \left(P(A | B) - \frac{P(A | B^c)(1 - P(B^c))}{P(B)} \right)$$

$$P(A | B) - P(A) = P(B^c)(P(A | B) - P(A | B^c))$$

□

1.2.1. Generalization of the Bayes' Theorem

Let B_1, \dots, B_n be a set of mutually exclusive events. Then

$$(1.10) \quad P(B_r | A) = \frac{P(B_r)P(A | B_r)}{P(A)} = \frac{P(B_r)P(A | B_r)}{\sum_{i=1}^n P(B_r)P(A | B_r)}.$$

- Assuming that $P(B_r) > 0, P(A | B) > P(A)$ if and only if $P(A | B) > P(A | B^c)$.

- In Bayesian inference we use Bayes' theorem in a particular way.
- Z is the parameter (vector) θ .
- Y is the data (vector) X .

So we have

$$(1.11) \quad f(\theta | X) = \frac{f(\theta)f(X | \theta)}{f(X)}$$

$$(1.12) \quad f(X) = \int f(\theta)f(X | \theta)d\theta.$$

$$(1.13) \quad f(\theta) = \textit{prior}.$$

$$(1.14) \quad f(\theta | X) = \textit{posterior}.$$

$$(1.15) \quad f(X | \theta) = \textit{likelihood}.$$

1.2.2. Interpreting our sense

How do we interpret the things we see, hear, feel, taste or smell?

Example 1.2.1. I hear a song on the radio I identify the singer as Robbie Williams. Why do I think it's Robbie Williams?. Because he sounds like that. Formally, $P(\text{What I hear Robbie Williams}) \gg P(\text{What I hear someone else})$

Example 1.2.2. I look out of the window and see what appears to be a tree. It has a big, dark coloured part sticking up out of the ground that branches into thinner sticks and on the ends of these are small green things. Clearly, $P(\textit{view} | \textit{tree})$ is high and $P(\textit{view} | \textit{car})$ or $P(\textit{view} | \textit{Robbie Williams})$ are very small. But $P(\textit{view} | \textit{cardboard cutout cunningly painted to look like a tree})$ is also very high. Maybe even higher than $P(\textit{view} | \textit{tree})$ in the sense that what I see looks almost like a tree.

Does this mean I should now believe that I am seeing a cardboard cut-out cunningly painted to look like a tree? No because it is much less likely to begin with than a red tree.

In statistical terms, consider some data X and some unknown parameter θ . The first step in any statistical analysis is to build a model that links the data to unknown parameters and the main function of this model is to allow us to state the probability of observing any data given any specified values of the parameters. That is the model defines $f(x | \theta)$.

When we think of $f(x | \theta)$ as a function of θ for fixed observed data X , we call it likelihood function and it by $L(\theta, X)$.

- So how can we combine this with our example?

This perspective underlies the differences between the two main theories of statistical inference.

- Frequentist inference essentially uses only the likelihood, it does not recognize $f(\theta)$.
- Bayesian inference uses both likelihood and $f(\theta)$.

The principal distinguishing feature of Bayesian inference as opposed to frequentist inference is its use of $f(\theta)$.

1.3. Prior to Posterior

We refer to $f(\theta)$ as the prior distribution of θ . It represents knowledge about θ prior to observing the data X . We refer to $f(\theta | X)$ as the posterior distribution of θ and it represents knowledge about θ after observing X .

- So we have two sources of information about θ .
- Here $f(x)$ does not depend on θ . Thus $\int f(\theta | x) d\theta = 1$. Since $f(x)$ is a constant within the integral, we can take it outside to get $1 = f^{-1}(x) \int f(\theta) f(x | \theta) d\theta$.
- $f(\theta | x) \propto f(\theta) f(x | \theta) \propto f(\theta) L(\theta; x)$ (the posterior is proportional to the prior times the likelihood).
- The constant that we require to scale the right hand side to integrate to 1 is usually called the normalizing constant. If we haven't dropped any constants from $f(\theta)$ or $f(x | \theta)$, then the normalising constant is just $f^{-1}(x)$, otherwise it also restores any dropped constants.

1.4. Triplot

If for any value of θ , we have either $f(\theta) = 0$ or $f(x | \theta) = 0$, then we will also have $f(\theta | x) = 0$. This is called the property of zero preservation. So if either:

- the prior information says that this θ value is impossible
- the data say that this value of θ is impossible because if it were the true value, then the observed data would have been impossible, then the posterior distribution confirms that this value of θ is impossible.

Definition 1.16. Crowell's Rule: If either information source completely rules out a specific θ , then the posterior must rule it out too.

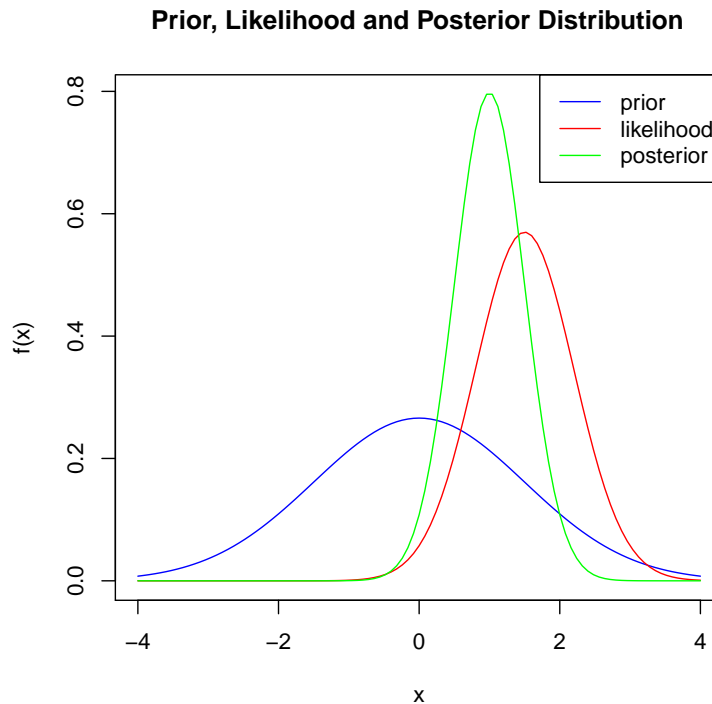


Figure 1. Triplot of prior, likelihood and posterior.

This means that we should be very careful about giving zero probability to something unless it is genuinely impossible. Once something has zero probability then no amount of further evidence can cause it to have a non-zero posterior probability.

- More generally, $f(\theta | x)$ will be low if either $f(\theta)$ is very small. We will tend to find that $f(x | \theta)$ is large when both $f(\theta)$ and $f(x | \theta)$ are relatively large, so that this θ value is given support by both information sources.

When θ is a scalar parameter, a useful diagram is the triplot, which shows the prior, likelihood and posterior on the same graph. An example is in Figure 1.¹

A strong information source in the triplot is indicated by a curve that is narrow (and therefore, because it integrates to one, also has a high peak). A narrow curves concentrates on a small range of θ values, and thereby "rules out" all values of θ outside that range.

¹All plots are generated in R, relevant codes are provided in Appendix R Codes

- Over the range $\theta < -1$, the likelihood: low
- Over the range $\theta > 3$, the likelihood: low
- Values of θ between -1 and 3 , the likelihood: high
- The maximum value of the posterior at: 1
- The Maximum likelihood estimation (MLE) of θ is ≈ 2

1.4.1. Normal Mean

For example, suppose that X_1, X_2, \dots, X_n are iid $\mathcal{N}(\mu, \sigma^2)$ and σ^2 is known. Then the likelihood is :

$$(1.17) \quad f(x | \mu) = \prod_{i=1}^n f(x_i | \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \\ \propto \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right).$$

As,

$$(1.18) \quad \sum (x_i - \bar{x} + \bar{x} - \mu)^2 = \sum (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 + 2(\bar{x} - \mu) \sum (x_i - \bar{x}) \\ = \sum (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \\ \propto \exp\left(-\frac{1}{2\sigma^2}n(\bar{x} - \mu)^2\right).$$

Note that $2(\bar{x} - \mu) \sum (x_i - \bar{x}) = 0$ as $\sum (x_i - \bar{x}) = 0$. Suppose the prior distribution for μ is normal:

$$(1.19) \quad \mu \sim \mathcal{N}(m, v).$$

Then applying Bayes' theorem we have:

$$(1.20) \quad f(\mu | x) \propto \underbrace{\exp\left(-\frac{1}{2\sigma^2}n(\bar{x} - \mu)^2\right)}_{f(x|\mu)} \underbrace{\exp\left(-\frac{1}{2v}n(\mu - m)^2\right)}_{f(\mu)} \\ = \exp\left(-\frac{\theta}{2}\right).$$

Note that

$$(1.21) \quad \theta = n\sigma^{-2}(\bar{x} - \mu) + v^{-1}(\mu - m)^2 = (v^*)^{-1}(\mu - m^*)^2 + R$$

and

$$(1.22) \quad v^* = (n\sigma^{-2} + v^{-1})^{-1}$$

$$(1.23) \quad m^* = (n\sigma^{-2} + v^{-1})^{-1}(n\sigma^{-2}\bar{x} + v^{-1}m) = a\bar{x} + (1 - a)m$$

where $a = n\sigma^{-2}/(n\sigma^{-2} + v^{-1})$

$$(1.24) \quad R = (n^{-1}\sigma^2 + v)(\bar{x} - m)^2$$

Therefore,

$$(1.25) \quad f(\mu | x) \propto \exp\left(-\frac{1}{2v}n(\mu - m)^2\right)$$

and we have shown that the posterior distribution is normal too: $\mu | x \sim \mathcal{N}(m^*, v^*)$

- m^* = weighted average of the mean m and the usual frequentist data-only estimate \bar{x} .
The weights \propto :
- Bayes' theorem typically works in this way. We usually find that posterior estimates are compromises between prior estimates and data based estimates and tend to be closer whichever information source is stronger. And we usually find that the posterior variance is smaller than the prior variance.

1.4.2. Weak Prior Information

It is the case where the prior information is much weaker than the data. This will occur, for instance, if we do not have strong information about Q before seeing the data, and if there are lots of data. Then in triplot, the prior distribution will be much broader and flatter than the likelihood. So the posterior is approximately proportional to the likelihood.

Example 1.4.1. In the normal mean analysis, we get weak prior information by letting the prior precision of v^{-1} become small. Then $m^* \rightarrow \bar{x}$ and $v^* \rightarrow \sigma^2/n$ so that the posterior distribution of μ corresponds very closely with standard frequentist theory.

LECTURE 2

Some Common Probability Distributions

Binomial on $Y \in \{0, 1, \dots, n\}$ with parameters $n \in \{1, 2, 3, \dots\}$ and $p \in (0, 1)$ is denoted by $Bi(n, p)$ and

$$(2.1) \quad f(y | n, p) = \binom{n}{y} p^y (1-p)^{n-y}$$

for $y = 0, 1, \dots, n$. The mean is given as:

$$(2.2) \quad E(y) = np.$$

Also the variance is given as:

$$(2.3) \quad v(y) = np(1-p).$$

Beta on $Y \in \{0, 1\}$ with parameters $a, b > 0$ is denoted by $Beta(p, q)$ and the density function is:

$$(2.4) \quad f(y | p, q) = \frac{\Gamma(p+q)}{\Gamma(p)\Gamma(q)} y^{p-1} (1-y)^{q-1}$$

for $y \in (0, 1)$. The mean is given as:

$$(2.5) \quad E(y) = \frac{p}{p+q},$$

Also the variance is given as:

$$(2.6) \quad v(y) = \frac{pq}{(p+q)^2(p+q+1)}.$$

$B(p, q) = \int_0^1 y^{p-1} (1-y)^{q-1} dy$ is the beta function and defined to be the normalizing constant for this density.

- In beta distribution, p and q change the shape of the distribution. Discuss!

Uniform distribution on $Y \in \{l, r\}$ where $-\infty < l < r < \infty$ is denoted by uniform (l, r) and its pdf is:

$$(2.7) \quad f(y | l, r) = \frac{1}{r - l}$$

for $y \in \{l, r\}$. The mean is given as:

$$(2.8) \quad E(y) = \frac{l + r}{2},$$

Also the variance is given as:

$$(2.9) \quad v(y) = \frac{(r - l)^2}{12}.$$

Poisson distribution on $Y \in \{0, 1, 2, \dots\}$ with parameter $\theta > 0$ is denoted by $Poisson(\theta)$ and its pdf is:

$$(2.10) \quad f(y | \theta) = \frac{\exp(-\theta)\theta^y}{y!}$$

for $y = 0, 1, 2, \dots$. The mean and the variance are given as[4]:

$$(2.11) \quad E(y) = v(y) = \theta.$$

Gamma distribution on $Y > 0$ with shape parameter $\alpha > 0$ and rate parameter $\lambda > 0$ is denoted by $Gamma(\alpha, \lambda)$ and the corresponding density is:

$$(2.12) \quad f(y | \alpha, \lambda) = \frac{\lambda^\alpha}{\Gamma(\alpha)} y^{\alpha-1} \exp(-\lambda y)$$

for $y > 0$. The mean is given as:

$$(2.13) \quad E(y) = \frac{\alpha}{\lambda},$$

Also the variance is given as:

$$(2.14) \quad v(y) = \frac{\alpha}{\lambda^2}.$$

Note that

$$(2.15) \quad \exp(\lambda) = Gamma(1, \lambda).$$

Univariate normal distribution on $Y \in \mathbb{R}$ with $\mu \in \mathbb{R}$ and variance $\sigma^2 > 0$ is denoted by $\mathcal{N}(\mu, \sigma^2)$ and its pdf is:

$$(2.16) \quad f(y | \mu, \sigma^2) = \frac{1}{\sigma} \left(\frac{1}{2\pi} \right)^{1/2} \exp \left\{ -\frac{1}{2\sigma^2} (y - \mu)^2 \right\}.$$

The mean is given as:

$$(2.17) \quad E(y) = \mu,$$

Also the variance is given as:

$$(2.18) \quad v(y) = \sigma^2.$$

K-variate normal distribution on $Y \in \mathbb{R}^k$ with vector $\mathbf{b} \in \mathbb{R}^k$ and positive definite symmetric (PDS) covariance matrix \mathbf{C} is denoted by $\mathcal{N}_k(\mathbf{b}, \mathbf{C})$ and the corresponding density function is:

$$(2.19) \quad f(y | \mathbf{b}, \mathbf{C}) = \frac{1}{\underbrace{|\mathbf{C}|^{1/2}}_{\text{determinant}}} \frac{1}{(2\pi)^{k/2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - \mathbf{b})^T \mathbf{C}^{-1} (y - \mathbf{b}) \right\}.$$

The mean is given as:

$$(2.20) \quad E(y) = \mathbf{b},$$

And the covariance matrix is given as:

$$(2.21) \quad Cov(y) = \mathbf{C}.$$

2.1. Posterior

Not only the beta distributions are the simplest and the most convenient distributions for a random variable confined to $[0, 1]$, they also work very nicely as prior distribution for a binomial observation. If $\theta \sim Be(p, q)$ then

$$(2.22) \quad f(x | \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}.$$

for $x = 1, 2, \dots, n$.

$$(2.23) \quad f(\theta) = \frac{1}{Be(p, q)} \theta^{p-1} (1 - \theta)^{q-1}$$

where $0 \leq \theta \leq 1$ and $p, q > 0$.

$$(2.24) \quad \begin{aligned} f(x) &= \int f(\theta) f(x | \theta) d\theta = \binom{n}{r} \frac{1}{Be(p, q)} \int_0^1 \theta^{p+x-1} (1 - \theta)^{q+n-x-1} d\theta \\ &= \binom{n}{r} \frac{Be(p+x, q+n-x)}{Be(p, q)}. \end{aligned}$$

From

$$(2.25) \quad f(\theta | x) = \frac{f(\theta) f(x | \theta)}{f(x)}.$$

(2.26)

$$f(\theta | x) = \frac{\theta^{p+x-1}(1-\theta)^{q+n-x-1}}{Be(p+x, q+n-x)} \propto \underbrace{\theta^{p-1}(1-\theta)^{q-1}}_{\text{Beta part}} \underbrace{\theta^x(1-\theta)^{n-x}}_{\text{Binomial part}}.$$

So $(\theta | x) \propto \text{Beta}(p+x, q+n-x)$. The posterior mean is:

$$(2.27) \quad E(\theta | x) = \frac{p+x}{p+q+n} = \frac{p+q}{p+q+n} E(\theta) + \frac{n}{p+q+n} \hat{\theta}$$

where $\hat{\theta} = x/n$. The posterior variance is:

$$(2.28) \quad v(\theta | x) = \frac{(p+x)(q+n-x)}{(p+q+n)^2(p+q+n+1)} \\ = \frac{E(\theta)(1-E(\theta | x))}{p+q+n+1}$$

But,

$$(2.29) \quad v(\theta) = \frac{E(\theta)(1-\theta)}{p+q+1}$$

So the posterior has higher relative precision than the prior.

SPECIAL NOTE:

The classical theory of estimation regards an estimator as good if it is unbiased and has small variance, or more generally if its mean-square-error is small. The MSE is an average squared error where the error is the difference between θ , i.e. y in previous notation, and the estimate t . In accordance with classical theory, the average is taken with respect to the sampling distribution of the estimator.

In Bayesian inference, θ is a random variable and it is therefore appropriate to average the squared error with respect to the posterior distribution of θ . Consider

$$(2.30) \quad E\{(t-\theta)^2 | x\} = E(t^2 | x) - E(2t\theta | x) + E(\theta^2 | x) \\ = t^2 - E(2t\theta | x) + E(\theta^2 | x) \\ = \{t - E(\theta | x)\}^2 + v(\theta | x).$$

Therefore the estimate t which minimizes posterior expected square error is $t = E(\theta | x)$, the posterior mean. The posterior mean can therefore be seen as an estimate of θ which is the best in the sense of minimizing expected squared error. This is distinct from, but clearly related to, its more natural role as a useful summary of location of the posterior distribution.

2.2. Weak Prior

If we reduce the prior relative precision to zero by setting $p = q = 0$, we get $\theta \mid x \sim Be(x, n - x)$. Then $E(\theta \mid x) = \hat{\theta}$ and $v(\theta \mid x) = \hat{\theta}(1 - \hat{\theta})/(n + 1)$ results which nicely parallel standard frequentist theory.

- Notice that we are not really allowed to let either parameter of the beta distribution be zero. However, by making p and q extremely small, we get as close to these results as we like. We can think of $p = q = 0$ as a defining limiting (if strictly improper) case of weak prior information.

Example 2.2.1. A doctor proposes a new treatment protocol for a certain kind of cancer. With current methods about 40% of patients with this cancer survive six months after diagnosis. After one year of using the new protocol, 15 patients with diagnosis, of whom 6 survived. After two years a further 55 patients have been followed to the six months mark, of whom 28 survived. So in total we have 34 patients surviving out of 70.

Let θ be the true success rate of the new treatment protocol, i.e. the true proportion of patients who survive 6 months and we wish to make comparison of θ with the current survival rate 40%.

Suppose that the doctor in charge has prior information leading her to assign a prior distribution with expectation $E(\theta) = 0.45$, i.e. expects a slight improvement over the existing protocol, from 40% to 45%, however her prior standard deviation is 0.07, $v(\theta) = 0.07^2 = 0.0049$

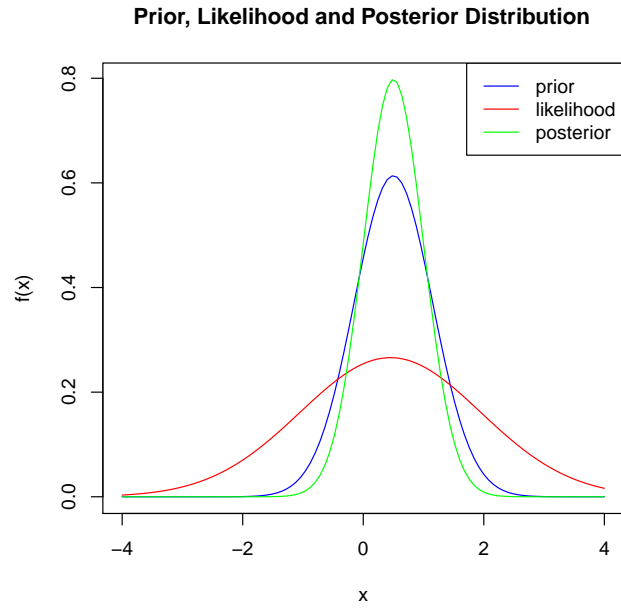


Figure 1. The triplot from first year's data.

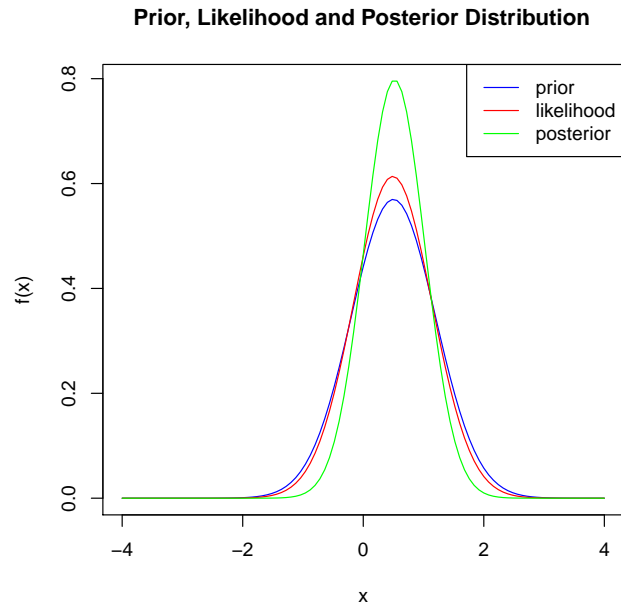


Figure 2. The triplot from two years' data.

2.3. Sequential Updating

In the last example we pooled the data from the two years and went back to the original prior distribution to use Bayes' theorem. We did not need to do this. A nice feature of Bayes' theorem is the possibility of updating sequentially, incorporating data as they arrive. In this case, consider the data to be just the new patients observed to a six months follow-up during the second year. These comprise 55 patients, of whom 28 had survived. The doctor could consider these as the data x with $n = 55$ and $r = 28$. What would the prior information be?

Clearly, the prior distribution should express her information prior to obtaining these new data, i.e. after the first years' data, so her prior for this second analysis is her posterior distribution from the first. This was $Be(28.28, 36.23)$. Combining this prior with the new data gives the same posterior $Be(28.28 + 28, 36.23 + 27) = Be(56.28, 63.23)$ as before. This simply confirms that we can get to the posterior distribution.

- In a single step, combining all the data with a prior distribution representing information available before any of the data were obtained.
- Sequentially, combining each item or block of new data with a prior distribution representing information available just before the new data were obtained (but after getting data previously received).

2.4. Normal Sample

Let X_1, X_2, \dots, X_n be from $\mathcal{N}(\mu, \sigma^2)$. $\theta = (\mu, \sigma^2) \rightarrow$ unknown parameters. The likelihood is:

$$(2.31) \quad \begin{aligned} f(x | \mu, \sigma^2) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right\} \\ &\propto \sigma^{-n} \exp \left[-\frac{1}{2\sigma^2} \{n(\bar{x} - \mu)^2 + S^2\} \right] \end{aligned}$$

where $S^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

2.5. NIC distributions

For the prior distribution, we now need a joint distribution for μ and σ^2

Definition 2.32. The normal-inverse-chi-squared distribution(NIC) has density:

$$(2.33) \quad f(x | \mu, \sigma^2) \propto \sigma^{-(d+3)/2} \exp \left[-\frac{1}{2\sigma^2} \{v^{-1}(\mu - m)^2 + a\} \right]$$

where $a > 0$, $d > 0$ and $v > 0$.

The following facts are easy to derive about $NIC(m, v, a, d)$ distribution.

- (a) The conditional distribution of μ given σ^2 is $\mathcal{N}(\mu, v\sigma^2)$ so $E(\mu | \sigma^2) = m$, $v(\mu | \sigma^2) = v\sigma^2$.
- (b) The marginal distribution of σ^2 is such that $a\sigma^{-2} \sim \chi_d^2$. We say that σ^2 has the inverse-chi-square distribution $IC(a, d)$. We have $E(\sigma^2) = a/(d-2)$ if $d > 2$ and $v(\sigma^2) = 2a^2 / \{(d-2)^2(d-4)\}$ if $d > 4$.
- (c) The conditional distribution of σ^2 given μ is $IC(v^{-1}(\mu - m)^2 + a, d+1)$ and in particular $E(\sigma^2 | \mu) = (v^{-1}(\mu - m)^2 + a)/(d-1)$ provided $d > 1$.
- (d) The marginal distribution of μ is such that $(\mu - m)\sqrt{d}/\sqrt{av}\mu + d$. We say that μ has t -distribution $t_d(m, av/d)$. We have $E(\mu) = m$ if $d > 1$, and $v(\mu) = av/(d-2)$ if $d > 2$.

2.6. Posterior

Supposing then that the prior distribution is $NIC(m, v, a, d)$, we find

$$(2.34) \quad f(\mu, \sigma^2 | x) \propto \sigma^{d+n+3} \exp \left[-\frac{1}{2\sigma^2} \theta \right]$$

where $\theta = v^{-1}(\mu - m)^2 + a + n(\bar{x} - \mu) + s^2$ is a quadratic expression in μ . After completing the square, we see that $\mu, \sigma^2 | x \propto NIC(m^*, v^*, a^*, d^*)$ where $m^* = (v^{-1}m + n\bar{x})/(v^{-1} + n)$, $v^* = (v^{-1} + n)^{-1}$, $a^* = a + S^2 + (\bar{x} - m)^2/(n^{-1} + v)$, $d^* = d + n$. To interpret these results, note first that the posterior mean of μ is m^* which is a weighted average of the prior mean m and the usual data only-estimate \bar{x} with weights v^{-1} and n .

The posterior mean of σ^2 is $a^*/(d^* - 2)$ which is a weighted average of three terms: the prior mean $a/(d-2)$ with weight $(d-2)$, the usual data-only estimate $S^2/(n-1)$ with weight $(n-1)$ and $(\bar{x} - m)/(n^{-1} + v)$ with weight 1.

2.7. Weak prior

We clearly obtain weak prior information about μ by letting v go to infinity or $v^{-1} \rightarrow 0$. Then $m^* = \bar{x}$, $v^* = 1/n$, $a^* = a + S^2$, because the third term disappears.

To obtain weak prior information also about σ^2 , if it is usual to set $a = 0$ and $d = 1$. Then $a^* = S^2$ and $d^* = n - 1$. The resulting inference match the standard frequentist results very closely with these parameters, since we have:

$$(2.35) \quad \frac{(\mu - \bar{x})\sqrt{n}}{S/\sqrt{n-1}} \propto t_{n-1},$$

$$(2.36) \quad \frac{S^2}{\sigma^2} \propto \chi_{n-1}^2$$

Exactly the same distribution statements underlie standard frequentist inference in this problem.

LECTURE 3

Inference

Summarisation: Here all inferences are derived from the posterior distribution. In frequentist statistics, there are three kinds of inference:

- Point estimation: unbiasedness, minimum variance estimation
- Interval estimation: Credible intervals in Bayesian inference, confidence interval in frequentist approach
- Hypothesis testing: Significance test.

In Bayesian inference, the posterior distribution expresses all that is known about θ . So it uses an appropriate summaries of the posterior distribution to describe the main features of what we now know about θ .

Plots: To draw the posterior density

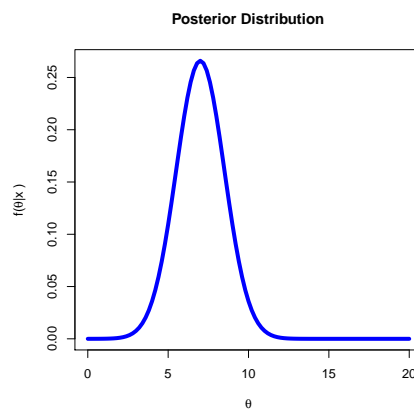


Figure 1. A posterior density plot.

For a bivariate parameter, we can still usefully draw the density as a perspective plot or a contour plot.

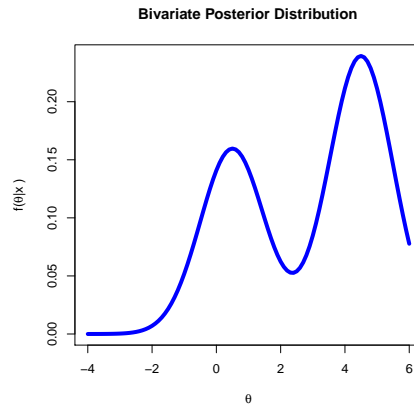
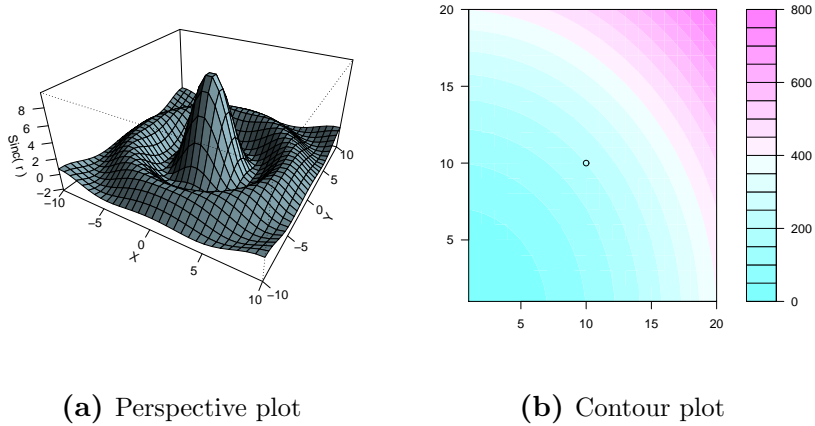


Figure 3. Marginal Densities.

3.1. Shape

In general, plots illustrate the shape of the posterior distribution. Important features of shape are modes (and antimodes), skewness and kurtosis (peakedness or heavy tails). The quantitative summaries of shape are needed to supplement like the view of mode (antimode).

The first task is to identify turning points of the density, i.e. solutions of $f'(\theta) = 0$. Such points include local maxima and minima of $f(\theta)$ which we call mode and antimode, respectively.

A point θ_0 is characterized as a mode if $f'(\theta_0) = 0$ and $f''(\theta_0) < 0$ whereas it is an antimode if $f'(\theta_0) = 0$ and $f''(\theta_0) > 0$. Any point θ_0 for which $f''(\theta_0) = 0$ is a point of inflection of the density (whether or not $f'(\theta_0) = 0$).

Example 3.1.1. Consider the gamma density

$$(3.1) \quad f(\theta) = \frac{a^b}{\Gamma(b)} \theta^{b-1} e^{-a\theta} \quad ; \quad \theta > 0$$

where a, b are positive constants.

$$(3.2) \quad f'(\theta) = \frac{a^b}{\Gamma(b)} \{(b-1) - a\theta\} \theta^{b-2} e^{-a\theta}$$

$$(3.3) \quad f''(\theta) = \frac{a^b}{\Gamma(b)} \{a^2\theta^2 - 2a(b-1)\theta - (b-1)(b-2)\} \theta^{b-3} e^{-a\theta}$$

So from $f'(\theta)$, the turning point at $\theta = (b-1)/a$. For $b \leq 1$, $f'(\theta) < 0$ for all $\theta \geq 0$, so $f(\theta)$ is monotonic decreasing and the mode is at $\theta = 0$. For $\theta > 1$, $f(\theta) \rightarrow 0$ as $\theta \rightarrow 0$, so the turning point $\theta = 0$ is not a mode. In this case, $f'(\theta) > 0$ for $\theta < (b-1)/a$ and $f'(\theta) < 0$ for $\theta > (b-1)/a$. Therefore $\theta = (b-1)/a$ is the mode. Looking at $f''(\theta)$, the quadratic expression has roots at $\theta = \frac{b-1}{a} \mp \frac{(b-1)^{1/2}}{a}$. Therefore $b > 1$, these are the points of inflection.

Example 3.1.2. Consider the mixture of two normal distributions

$$(3.4) \quad f(\theta) = \underbrace{\frac{0.8}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\theta^2\right\}}_{\sim \mathcal{N} \text{ with weight}} + \underbrace{\frac{0.2}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(\theta-4)^2\right\}}_{\sim \mathcal{N} \text{ with weight}}$$

$$(3.5) \quad f'(\theta) = -\frac{0.8\theta}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\theta^2\right\} - \frac{0.2(\theta-4)}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(\theta-4)^2\right\}$$

For $\theta \leq 0$, $f'(\theta) > 0$ and for $\theta \geq 0$, $f'(\theta) < 0$, the turning points at $\theta = 0.00034, 2.46498$ and 3.9945 .

$$(3.6) \quad f''(\theta) = -\frac{0.8(\theta^2-1)}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\theta^2\right\} - \frac{0.2((\theta^2-8\theta-15))}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(\theta-4)^2\right\}$$

This is positive for $\theta \leq -1$, for $1 \leq \theta \leq 3$ and $\theta \geq 5$, confirming that the middle turning point is an antimode. Calculating $f''(\theta)$ at the other points confirms them to be modes. Finally points of inflection are at $\theta = -0.99998, \theta = 0.98254, \theta = 3.17903, \theta = 4.99971$.

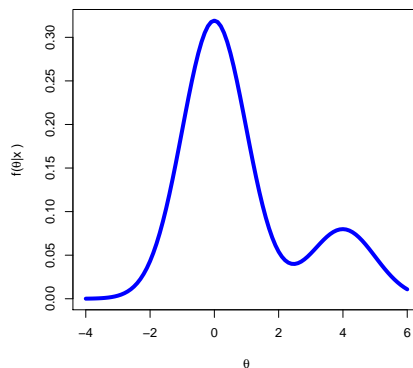


Figure 4. Plot of mixture of two normal distributions

3.2. Visualizing multivariate densities

Turning points:

At a mode, or at a turning point generally, the gradients of function in all directions are zero. Therefore, the turning points are solutions of the simultaneous equations:

$$(3.7) \quad 0 = \frac{\partial f(\theta, \phi)}{\partial \phi} = \frac{\partial f(\theta, \phi)}{\partial \theta}.$$

The turning points may be classified by examining the symmetric matrix $F''(\theta, \phi)$ of the second order partial derivatives. For instance, for a bivariate density $f(\theta, \phi)$,

$$(3.8) \quad F''(\theta, \phi) = \begin{bmatrix} \frac{\partial^2 f(\theta, \phi)}{\partial \theta^2} & \frac{\partial^2 f(\theta, \phi)}{\partial \theta \partial \phi} \\ \frac{\partial^2 f(\theta, \phi)}{\partial \phi \partial \theta} & \frac{\partial^2 f(\theta, \phi)}{\partial \phi^2} \end{bmatrix}$$

$F''(\theta, \phi)$ is known as the Hessian matrix. The second derivative of $f(\theta)$, in a direction t is after differentiating. At a mode, this must be negative in all directions. So that the Hessian matrix is negative definite. Similarly, at an antimode, it is positive definite. In the intermediate case, where $F''(\theta, \phi)$ is indefinite, i.e. has both positive and negative eigenvalues, we have a saddle point.

If $F''(\theta, \phi)$ is positive definite, regions are:

$F''(\theta, \phi)$ are negative definite, regions are:

$F''(\theta, \phi)$ is indefinite, indefinite curvature. On the boundaries between these regions, one eigenvalue of $F''(\theta, \phi)$ is zero, so all points on such boundaries are inflection points.

A point of inflection corresponds to the second derivative being zero

in some direction t , therefore inflection points are characterized by $F''(\theta, \phi)$ being singular. In fact, in more than two dimensions, we can further subdivide the regions of indefinite curvature according to how many positive eigenvalues of $F''(\theta, \phi)$ has, and all these subregions are also separated by inflection boundaries.

Location:

A plot gives a good idea of location, but the conventional location measures for distributions are also useful. These include the mean, mode and median.

Dispersion:

The usual dispersion measure is the variance, or for a multivariate distribution the variance-covariance matrix.

Dependence:

It is important with multivariate distributions to summarize the dependence between individual parameters. This can be done with correlation coefficients, but plots of regression functions (conditional mean functions) can be more informative.

3.3. Informal Inferences

- (a) Point estimation: The obvious posterior estimate of θ is its posterior mean $\hat{\theta} = E(\theta | x)$. Modes and medians are also natural point estimates, and they all have intuitively different interpretations. The mean is the expected value, the median is the central value and the mode is the most probable value.
- (b) Interval estimation: If asked to provide an interval in which θ probably lies, we can readily derive such a thing from its posterior distribution. For instance, in the density shown on page 1, there is a probability 0.05 to the left of $\theta = 3.28$ and also 0.05 to the right of $\theta = 11.84$. So the interval $(3.28, 11.84)$ is a 90% posterior probability for θ . We call such an interval a credible interval.
 - If a frequentist had found this interval, it means that it would say that if we repeatedly draw samples of data from the same population, and applied the rule that was used to derive this particular interval to each off those datasets, then 90% of those intervals would contain θ .

If a Bayesian approach, there is a posterior probability 0.9 that θ lies between 3.28 between 11.84.

Definition 3.9. A $100(1 - \alpha)\%$ credible set for θ is a subset C such that:

$$(3.10) \quad 1 - \alpha \leq P(C | y) = \int_C p(\theta | y) d\theta.$$

where integration is replaced by summation for discrete components.

Definition 3.11. The exact possible coverage of $(1 - \alpha)$ can be found by the highest posterior density of HPD credible set as the set:

$$(3.12) \quad C = \{\theta \in \Theta : p(\theta | y) \geq k(\alpha)\}.$$

where $k(\alpha)$ is the largest constant satisfying $P(C | y) \geq 1 - \alpha$.

For 2-sided credible set, we can generally take the $\alpha/2$ and $(1 - \alpha/2)$ quantiles of $p(\theta | y)$ as our $100(1 - \alpha)\%$ credible set for θ . This equal tail credible set will be equal to the HPD credible set if the posterior is symmetric and unimodal, but will be a bit wider otherwise.

- (c) Evaluating hypothesis: Suppose we wish to test a hypothesis H which asserts that θ lies in some region A . The Bayesian way to test to the hypothesis is simply to calculate the (posterior) probability that it is true: $P(\theta \in A | x)$.

3.4. Multivariate inference

All the above treated θ as a scalar parameter. If we have a vector θ , then in general we can consider inference of the above forms about any scalar function $\phi = g(\theta)$. The inferences are then derived simply from the marginal posterior distribution of ϕ .

LECTURE 4

Formal Inference

Suppose that we want to answer a question that falls neatly into the frequentist point estimation framework, "What is the best estimate of θ ?". In the frequentist theory, we need to be explicit about what we would regard as good properties for an estimator in order to identify a best one.

The Bayesian approach also needs to know what makes a good estimate before an answer can be given. This falls into the framework of formal inference.

Formally, we would seek the estimate that minimises expected square error. So the expectation is derived from the posterior distribution. We want to minimise:

$$(4.1) \quad \begin{aligned} E((\hat{\theta} - \theta)^2 | x) &= \hat{\theta}^2 - 2\hat{\theta}E(\theta | x) + E(\theta^2 | x) \\ &= (\hat{\theta} - E(\theta | x))^2 + v(\theta | x). \end{aligned}$$

So the estimate $\hat{\theta}$ that minimises this expected squared error is $\hat{\theta} = E(\theta | x)$.

4.1. Utility and decisions

Formal inference aims to obtain optimal answer to inference questions. This is done with reference to a measure of how good or bad the various possible inferences would be deemed to be if we knew the true value of θ . This measure is a utility function. Formally, $u(d, \theta)$ defines the value of inference d if the true value of the parameter is θ . Formal inference casts an inference problem as a decision problem. A decision problem is characterized by:

- a set Ω of possible parameter values
- a probability distribution for $\theta \in \Omega$
- a set D of possible decisions

- utility function $u(d, \theta)$ for $d \in D$ and $\theta \in \Omega$.

The solution is

$$(4.2) \quad \boxed{d_{opt} = \arg \max_d E_\theta(u(d, \theta))}$$

Here the distribution of θ is its prior distribution.

- In inference problems, we generally define a measure of badness of an inference, which we call a loss function $L(d, \theta)$. We can simply define utility to be negative loss, and then the optimal inference is the one which minimises posterior expected loss.

4.1.1. Formal Point Estimation

The set D is the set of all possible values of θ . We have seen that if we use squared error loss (which is implicitly the measure used by frequentist in considering mean-squared-error of the variance of an unbiased estimator), formally defining $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$, then the posterior mean is the optimal estimator. If we use absolute error loss $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$, then the optimal estimator is the posterior median.

4.1.2. Formal Interval Estimation

The possible inferences now are interval or, more generally, subsets of possible values of θ . A loss function will penalize an interval if it fails to contain the true value of θ . So the optimal interval is found from this form:

$$(4.3) \quad \boxed{d_{opt} = \{\theta : f(\theta |) \geq t\}}$$

where t is chosen to obtain the desired $P(\theta \in d_{opt} | x)$. Such a credible interval is called a highest (posterior) density interval (HPD interval or HDI).

Example 4.1.1. From the following figure 1 which shows 90% credible interval, it is not HPI. Because here the upper and lower limits should have the same posterior density, and shows that the density at $\theta = 3.28$ is higher than that at $\theta = 11.84$. So the 90% highest density interval is actually (2.78, 11.06).

4.2. Formal Hypothesis Testing

If we really need to decide whether (to act as if) hypothesis that $\theta \in A$, is true, there are just two inferences. d_0 is to say it is true, while d_1

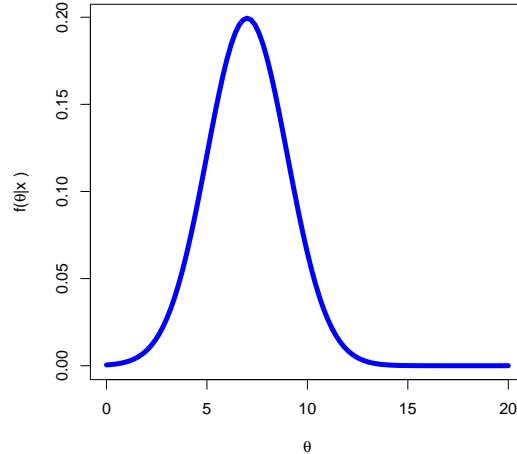


Figure 1. A posterior density plot

says it is false. The loss function will take the form:

$$(4.4) \quad \begin{aligned} L(d_0, \theta) &= 0 & \text{if } \theta \in A \\ &= 1 & \text{if } \theta \notin A \end{aligned}$$

$$(4.5) \quad \begin{aligned} L(d_1, \theta) &= k & \text{if } \theta \in A \\ &= 0 & \text{if } \theta \notin A \end{aligned}$$

where k defines the relative seriousness of the first kind of error relative to the second. Then $E_\theta(L(d_0, \theta)) = P(\theta \notin A | x)$, while $E_\theta(L(d_1, \theta)) = kP(\theta \in A | x)$. The optimal decision is to select d_0 (say that H is true) if its probability $P(\theta \in A | x)$ exceeds $\frac{1}{k+1}$. The greater the relative seriousness k of the first kind of error, the more willing we are:

4.3. Nuisance Parameter

In any inference problem, the parameter(s) that we wish to make inference about is (are) called the parameter of interest, and the remainder of components of θ is (are) called nuisance parameters.

Example 4.3.1. If we have a sample from $\mathcal{N}(\mu, \sigma^2)$ and we wish to make inference about μ , then nuisance parameter is σ^2 .

If $\theta = (\phi, \psi)$ and ψ is the vector of nuisance parameters, then the inference about ϕ is made from marginal value posterior distribution.

4.4. Transformation

If $\hat{\theta}$ is an estimate of θ , is $g(\hat{\theta})$ the appropriate estimate of ϕ ?

This depends on the kind of inference being made. In the particular case of point estimation, then the posterior mean is not invariant in this way.

Example 4.4.1. If $\phi = \theta^2$, then

$$(4.6) \quad E(\phi | x) = E(\theta^2 | x) = v(\theta | x) + E(\theta | x)^2 \geq E(\theta | x)^2$$

The mode is not invariant to transformations but the median is invariant, at least to 1 – 1 transformations.

Interval estimates are also invariant to 1 – 1 transformations in the sense that if $[a, b]$ is a 90% interval, say for θ , then $[g(a), g(b)]$ is a 90% interval for ϕ if g is a monotone increasing function. If $[a, b]$ is a 90% HPD interval for θ , then $[g(a), g(b)]$ is an HPD interval for ϕ ?

4.5. The prior distribution

The nature of probability:

- (a) Frequency probability: Frequentist statistics uses the familiar idea that the probability of an event is the limiting relative frequency with which that event would occur in an infinite sequence of repetitions. For this definition of probability to apply, it is necessary for the event to be, at least in principle, repeatable.
- (b) Personal probability: Bayesian statistics is based on defining the probability of a proposition to be a measure of a person's degree of belief is the truth of that proposition.

In the Bayesian framework, wherever there is uncertainty there is probability.

In particular, parameters have probability distributions.

4.6. Subjectivity

The main critic to Bayesian methods is the subjectivity due to the prior density.

If the data are sufficiently strong, the remaining element of personal judgement will not matter, because all priors based on reasonable interpretation of the prior information will lead to effectively the same posterior inferences. Then we can claim robust conclusion on the basis of the synthesis of prior information and data.

If the data are not that strong, then we do not yet have enough scientific evidence to reach an objective conclusion. Any method which

claims to produce a definitive answer in such a situation is misleading, so this is actually a strength of the Bayesian approach.

4.7. Noninformative Priors

The basis of this is that if we have a completely flat prior distribution such that $f(\theta)$ is a constant, then the posterior density is proportional to the likelihood and inferences will be based only on the data. If we can do this, we can get the other benefits of Bayesian analysis, such as having more meaningful inferences that actually answer the question, but without the supposed disadvantage of subjectivity.

The main problem with this neat solution is that it can not be applied consistently.

Example 4.7.1. $f(\theta) = 1$ for all $\theta \in [0, 1]$ uniform distribution represents complete ignorance about θ .

If θ is ignored, then $\phi = \theta^2$, which also takes values in $[0, 1]$, is also completely ignored.

But the implied distribution $f(\phi) = 1$ is not consistent with the previous specification of $f(\theta) = 1$. The uniform prior distribution for θ implies that $\phi = \theta^2$ should have the density $f(\phi) = \frac{1}{2\sqrt{\phi}}$. Conversely, if ϕ has a uniform prior distribution, then the implied prior for θ has density $f(\theta) = 2\theta$.

In general, a uniform prior for θ translates into a non-uniform prior for any function of θ . Another complication is that if the range of possible values of θ

Another complication is that if the range of possible values of θ is bounded then we can not properly give it a uniform distribution. For instance, if $\theta \in [0, \infty)$ and we try to define a prior distribution $f(\theta) = c$ for some constant c , then there is no value of c that will make this density integrate to 1. For $c = 0$, it integrates to 0, and for any positive c it integrates to infinity. In these situations, we appeal to proportionality and simply write $f(\theta) \propto 1$

A distribution expressed as $f(\theta) \propto h(\theta)$ when there can not be any proportionality constant that would make this into a proper density function, is called an improper distribution. This arises whenever the integral of $h(\theta)$ over the range of possible values of θ does not change.

Example 4.7.2. For a parameter $\theta \in [0, 1]$, three favourite recommendations are $f(\theta) = 1$, $f(\theta) = \pi^{-1}\theta^{-1/2}(1 - \theta)^{-1/2}$ and $f(\theta) \propto \theta^{-1}(1 - \theta)^{-1}$, the last of these being improper. We can identify these as the $Be(1, 1)$, $Be(1/2, 1/2)$ and $Be(0, 0)$ distributions.

Improper distributions are not in fact usually much of a problem, since we can appeal to proportionality. That is, the absence of a well-defined proportionality constant is ignored and assumed to cancel in the proportionality constant of Bayes' theorem [2]. In effect, we are obtaining the limit of the posterior distribution as we go through a range of increasingly flat priors towards the uniform limit.

- (a) The posterior distribution may also be improper. In this case, technically, the limit of the above process is not well-defined. Improper prior distribution should never be used when the resulting posterior distribution is improper, so it is important to verify propriety of the posterior.
- (b) When comparing different models for the data, improper distributions always lead to undefined model comparisons. This is an area outside the scope of this course, but very important in practice.

4.8. Informative Priors

So in specifying an informative prior distribution, (a) we specify values for whatever summaries best express the features of the prior information, then (b) we simply choose any conventional $f(\theta)$ that has those summaries.

Example 4.8.1. I wish to formulate my prior beliefs about number N of students who will turn up to one of my lectures. I first ask myself what my best estimate would be, and I decide on 38, so I set $E(N) = 38$. I next ask myself how far wrong this estimate might be. I decide that the actual number could be as high as 48 or as low as 30, but I think the probability of the actual number being outside that range is small, maybe only 10%. Now a convenient prior distribution that matches these summaries is the Poisson distribution.

So it has mean 38 and $P(30 \leq N \leq 48) = 0.87$, which seems a good enough fit to my specified summaries.

4.9. Prior Choices

There are several alternatives to overcome the problem of improper prior.

- (a) Jeffrey's Prior:

Let $I(\theta)$ be the Fisher information:

$$(4.7) \quad I(\theta) = -E \left\{ \frac{\partial^2 \log f(x | \theta)}{\partial \theta^2} \right\}.$$

In the case of a vector parameter, $I(\theta)$ is the matrix formed as minus the expectation of the matrix of second order partial derivatives of $\log f(x | \theta)$. The Jeffrey's prior distribution is then:

$$(4.8) \quad \boxed{f_0(\theta) \propto |I(\theta)|^{1/2}}$$

Example 4.9.1. If x_1, x_2, \dots, x_n are normally distributed with mean θ and known variance v , then

$$(4.9) \quad f(x | \theta) \propto \exp \left\{ -\frac{n(\bar{x} - \theta)^2}{2v} \right\}$$

What is the Jeffrey's prior for this distribution?

Solution:

$$(4.10) \quad \log f(x | \theta) = -\frac{n(\bar{x} - \theta)^2}{2v}$$

$$(4.11) \quad \frac{d^2}{d\theta^2} \log f(x | \theta) = -\frac{n}{2v}$$

Therefore,

$$(4.12) \quad I(\theta) = -E\left(\frac{d^2}{d\theta^2} \log f(x | \theta)\right) = \frac{n}{2v}$$

As a result,

$$(4.13) \quad f_0(\theta) = \sqrt{I(\theta)} = \sqrt{\frac{n}{2v}}.$$

Example 4.9.2. If x_1, x_2, \dots, x_n are distributed as $\mathcal{N}(\mu, \sigma^2)$ with $\theta = (\mu, \sigma^2)$, then

$$(4.14) \quad f(x | \theta) \propto \sigma^{-n} \exp \left\{ -\frac{-n(s + (\bar{x} - \theta)^2)}{2\sigma^2} \right\}$$

where $s = \frac{\sum(x_i - \bar{x})^2}{n}$, Then what is the Jeffrey's prior of $f(\mu, \sigma^2)$?

Solution:

Example 4.9.3. If x_1, x_2, \dots, x_n are normally distributed with known mean m and variance θ , then

$$(4.15) \quad f(x | \theta) \propto \theta^{-n/2} \exp \left\{ -\frac{-s}{(2\theta)^2} \right\}$$

where $s = \sum (x_i - m)^2$, then what is the Jeffrey's prior for θ ?

Solution:

$$(4.16) \quad \log f(x | \theta) =$$

A number of objections can be made to the Jeffrey's prior, the most important of which is that it depends on the form of the data. The prior distribution should only represent the prior information, and not be influenced by what data are to be collected.

(b) Maximum Entropy:

The entropy $H(f) = - \int_{-\infty}^{\infty} f(\theta) \log f(\theta) d\theta$ of the density $f(\theta)$ can be thought of a measure of how uninformative $f(\theta)$ is about θ . For if we try to convert our information about θ as a general form of inference in the scoring rule framework, $H(f)$ is the lowest obtainable expected loss. If $H(f)$ is high, then the best decision is still poor. Now to represent prior ignorance we could use the prior density $f(\theta)$ which maximizes the entropy.

Example 4.9.4. Suppose that θ is discrete with possible values $\theta_1, \theta_2, \dots, \theta_k$. The prior distribution with maximum entropy will then maximize $\sum_{i=1}^k p_i \log p_i + \lambda \sum p_i$, where λ is a Lagrange multiplier, $\partial F / \partial p_i = \log p_i + 1 + \lambda$. Equating this to zero yields the solution $p_i = k^{-1}$, $i = 1, 2, 3, \dots, k$. That is the maximum entropy prior is the uniform distribution.

The primary criticism of this approach is that it is not invariant under change of parametrization, the problem which the Jeffrey's prior was designed to avoid. In general, unrestricted maximization of entropy leads to a uniform prior distribution, which was shown to be sensitive to parametrization.

(c) Reference Prior: The expected amount of information provided by observing x is given by

$$(4.17) \quad H \{f(\theta)\} - E [H \{f(\theta | x)\}],$$

where the expectation is over the preposterior distribution of $f(x)$ of x . If the experiment yielding x were to be repeated, giving a new observation independent of x given θ and with the same distribution, the posterior distribution would be expected to show a further reduction in entropy, representing the expected information in the second observation. If this were

repeated indefinitely we would eventually learn θ exactly and so remove all the entropy in the original prior distribution.

In the case of discrete θ taking a finite of possible values, this process reduces ??? maximizing prior entropy, and so gives the uniform distribution. This is not the case for continuous θ . It is shown that under appropriate regularity conditions, the reference prior distribution is the Jeffrey prior.

LECTURE 5

Structuring Prior Information

Independence: Suppose that x_1, x_2, \dots, x_n are a sample from the $\mathcal{N}(\mu, 1)$ distribution, which we write formally as:

$$(5.1) \quad x_i \mid \mu \sim \mathcal{N}(\mu, 1)$$

, independent. From a Bayesian perspective, what is meant here is conditional dependence. That is the x_i 's are independent given μ .

Exchangeability: It is same as the independence in frequentist approach.

Definition 5.2. Random variables x_1, x_2, \dots, x_m are said to be exchangeable if their joint distribution is unaffected by permuting the order of the x_i 's.

So first consider the (marginal) distribution of x_1 . The definition says that every one of the x_i 's must have the same marginal distribution, because we can permute them so that any desired x_i comes into the first position in the sequence. So one implication of exchangeability is that the random variables in question are identically distributed.

next consider the joint distribution of x_1 and x_2 Exchangeability means that every pair (x_i, x_j) (for $i \neq j$) has the same bivariate distribution as (x_1, x_2) . In particular, the correlation between any pair of random variables is the same.

And so it goes to higher order joint distributions. The joint distributions (x_1, x_2, \dots, x_k) is the same as that of any other collection of k distinct x_i 's. This is the meaning of exchangeability.

- In general, suppose that x_i 's have a common distribution $g(x | \theta)$, and are independent, given θ . Then the joint density is:

$$\begin{aligned}
 f(x_1, x_2, \dots, x_m) &= \int f(x_1, x_2, \dots, x_m, \theta) d\theta \\
 (5.3) \qquad \qquad \qquad &= \int f(x_1, x_2, \dots, x_m) f(\theta) d\theta \\
 &= \int \prod_{i=1}^m g(x_i | \theta) f(\theta) d\theta.
 \end{aligned}$$

which is unaffected by permuting the x_i 's. Their common marginal density is:

$$(5.4) \qquad \qquad \qquad f(x) = \int g(x | \theta) f(\theta) d\theta,$$

and the common distribution of any pair of the x_i 's is

$$(5.5) \qquad \qquad \qquad f(x, y) = \int g(x | \theta) g(y | \theta) f(\theta) d\theta,$$

This is generally what a frequentist means by the x_i 's being iid, or being a random sample from the distribution $g(x | \theta)$. From the Bayesian perspective all frequentist statements are conditional on the parameters. So exchangeability is the Bayesian concept that corresponds very precisely with the frequentist idea of a random sample.

5.1. Binary Exchangeability

Theorem 5.6 (De Finetti, 1937). *Let x_1, x_2, x_3, \dots be an infinite sequence of exchangeable binary random variables. Then their joint distribution is characterised by a distribution $f(\theta)$ for a parameter $\theta \in [0, 1]$ such that the x_i 's are independent of θ with $P(x_i = 1 | \theta) = \theta$.*

De Finetti's theorem says that if we have a sequence of binary variables, so that each x_i takes the value 1 (success) or 0 (failure), then we can represent them as independent Bernoulli trials with probability θ of success in each trial.

5.2. Exchangeable Parameters

Consider the simple one-way analysis of variance model:

$$(5.7) \qquad y_{ij} \sim \begin{cases} \mathcal{N}(\mu_i, \sigma^2), & i = 1, 2, \dots, k \\ j = 1, 2, \dots, n_i \end{cases}, \text{ independent}$$

where $\theta = (\mu_1, \mu_2, \dots, \mu_k, \sigma^2)$. This says that we have k independent normal samples, where the i -th sample has mean μ_i and size n_i , and where all the observations have a common variance σ^2 .

For this model, we need to specify a joint prior distribution for $\mu_1, \mu_2, \dots, \mu_k$ and σ^2 . Now when it comes to formulating a joint prior distribution for many parameters, it is much easier if we can regard them as independent. Then we could write:

$$(5.8) \quad f(\mu_1, \mu_2, \dots, \mu_k, \sigma^2) = f(\sigma^2) \prod_{i=1}^k f(\mu_i),$$

and we would only need to specify the prior distribution of each parameter separately. Unfortunately, this is unlikely to be the case in practice with this model.

The model is generally used when the samples are from a related or similar populations. For instance, they might be weight gains of pigs given k different diets, and μ_i is the mean weight gain with diet i . In this sort of situation, the μ_i 's would not be independent. Then we could add to the "prior model"

$$(5.9) \quad \mu_i \mid \xi, \tau^2 \sim \mathcal{N}(\xi, \tau^2)$$

, independent, which says that the μ_i 's are drawn from a normal (assumed) population with unknown mean ξ and unknown variance τ^2 .

The prior distribution could then be completed by specifying a joint prior distribution $f(\xi, \sigma^2, \tau^2)$ for the remaining parameters.

5.3. Hierarchical Models

The kind of modelling seen in the previous section is called "hierarchical". In general, we can consider a model of the form:

- Data model: $f(x \mid \theta)$
- First level of prior: $f(\theta \mid \phi)$
- Second level of prior: $f(\phi)$

We often refer to ϕ as the hyperparameter(s)

- If we are only interested in the parameter θ of the original data model,

$$(5.10) \quad f(\theta) = \int f(\theta \mid \phi) f(\phi) d\phi.$$

- Let us be actually interested in the hyperparameters ϕ ,

$$(5.11) \quad f(x \mid \phi) =$$

$$(5.12) \quad f(\theta, \phi) =$$

$$(5.13) \quad f(\theta, \phi | x) \propto$$

$$(5.14) \quad f(\theta | x) =$$

Shrinkage: It means that the posterior distribution and posterior estimates of these parameters will generally be closer together than their corresponding data estimates. This is a phenomenon known as "shrinkage". Let $w = (\xi, \sigma^2, \tau^2)$ and $\mu = (\mu_1, \mu_2, \dots, \mu_k)$ and

$$(5.15) \quad \begin{aligned} f(\mu | x, w) &\propto f(x | \mu, \sigma^2) f(\mu | w) \\ &= \prod_{i=1}^k \left\{ \prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{(y_{ij} - \mu_i)^2}{2\sigma^2} \right] \right\} \exp \left[-\frac{(\mu_i - \xi)^2}{2\tau^2} \right] \\ &= \prod_{i=1}^k f(\mu_i | x, w) \end{aligned}$$

where each of the $f(\mu_i | x, w)$ comes from the analysis of a single normal sample given in Lecture 1. That is conditional on w , the μ_i 's are independent $\mathcal{N}(m_i^*, v_i^*)$, where:

$$(5.16) \quad v_i^* = (n_i\sigma^{-2} + \tau^{-2})^{-1}$$

$$(5.17) \quad m_i^* = (n_i\sigma^{-2} + \tau^{-2})^{-1}(n_i\sigma^{-2}\bar{y}_i + \tau^{-2}\xi)$$

We can already see the shrinkage in the model, because the posterior mean of each μ_i is a weighted average of its own data estimate \bar{y}_i and the common value ξ . So they are shrunk towards this common value.

Example 5.3.1. Let's consider a simple regression situation in which we have observations y_1, y_2, \dots, y_n at values $x_1 < x_2 < \dots < x_n$ of the explanatory variables. The usual linear regression model specifies

$$(5.18) \quad y_i | \alpha, \beta, \sigma^2 \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2)$$

Instead we can create the hierarchical model:

- Data model:
- First level of prior:
- Second level of prior:

LECTURE 6

Sufficiency and Ancillary

What happens if $f(x | \theta)$ does not depend on θ ? If $f(x | \theta)$ does not depend on θ , the data are completely uninformative.

Definition 6.1. $t(x)$ is sufficient if, for any given θ , $f(x | \theta)$ is a function only of $t(x)$, apart from a multiplicative factor that can be any function of x

It means that we only need to know $t(x)$ in order to obtain the posterior distribution. It is sufficient to know $t(x)$. Therefore the posterior distribution is the same as if we only had observed $t(x)$ rather than the whole of x .

Suppose that $s(x)$ represents all other information in x , $x = (t(x), s(x))$. Then

$$(6.2) \quad f(x | \theta) = f(t(x) | \theta) f(s(x) | t(x), \theta).$$

$t(x)$ is sufficient $f(s(x) | t(x), \theta)$ must not depend on θ . So once we know $t(x)$ there is no information in $s | x$.

Definition 6.3. $s(x)$ is ancillary if $f(s(x) | \theta)$ does not depend on θ .

Example 6.0.1. Let $x_i | \sim \mathcal{N}(\mu, \sigma^2)$ with σ^2 known. Then $x_i - x_j \sim \mathcal{N}(0, 2\sigma^2)$ is ancillary for any $i \neq j$.

- Let $t(x)$ be the rest of the information in x , so that again we have $x = (t(x), s(x))$. Then

$$(6.4) \quad f(x | \theta) = f(s(x) | \theta) f(t(x) | s(x), \theta),$$

but now this implies that $f(x | \theta) \propto f(t(x) | s(x), \theta)$

6.1. The Likelihood Principle

The likelihood principle asserts that inference should be based only on the likelihood.

Example 6.1.1. Let $x \sim Be(1, \theta)$ and the results of a fixed number n of independent Bernoulli trials.

$$(6.5) \quad f_x(x | \theta) = \theta^r (1 - \theta)^{n-r}$$

where r is the number of success. If $x \sim Bi(n, \theta)$, then $f_R(r | \theta) = \binom{n}{r} \theta^r (1 - \theta)^{n-r} \propto \theta^r (1 - \theta)^{n-r}$.

If we keep observing independent Bernoulli trials until we get a fixed number r of success.

$$(6.6) \quad f_y(y | \theta) = \theta^r (1 - \theta)^{n-r}$$

where r -th success on the n -th trial. If the distribution is negative binomial

$$(6.7) \quad f_N(n) = \binom{n-1}{r-1} \theta^r (1 - \theta)^{n-r} \propto \theta^r (1 - \theta)^{n-r}$$

6.2. Identifiability

Let $\theta = (y(\theta), h(\theta))$ and $f(x | \theta)$ depend only on $g(\theta)$.

$$(6.8) \quad \begin{aligned} f(\theta | x) &\propto f(x | \theta) f(\theta) = f(x | g(\theta)) f(\theta) \\ &= f(x | g(\theta)) f(g(\theta)) f(h(\theta) | g(\theta)) \\ &\propto f(g(\theta) | x) f(h(\theta) | g(\theta)). \end{aligned}$$

This says that the posterior distribution of θ is made up of the distribution of $g(\theta)$ and the prior distribution of $h(\theta)$ given $g(\theta)$. So it is the conditional posterior distribution of $h(\theta)$ given $g(\theta)$ that is the same as the prior. That is

$$(6.9) \quad f(h(\theta) | x, g(\theta)) = f(h(\theta) | g(\theta)).$$

We say that $h(\theta)$ is not identifiable from these data. No matter how much data we get, we can not learn exactly what $h(\theta)$ is. With sufficient data we can learn $g(\theta)$, but not $h(\theta)$.

6.3. Asymptotic Theory

Suppose that we have a sequence of iid observations x_1, x_2, x_3, \dots and suppose that $\mathbf{x}_n = (x_1, x_2, \dots, x_n)$ comprises the first n observations.

We can now consider a sequence of posterior distributions $f(\theta | x_1), f(\theta | x_2), f(\theta | x_3), \dots, f(\theta | x_n), f(\theta | x_{n+1})$. We wish to know how the posterior distribution $f(\theta | x_n)$ behaves as $n \rightarrow \infty$.

So as we get more data, we expect that the posterior will in some sense converge to the true value. Also, the weight the posterior gives to the data increases, and therefore we can expect that in the limit the posterior will be insensitive to the prior.

Subject to some regularity conditions. Regularity conditions:

- (a) The whole θ needs to be identifiable.
- (b) Prior possibility condition: the prior probability does not give zero probability to the true value of θ .
- (c) Continuity condition: we need a continuity condition for θ .

6.4. Preposterior Properties

Let X and Y be any two random variables. Then

$$(6.10) \quad E(Y) = E\{E(Y | X)\},$$

$$(6.11) \quad v(Y) = E\{v(Y | X)\} + v\{E(Y | X)\}$$

Let's replace Y by the parameter vector θ and X by the data vector X .

Remark. If we use the posterior mean $E(\theta | x)$ to estimate θ , the its expected bias 0.

6.5. Conjugate Prior Forms

The conjugacy is a joint property of the prior and the likelihood function that provides a posterior from the same distributional family as the prior.

Example 6.5.1. Conjugacy in exponential specifications.

$$(6.12) \quad E(x | \theta) = \theta \exp\{-\theta x\}$$

where $0 \leq x$, $0 < \theta$. If $\theta \sim \text{Gamma}(\alpha, \beta)$, then

$$(6.13) \quad f(\theta | \alpha, \beta) = \frac{1}{\Gamma(\alpha)} \beta^\alpha \theta^{\alpha-1} \exp\{-\beta\theta\}$$

where $\theta, \alpha, \beta > 0$

Suppose we now observe $x_1, x_2, \dots, x_n \sim iid$. The likelihood is

$$(6.14) \quad L(\theta | x) = \prod_{i=1}^n \theta e^{-\theta x_i} = \theta^n \exp\left\{-\theta \sum x_i\right\}$$

Thus,

$$(6.15) \quad \begin{aligned} \pi(\theta | x) &\propto E(x | \theta)L(\theta | x) \\ &= \theta^n \exp\left\{-\theta \sum x_i\right\} \frac{1}{\Gamma(\alpha)} \beta^\alpha \theta^{\alpha-1} \exp\{-\beta\theta\} \\ &\propto \theta^{\alpha+n-1} \exp\left\{-\theta(\sum x_i + \beta)\right\} \end{aligned}$$

This is the kernel of a $Gamma(\alpha + n, \sum x_i + \beta)$ and therefore the gamma distribution is shown to be conjugate to the exponential likelihood function.

Table 1. Conjugate Prior Distribution Table

Likelihood Form	Conjugate Prior Distribution	Hyperparameters
Bernoulli	Beta	$\alpha > 0, \beta > 0$
Binomial	Beta	$\alpha > 0, \beta > 0$
Multinomial	Dirichlet	$\theta_j > 0, \sum \theta_j = \theta_0$
Negative Binomial	Beta	$\alpha > 0, \beta > 0$
Poisson	Gamma	$\alpha > 0, \beta > 0$
Exponential	Gamma	$\alpha > 0, \beta > 0$
Gamma (ind χ^2)	Gamma	$\alpha > 0, \beta > 0$
Normal for μ	Normal	$\mu \in \mathbb{R}, \sigma^2 > 0$
Normal for σ^2	Inverse Gamma	$\alpha > 0, \beta > 0$
Pareto for α	Gamma	$\alpha > 0, \beta > 0$
Pareto for β	Pareto	$\alpha > 0, \beta > 0$
Uniform	Pareto	$\alpha > 0, \beta > 0$

LECTURE 7

Tackling Real Problems

There are various computational tools that are widely used in practical Bayesian statistics, the most well-known one is Markov Chain Monte Carlo or MCMC. The basic idea is that we randomly draw a very large sample $Q^{(1)}, Q^{(2)}, \dots$ from the posterior distribution. Given such a sample, we can compute any inference we wish. If we want to make inference about some derived parameter $\phi = g(\theta)$ then $\phi^{(1)} = g(Q^{(1)})$, $\phi^{(2)} = g(Q^{(2)}), \dots$, is a sample from its posterior distribution $f(\theta | x)$. The sample mean $\bar{\phi}$ is an estimate of $E(\phi | x)$. In principle, we could draw such a sample using simple Monte Carlo sampling. That is each $Q^{(j)}$ is independently drawn from $f(\theta | x)$. There are algorithms for efficiently drawing random samples from a wide variety of standard distributions.

7.1. What is a Markov Chain?

”A stochastic process” is a consecutive set of random quantities defined on some known state space Q , indexed so that the order is known. $\{Q^{[t]}, t \in T\}$. Here the state space (which is parameter space for us) is just the allowable range of values for the random vector of interest. The state space Q is either discrete or continuous depending on how the variable of interest is measured.

A Markov chain is a stochastic process with the property that at time t in the series, the probability of making a transition to any new state is dependent only on the current state of the process.

$$(7.1) \quad p(Q^{[t]} \in A | Q^{[0]}, Q^{[1]}, \dots, Q^{[t-2]}, Q^{[t-1]})$$

where A is the identified set on the complete state space.

A fundamental concern is the transition process that defines the probabilities of moving to other points in the state space, given the current location of the chain. This structure is defined via the transition kernel K as a general mechanism for describing the probabilities of moving to some other specified state based on the current chain status. When the state space is discrete, then K is a matrix, $k \times k$ for k discrete elements in A , where each cell defines the probability of a state transition from the first term in the parentheses to all possible states:

$$(7.2) \quad P_A = \begin{bmatrix} p(\theta_1, \theta_1) & \dots & p(\theta_1, \theta_k) \\ \vdots & & \\ p(\theta_k, \theta_1) & \dots & p(\theta_k, \theta_k) \end{bmatrix}$$

The row and columns indicate:

•
•

An important feature of the transition kernel is that the transition probabilities between two selected states for arbitrary numbers of steps m can be calculate multiplicative:

$$(7.3) \quad p^m(\theta_j^{[m]} = y \mid \theta_i^{[0]} = x) = \underbrace{\sum_{\theta_1} \sum_{\theta_2} \dots \sum_{\theta_{m-1}}}_{\text{all possible paths}} \underbrace{p(\theta_i, \theta_1)p(\theta_1, \theta_2)\dots p(\theta_{m-2}, \theta_{m-1})}_{\text{transition products}}$$

So $p^m(\theta_j^{[m]} = y \mid \theta_i^{[0]} = x)$ is also a stochastic transition matrix.

Example 7.1.1.

$$(7.4) \quad \underbrace{P}_{\text{current period}} = \begin{cases} \theta_1 & \begin{bmatrix} 0.8 & 0.2 \end{bmatrix} \\ \theta_2 & \begin{bmatrix} 0.6 & 0.4 \end{bmatrix} \end{cases}$$

Let the starting point $S_0 = [0, 5 \quad 0, 5]$, to get the first state

$$(7.5) \quad S_1 =$$

To get the second state,

$$(7.6) \quad S_2 =$$

$$(7.7) \quad S_3 =$$

$$(7.8) \quad S_4 =$$

So the choice proportions are converging to $[0, 75 \quad 0, 25]$ since the transition matrix is pushing toward a steady state or stationary distribution of the proportions. So when we reach this distribution, all future states are constant, that is stationary.

7.2. The Chapman-Kolmogorov Equations

These equations specify how successive events are bound together probabilistically. If we abbreviate the hand side of expression 7.3

$$(7.9) \quad p^{m_1+m_2}(x, y) = \sum_{\text{all } z} p^{m_1}(x, z)p^{m_2}(z, y) \quad \text{discrete case.}$$

$$(7.10) \quad p^{m_1+m_2}(x, y) = \int_{\text{range } z} p^{m_1}(x, z)p^{m_2}(z, y)dz \quad \text{continuous case.}$$

This is also equal to:

$$(7.11) \quad p^{m_1+m_2} = p^{m_1}p^{m_2} = p^{m_1}p^{m_2-1}p = p^{m_1}p^{m_2-2}p^2 = \dots$$

For discrete case. Thus iterative probabilities can be decomposed into segmented products in any way we like, depending on the interim step.

7.3. Marginal Distributions

The marginal distributions at some step m -th from the transition kernel is found via

$$(7.12) \quad \underbrace{\pi^m(\theta)}_{\text{current value of the chain}} = \underbrace{[p^m(\theta_1), p^m(\theta_2), \dots, p^m(\theta_k)]}_{\text{the row of the transition kernel for the } m\text{-th step}}$$

So the marginal distribution at the first step of the Markov chain is given by

$$(7.13) \quad \boxed{\pi^1(\theta) = \pi^0(\theta)p^1}.$$

where π^0 = initial starting value assigned to the chain, $p^1 = P$ = simple transition matrix, Thus

$$(7.14) \quad \pi^n =$$

$$(7.15) \quad \pi^m(\theta_j) =$$

7.4. General Properties of Markov Chains

- (a) Homogeneity: A homogeneous Markov chain at step n -th transition probability that does not depend on the value of so the decision to move at this step is independent of this being the current point in time.
- (b) Irreducibility: A Markov chain is irreducible if every point or collection of points can be reached from every other point or collection of points. So irreducibility implies the existence of a path between any two points in the subspace

(c) Recurrence: If a subspace is closed, finite, and irreducible, then all states within the subspace are recurrent. An irreducible Markov chain is called recurrent with regard to a given state, A , which is a single point or a defined collection of points, if the probability that the chain occupies A infinitely often over unbounded time is non-zero.

(d) Stationarity: Let $\pi(\theta)$ = stationary distribution of the Markov chain for θ of the state space.

$p(\theta_i, \theta_j)$ = the probability that the chain will move from θ_j to θ_i at some arbitrary step t from the transition kernel.

$\pi^t(\theta)$ = the marginal distribution, thus the stationary distribution is defined as:

$$(7.16) \quad \sum_{\theta_i} \pi^t(\theta_i) p(\theta_i, \theta_j) = \pi^{t+1}(\theta_j) \quad \text{discrete case}$$

$$(7.17) \quad \int \pi^t(\theta_i) p(\theta_i, \theta_j) d\theta_i = \pi^{t+1}(\theta_j) \quad \text{continuous case}$$

That is $\pi = \pi p$. The marginal distribution remains fixed when the chain reaches the stationary distribution.

Once the chain reaches its stationary distribution (also its invariant distribution, equilibrium distribution or limiting distribution), it stays in this distribution and moves about or "mixes" throughout the subspace according to marginal distribution, $\pi(\theta)$, forever.

(e) Ergodicity: If the chain is irreducible, positive Harris recurrent (i.e. recurrence under unbounded continuous state space), and aperiodic, then we call it ergodic. Ergodic Markov chains have the property:

$$(7.18) \quad \lim_{n \rightarrow \infty} p^n(\theta_i, \theta_j) = \pi(\theta_j)$$

for all θ_i , and θ_j in the subspace. Therefore in the limit, the marginal distribution at one step is identical to the marginal distributions at all other steps. The ergodic theorem is analogous to the strong law of large numbers but for Markov chains. Thus suppose that $\theta_{i+1}, \dots, \theta_{i+n}$ are n values from a Markov chain that has reached its ergodic distribution. A statistic of interest, $h(\theta)$, can be calculated empirically:

$$(7.19) \quad h(\theta_i) = \frac{1}{n} \sum_{j=i+1}^{i+n} h(\theta_j) \approx h(\theta)$$

For a given empirical estimator $\hat{h}(\theta_i)$ with bounded limiting variance, we get the central limit theorem results

$$(7.20) \quad \sqrt{n} \frac{\hat{h}(\theta_i) - h(\theta)}{\sqrt{v(\hat{h}(\theta_i))}} \rightarrow_{n \rightarrow \infty} \mathcal{N}(0, 1).$$

7.5. Noniterative Monte Carlo Methods

(a) Direct Methods:

The most basic definition of Monte Carlo integration is that $\theta \sim h(\theta)$ and we seek $\gamma = E[f(\theta)] = \int f(\theta)h(\theta)d\theta$. Then if $\theta_1, \dots, \theta_N \sim^{iid} h(\theta)$, we have

$$(7.21) \quad \hat{\gamma} = \frac{1}{N} \sum_{j=1}^N f(\theta_j)$$

with converges to $E[f(\theta)]$ with probability 1 as $N \rightarrow \infty$ by the strong law of large numbers. In our case, $h(\theta)$ is a posterior distribution and γ is the posterior mean of $f(\theta)$. Hence the computation of posterior expectations require only a sample of size N from the posterior distribution.

(b) Indirect methods:

If we can't directly sample from distribution? In this case, we can use one of the following methods

- Importance sampling: Suppose we wish to approximate a posterior expectation, say

$$(7.22) \quad E(f(\theta) | y) = \frac{\int f(\theta)L(\theta)\pi(\theta)d\theta}{\int L(\theta)\pi(\theta)d\theta}$$

where f is the function of interest, L is the likelihood or the data y . By defining the weight function $w(\theta) = \frac{L(\theta)\pi(\theta)}{g(\theta)}$, we have

$$(7.23) \quad \begin{aligned} E(f(\theta) | y) &= \frac{\int f(\theta)w(\theta)g(\theta)d\theta}{\int w(\theta)g(\theta)d\theta} \\ &\approx \frac{\frac{1}{N} \sum_{j=1}^N f(\theta_j)w(\theta_j)}{\frac{1}{N} \sum_{j=1}^N w(\theta_j)} \end{aligned}$$

where $\theta_j \sim^{iid} g(\theta)$. Here $g(\theta)$ is called the importance function, how closely it resembles $L(\theta)\pi(\theta)$ controls how good the approximation in the equation is. If $g(\theta)$ is a good approximation, the weights will all be roughly equal,

which in turn will minimize the variance of the numerator and denominator.

If $g(\theta)$ is a poor approximation, many of the weights will be close to zero, and thus a few θ_j 's will dominate the sums, producing an inaccurate approximation.

Thus in importance sampling, one chooses a known density function $g(\theta)$ that is easy to sample. The procedure works best if $g(\theta)$ is similar in shape to the known kernel of the posterior $L(\theta)\pi(\theta)$ with tails that do not decay more rapidly than the tails of the posterior.

- Rejection sampling: Here instead of trying to approximate the normalized posterior:

$$(7.24) \quad h(\theta) = \frac{L(\theta)\pi(\theta)}{\int L(\theta)\pi(\theta)d\theta},$$

we try to "blanket" it. That is suppose there exists an identifiable constant $\mu > 0$ and a smooth density $g(\theta)$, called the envelope function, such that $L(\theta)\pi(\theta) < \mu g(\theta)$ for all θ .

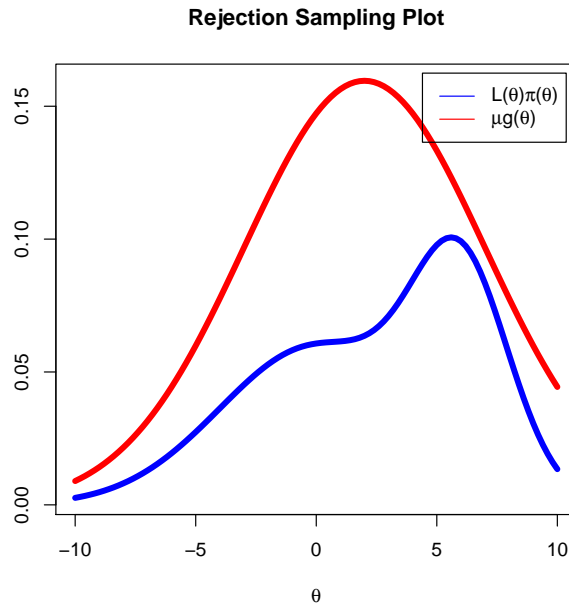


Figure 1. Plot of $L(\theta)\pi(\theta)$ and $\mu g(\theta)$

Example 7.5.1.

The rejection method proceeds as follows:

- (a) Generate $\theta_j \sim g(\theta)$
 - (b) Generate $U \sim Uniform(0, 1)$
 - (c) If $\mu u g(\theta_i) < L(\theta_j)\pi(\theta_j)$, accept θ_j , otherwise reject θ_j
 - (d) Return to step (i) and repeat until the desired sample $\{\theta_j, j = 1, \dots, N\}$ is obtained. The members of this sample will then be random samples from $h(\theta)$. Like an importance sampling density, the envelope density g should be similar to the posterior in general appearance, but with heavier tails and sharper infinite peaks, in order to assure that there are sufficiently many rejection candidates available across its entire domain. Also μg is actually an "envelope" for the unnormalized posterior $L\pi$.
- Weighted bootstrap: Suppose an μ appropriate for the rejection method is not readily available, but that we do have a sample $\theta_1, \dots, \theta_N$ from same approximating density $g(\theta)$, Define:

$$(7.25) \quad w_i = \frac{L(\theta_i)\pi(\theta_i)}{g(\theta_i)}$$

and

$$(7.26) \quad q_i = \frac{w_i}{\sum_{j=1}^N w_j}$$

Now draw θ^* from the discrete distribution over $\{\theta_1, \dots, \theta_N\}$ which places mass at θ_i . Then

$$(7.27) \quad \theta^* \sim h(\theta) = \frac{L(\theta)\pi(\theta)}{\int L(\theta)\pi(\theta)d\theta}$$

with the approximation improving as $N \rightarrow \infty$. This is a weighted bootstrap, since instead of resampling from the set $\{\theta_1, \dots, \theta_N\}$ with equally likely probabilities of selection, we are resampling some points more often than others due to the unequal weighting.

NOTE: Importance sampling, rejection sampling and the weighted bootstrap are all "one-off", or non-iterative methods: they draw a sample of size N , and stop. Hence there is no notion of the algorithm "converging" we simply require N sufficiently large. But for many problems, especially high-dimensional ones, it may be difficult or impossible to find an importance sampling density (or envelope function)

which is an acceptably accurate approximation to the lay posterior, but still easy to sample from.

Solution: MCMC based methods

- Gibb's sampling
- Metropolis Algorithm
- Metropolis-Hastling Algorithms

LECTURE 8

- There is a superficial resemblance between MCMC and the frequentist technique of bootstrapping.

An MCMC sampler is a Markov chain in θ space. We start the chain with an arbitrary $\theta^{(1)}$, which is therefore not random. Each subsequent $\theta^{(j)}$ drawn from a distribution $q(\theta^{(j)} | \theta^{(j-1)})$, so that it depends only on the previous point $\theta^{(j-1)}$ and not on the history of the chain up to that point. This is the Markov property. The distribution $q(\cdot | \cdot)$ defines the transition probabilities of the chain.

- Let $\theta^{(2)}$ be the distribution given the initial $\theta^{(1)}$ so, just $q(\theta^{(2)} | \theta^{(1)})$. The distribution of $\theta^{(3)}$ is:

$$(8.1) \quad f(\theta^{(3)} | \theta^{(1)}) = \int q(\theta^{(3)} | \theta^{(2)})q(\theta^{(2)} | \theta^{(1)})d\theta^{(2)}.$$

The distribution of $\theta^{(j)}$ can then be obtained in principle by iterating this convolution using

$$(8.2) \quad f(\theta^{(j)} | \theta^{(1)}) = \int q(\theta^{(j)} | \theta^{(j-1)})q(\theta^{(j-1)} | \theta^{(1)})d\theta^{(j-1)}.$$

The Markov chain theory then says that, subject to some conditions, there is a unique limiting distribution $p(\theta)$ such that for all sufficiently large j we have $f(\theta^{(j)} | \theta^{(1)}) \approx p(\theta^{(j)})$, and this limiting distribution is independent of the arbitrary starting value $\theta^{(1)}$.

8.1. The Gibbs Sampler

Procedure:

- (a) Choose starting values: $\theta^{[0]} = [\theta_1^{[0]}, \theta_2^{[0]}, \dots, \theta_k^{[0]}]$

(b) At the j -th starting at $j = 1$ complete the single cycle by drawing values from the k distributions given by:

$$\begin{aligned}\theta_1^{[j]} &\approx \pi(\theta_1 \mid \theta_2^{[j-1]}, \theta_3^{[j-1]}, \theta_4^{[j-1]}, \dots, \theta_{k-1}^{[j-1]}, \theta_k^{[j-1]}) \\ \theta_2^{[j]} &\approx \pi(\theta_2 \mid \theta_1^{[j]}, \theta_3^{[j-1]}, \theta_4^{[j-1]}, \dots, \theta_{k-1}^{[j-1]}, \theta_k^{[j-1]}) \\ \theta_3^{[j]} &\approx \pi(\theta_3 \mid \theta_1^{[j]}, \theta_2^{[j]}, \theta_4^{[j-1]}, \dots, \theta_{k-1}^{[j-1]}, \theta_k^{[j-1]}) \\ &\vdots \\ &\vdots \\ \theta_{k-1}^{[j]} &\approx \pi(\theta_{k-1} \mid \theta_1^{[j]}, \theta_2^{[j]}, \theta_3^{[j]}, \dots, \theta_{k-2}^{[j]}, \theta_k^{[j-1]}) \\ \theta_k^{[j]} &\approx \pi(\theta_k \mid \theta_1^{[j]}, \theta_2^{[j]}, \theta_3^{[j]}, \dots, \theta_{k-2}^{[j]}, \theta_{k-1}^{[j]})\end{aligned}$$

(c) Increment j and repeat until convergence.

Example 8.1.1. $\theta_1 \mid \theta_2 \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$

$$(8.3) \quad \theta_1^{[j]} \mid \theta_2^{[j-1]} \sim \mathcal{N}(\rho\theta_2^{[j-1]}, 1 - \rho^2)$$

$$(8.4) \quad \theta_2^{[j]} \mid \theta_1^{[j]} \sim \mathcal{N}(\rho\theta_1^{[j]}, 1 - \rho^2)$$

Example 8.1.2. Let x_1, x_2, \dots, x_n be a series of count data where there exists the possibility of a changepoint at some period k , along the series. Therefore there are two Poisson data-generating processes:

$$(8.5) \quad x_i \mid \lambda \sim P(\lambda) \quad i = 1, 2, \dots, k$$

$$(8.6) \quad x_i \mid \phi \sim P(\phi) \quad i = k + 1, \dots, n$$

The parameters to be estimated are λ , ϕ and k . Also the three independent priors applied to this model are:

$$\lambda \sim \text{Gamma}(\alpha, \beta),$$

$$\phi \sim \text{Gamma}(\gamma, \delta),$$

$$k \sim \text{Discrete uniform on } [1, 2, \dots, n].$$

So the joint posterior is

$$\begin{aligned}(8.7) \quad &\pi(\lambda, \phi, k \mid y) \propto L(\lambda, \phi, k \mid y)\pi(\lambda \mid \alpha, \beta)\pi(\phi \mid \gamma, \delta)\pi(k) \\ &= \left(\prod_{i=1}^k \frac{e^{-\lambda}\lambda^{y_i}}{y_i!}\right) \left(\prod_{i=k+1}^n \frac{e^{-\phi}\phi^{y_i}}{y_i!}\right) \left(\frac{\beta^\alpha}{\Gamma(\alpha)}\lambda^{\alpha-1}e^{-\beta\lambda}\right) \left(\frac{\delta^\gamma}{\Gamma(\gamma)}\phi^{\gamma-1}e^{-\delta\phi}\right) \frac{1}{n} \\ &\propto \lambda^{\alpha-1+\sum_{i=1}^k y_i} \phi^{\gamma-1+\sum_{i=k+1}^n y_i} e^{-(k+\beta)\lambda - (n-k+\delta)\phi}.\end{aligned}$$

So

$$\lambda \mid \phi, k \sim \text{Gamma}(\alpha + \sum_{i=1}^k y_i, \beta + k),$$

$$\phi \mid \lambda, k \sim \text{Gamma}(\gamma + \sum_{i=k+1}^n y_i, \delta + n - k).$$

Let λ and ϕ be fixed,

$$\begin{aligned} p(y \mid k, \lambda, \phi) &= \left(\prod_{i=1}^k \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} \right) \left(\prod_{i=k+1}^n \frac{e^{-\phi} \phi^{y_i}}{y_i!} \right) \\ (8.8) \quad &= \left(\prod_{i=1}^k \frac{1}{y_i!} \right) e^{k(\phi-\lambda)} e^{-n\phi} \lambda^{\sum_{i=1}^k y_i} \left(\prod_{i=k+1}^n \phi^{y_i} \right) \\ &= \left(\prod_{i=k+1}^n \frac{e^{-\phi} \phi^{y_i}}{y_i!} \right) \left(e^{k(\phi-\lambda)} \left(\frac{\lambda}{\phi} \right)^{\sum_{i=1}^k y_i} \right) \\ &= f(y, \phi) L(y \mid k, \lambda, \phi). \end{aligned}$$

Listing 8.1. MCMC Function in R

```

1 bcp<- function(theta.matrix, y, a, b, g, d) {
2
3 n<- length(y);
4 k.prob <- rep(0, length=n);
5
6 for (i in 1:nrow(theta.matrix)) {
7
8 lambda <- rgamma(1, a+sum(y[theta.matrix[(i-1), 3]:n])
9 , b+theta.matrix[(i-1), 3]);
10
11 phi <- rgamma(1, g+sum(y[theta.matrix[(i-1), 3]:n])
12 , d+length(y)-theta.matrix[(i-1), 3]);
13
14 for(j in 1:n) k.prob[j] <- ...
15     exp(j*(phi-lambda)) * (lambda/phi) ^ sum(y[1:j]);
16
17 k.prob <- k.prob/sum(k.prob);
18 k <- sample(1:n, , size=1, prob=k.prob);
19 theta.matrix[i, ] <-c(lambda, phi, k);
20 }
21 return(theta.matrix)
22 }
```

Example 8.1.3. The time series stored in the file gives the # of British coal mining disasters per year over the period 1851 – 1962.

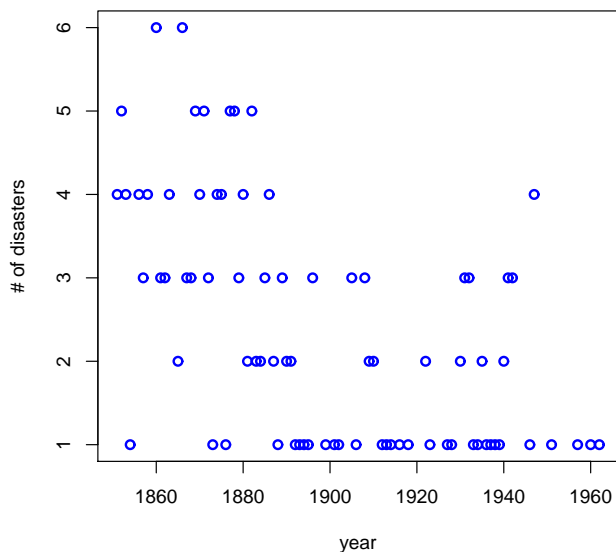


Figure 1. Plot of # of disasters over 1851 – 1962

There has been a reduction in the rate of disasters over the period. Let y_i denote the # of disasters in year $i = 1, \dots, 12$ (relabelling the years by numbers from 1 to $n = 112$). A model that has been proposed in the literature has the form:

$$\begin{aligned} y_i &\sim \text{Poisson}(\theta), & i = 1, \dots, k \\ y_i &\sim \text{Poisson}(\lambda), & i = k + 1, \dots, n \end{aligned}$$

Let

$$\begin{aligned} \theta &\sim \text{Gamma}(a_1, b_1) \\ \lambda &\sim \text{Gamma}(a_2, b_2) \\ k &\sim \text{discrete uniform over } \{1, \dots, n\} \\ b_1 &\sim \text{Gamma}(c_1, d_1) \\ b_2 &\sim \text{Gamma}(c_2, d_2) \end{aligned}$$

So

$$\begin{aligned} \pi(\theta, \lambda, k, b_1, b_2 \mid y) &\propto e^{-\theta k} \theta^{\sum_{i=1}^k y_i} e^{-\theta(n-k)} \lambda^{\sum_{i=k+1}^n y_i} b_1^{a_1} \theta^{a_1-1} e^{-b_1 \theta} \\ &\quad \times b_1^{c_1-1} e^{-d_1 b_1} b_2^{a_2} \lambda^{a_2-1} e^{-b_2 \lambda} b_2^{c_2-1} e^{-d_2 b_2} I[k \in \{1, 2, \dots, n\}]. \end{aligned}$$

$$\theta \mid y, \lambda, b_1, b_2, k \sim \text{Gamma}(a_1 + \sum_{i=1}^k y_i, b_1 + k)$$

$$\lambda \mid y, \theta, b_1, b_2, k \sim \text{Gamma}(a_2 + \sum_{i=k+1}^n y_i, b_2 + n - k)$$

$$b_1 \mid y, \theta, \lambda, b_2, k \sim \text{Gamma}(c_1 + a_1, d_1 + \theta)$$

$$b_2 \mid y, \theta, \lambda, b_1, k \sim \text{Gamma}(c_2 + a_2, d_2 + \lambda)$$

and

$$p(k \mid y, \theta_1, \lambda, b_1, b_2) = \frac{e^{(\lambda-\theta)k} (\theta/\lambda)^{\sum_{i=1}^k y_i} I[k \in \{1, 2, \dots, n\}]}{\sum_{j=1}^n \left\{ e^{(\lambda-\phi)j} \left(\frac{\theta}{\lambda}\right)^{\sum_{i=1}^k y_i} \right\}}$$

Burn-in period: The observations obtained after the chain has settled down to the posterior will be more useful in estimating probabilities and expectations for $p(\theta)$. If we throw out the early observations, taken while the process was settling down, the remainder of the process should be a very close approximation to one in which every observation is sampled from the posterior. Dropping the early observations is referred to as using a burn-in period.

Thinning is a process used to make the observations more nearly independent, hence more nearly random sample from the posterior distribution. Frankly, after a burn-in, there is not much point in thinnings unless the correlations are extremely large. If there is a lot of correlation between adjacent observations, a larger overall MC sample size is needed to achieve reasonable numerical accuracy, in addition to needing a much larger burn-in.

LECTURE 9

Summary of the properties of Gibbs sampler

- (a) Since the Gibbs sampler conditions on values from the last iteration of its chain values, it clearly constitutes a Markov chain.
- (b) The Gibbs sampler has the true posterior distribution of parameter vector as its limiting distribution.
- (c) The Gibbs sampler is a homogeneous Markov chain, the consecutive probabilities are independent of n , the current length of the chain.
- (d) The Gibbs sampler converges at a geometric rate: the total variation distance between an arbitrary time and the point of convergence decreases at a geometric rate in time (t).
- (e) The Gibbs sampler is an ergodic Markov chain.

9.1. Metropolis Algorithm

The Metropolis algorithm[3] is another type of accept-reject algorithm. It requires a candidate generating distribution, sometimes referred to as the proposal distribution. The algorithm begins with an initial value θ^1 . At the k -th iteration we have $(\theta^1, \theta^2, \dots, \theta^k)$. The $(k+1)$ st iteration first generates θ^* from a proposal density $h(\theta^* | \theta^k)$. This density should mimic the actual posterior distribution in some sense, but in theory, it can be any distribution with the same support as the posterior. Define:

$$(9.1) \quad \alpha(\theta^*, \theta^k) = \min \left\{ 1, \frac{p(\theta^*)h(\theta^k | \theta^*)}{p(\theta^k)h(\theta^* | \theta^k)} \right\} \cong \alpha.$$

We then simulate $u \sim u[0, 1]$ and we select $\theta^{k+1} = \theta^*$ if $u \leq \alpha$ and otherwise take $\theta^{k+1} = \theta^k$. Thus

$$(9.2) \quad \theta^{k+1} = \begin{cases} \theta^* & \text{with probability } \alpha(\theta^*, \theta^k) \\ \theta^k & \text{with probability } 1 - \alpha(\theta^*, \theta^k) \end{cases}$$

Here α only uses the ratio of two values of $p(\cdot)$, so it is enough to know the kernel of the posterior density. The acceptance ration can be also written as:

$$(9.3) \quad \alpha = \frac{p(\theta^* | y)/h_t(\theta^* | \theta^{t-1})}{p(\theta^{t-1} | y)/h_t(\theta^{t-1} | \theta^*)}.$$

The ratio α is always defined, because a jump from θ^{t-1} to θ^* can only occur if both $p(\theta^{t-1} | y)$ and $h_t(\theta^* | \theta^{t-1})$ are nonzero.

So here the proposal density is asymmetric.

9.2. Metropolis Algorithm

The original Metropolis algorithm assumes that $h(\theta^k | \theta^*) = h(\theta^* | \theta^k)$ so that $\alpha(\theta^*, \theta^k) = \min \{p(\theta^*)/p(\theta^k)\}$. This is called the random walk.

Various suggestions are made about how to choose $h(\theta^* | \theta^k)$. Often it is taken as a $\mathcal{N}(\theta^k, \Sigma^k)$ distribution with various suggestions for Σ^k .

The Metropolis algorithm is an adaptation of a random walk that uses an acceptance/rejection rule to converge to the specified target distribution. The algorithm proceeds as follows:

- (a) Draw a starting θ^0 , for which $p(\theta^0 | y) > 0$ from a starting distribution $p_0(\theta)$. The starting distribution might, for example, be based on an approximation or we may simply choose starting values dispersed around a crude approximate estimate.
- (b) For $t = 1, 2, \dots$
 - Sample a proposal θ^* from a jumping distribution (or proposal distribution) at time t , $R_t(\theta^* | \theta^{t-1})$. For the Metropolis algorithm, the jumping distribution must be symmetric, satisfying the condition $h_t(\theta_a | \theta_b) = h_t(\theta_b | \theta_a)$ for all θ_a, θ_b and t .
 - Calculate the ratio of the densities,

$$(9.4) \quad \alpha = \frac{p(\theta^* | y)}{p(\theta^{t-1} | y)}$$

- Set

$$(9.5) \quad \theta^t = \begin{cases} \theta^* & \text{with probability } \min(\alpha, 1) \\ \theta^{t-1} & \text{otherwise} \end{cases}$$

LECTURE 9. SUMMARY OF THE PROPERTIES OF GIBBS SAMPLER **63**

Given the current value θ^{t-1} , the transition distribution $h_t(\theta^t | \theta^{t-1})$ of the Markov chain is thus a mixture of a point mass at $\theta^t = \theta^{t-1}$, and a weighted version of the jumping distribution, $h_t(\theta^t | \theta^{t-1})$, that adjusts for the acceptance rate. The algorithm requires the ability to calculate the ratio α for all (θ, θ^*) , and to draw θ from the jumping distribution $h_t(\theta^* | \theta)$ for all θ and t . In addition, step (c) requires the generation of a uniform random number.

If $\theta^t = \theta^{t-1}$, that is the jump is not accepted, this counts as an iteration in the algorithm.

Interpretation of the Gibbs sampler as a special case of the Metropolis-Hastling algorithm:

Gibbs sampling can be viewed as a special case of the Metropolis-Hastling's algorithm in the following way. We first define iteration t to consist of a series of d steps, with step j of iteration t corresponding to an update of the subvector θ_j conditional on all other elements of θ . Then the jumping distribution $h_{j,t}(\cdot | \cdot)$ at step j of iteration t only jumps along the j -th subvector, and does so with the conditional posterior density of θ_j density θ_{-j}^{t-1} :

$$(9.6) \quad h_{j,t}^{Gibbs}(\theta^* | \theta^{t-1}) = \begin{cases} p(\theta_j^* | \theta_{-j}^{t-1}, y) & \text{if } \theta_{-j}^* = \theta_{-j}^{t-1} \\ 0 & \text{otherwise} \end{cases}$$

The only possible jumps are to parameter vectors that match θ^{t-1} on all components other than j -th. Under this jumping distribution, the ratio at the j -th step of iteration t is:

$$(9.7) \quad \begin{aligned} \alpha &= \frac{p(\theta^* | y) / h_{j,t}^{Gibbs}(\theta^* | \theta^{t-1})}{p(\theta^{t-1} | y) / h_{j,t}^{Gibbs}(\theta^{t-1} | \theta^*)} \\ &= \frac{p(\theta^* | y) / p(\theta_j^* | \theta_{-j}^{t-1}, y)}{p(\theta_j^{t-1} | y) / p(\theta_{-j}^{t-1} | \theta^*, y)} \\ &= \frac{p(\theta_{-j}^{t-1} | y)}{p(\theta_{-j}^{t-1} | y)} \\ &= 1 \end{aligned}$$

and thus every jump is accepted.

9.3. Data Augmentation

It is a technique that can be helpful in making problems amenable to Gibbs sampling.

- (a) In the real world, some of our data can be missing.
- (b) The likelihood function is not tractable for one reason or another but conditional on a collection of unobserved random variables, the likelihood becomes easy to handle.

Example 9.3.1. Suppose Y_0, Y_1, \dots, Y_n is a time series of random variables defined by $Y_0 = 0$ and for each $i = 1, \dots, n$, $Y_i = Y_{i-1} + s_i$ where $s_i \sim \text{Beta}(\theta, \theta)$, $\theta > 0$. Therefore:

$$(9.8) \quad Y_i | Y_0, Y_1, \dots, Y_{i-1} \sim Y_{i-1} + s_i$$

Likelihood of θ

$$(9.9) \quad \begin{aligned} f(y_0, \dots, y_n | \theta) &= f(y_0 | \theta) \prod_{i=1}^n f(y_i | y_0, \dots, y_{i-1}, \theta) \\ &= \prod_{i=1}^n f(y_i | y_{i-1}, \theta) \\ &= \prod_{i=1}^n \frac{\Gamma(2\theta)}{\{\Gamma(\theta)\}^2} (y_i - y_{i-1})^{\theta-1} \{1 - (y_i - y_{i-1})\}^{\theta-1} I[0 < y_i - y_{i-1} < 1]. \end{aligned}$$

However suppose that an observation y_i^* is missing then the likelihood is no longer available in closed form.

Solution:

z additional variables to be included in the model. $f(y | \theta)$ is not tractable but $f(y, z | \theta)$ is tractable. The posterior distribution of (θ, z) is proportional to:

$$(9.10) \quad \pi(\theta, z | y) \propto f(y, z | \theta) \pi(\theta).$$

$y = (y_0, y_i^*, y_n)$ where $z = y_i^*$.

$$(9.11) \quad \begin{aligned} f(y, y_i^* | \theta) &= f(y_0, \dots, y_n | \theta) \\ &= \prod_{i=1}^n \frac{\Gamma(2\theta)}{\{\Gamma(\theta)\}^2} (y_i - y_{i-1})^{\theta-1} \{1 - (y_i - y_{i-1})\}^{\theta-1} I[0 < y_i - y_{i-1} < 1]. \end{aligned}$$

Therefore the posterior density of θ given (y, y_i^*) is explicit:

$$(9.12) \quad \Gamma(\theta | y, y_i^*) \propto f(y, y_i^* | \theta) \pi(\theta).$$

To complete the Gibbs sampler, we also need to sample from the conditional posterior distribution of y_i^* . This has density:

(9.13)

$$f(y_i^*, y | \theta) \propto f(y, y_i^* | \theta) \\ \propto [(y_i^* - y_{i-1}^*) \{1 - (y_i^* - y_{i-1}^*)\} (y_{i+1}^* - y_i^*) \{1 - (y_{i+1}^* - y_i^*)\}]^{\theta-1}$$

on the region $y_i^* \in (y_{i-1}^*, y_{i-1}^* + 1) \cap (y_{i+1}^* - 1, y_{i+1}^*)$. This sampling can be carried out by rejection sampling.

GIBBS SAMPLING

$$\begin{aligned}y_i &\sim \text{Poisson}(\theta), \quad i = 1, \dots, k \\y_i &\sim \text{Poisson}(\lambda), \quad i = k + 1, \dots, n \\ \theta &\sim \text{Gamma}(a_1, b_1) \\ \lambda &\sim \text{Gamma}(a_2, b_2) \\ k &\sim \text{discrete uniform over } \{1, \dots, n\} \\ b_1 &\sim \text{Gamma}(c_1, d_1) \\ b_2 &\sim \text{Gamma}(c_2, d_2)\end{aligned}$$

Likelihoods:

$$\begin{aligned}f(y_I | \theta) &= \prod_{i=1}^k f(y_i | \theta) \\ f(y_J | \lambda) &= \prod_{j=k+1}^n f(y_j | \lambda) \\ f(y_i | \theta) &= \prod_{i=1}^k \frac{e^{-\theta} \theta^{y_i}}{y_i!} = \frac{e^{-k\theta} \theta^{\sum_{i=1}^k y_i}}{\prod_{i=1}^k y_i!} \\ f(y_j | \theta) &= \prod_{j=k+1}^n \frac{e^{-\theta} \theta^{y_j}}{y_j!} = \frac{e^{-(n-k)\theta} \theta^{\sum_{i=k+1}^n y_i}}{\prod_{i=k+1}^n y_i!}\end{aligned}$$

Priors:

$$\begin{aligned}\pi(\theta) &= \frac{1}{b_1^{a_1} \Gamma(a_1)} \theta^{a_1-1} e^{-\theta/b_1} \\ \pi(\lambda) &= \frac{1}{b_2^{a_2} \Gamma(a_2)} \theta^{a_2-1} e^{-\theta/b_2} \\ \pi(b_1) &= \frac{1}{d_1^{c_1} \Gamma(c_1)} b_1^{c_1-1} e^{-b_1/d_1} \\ \pi(b_2) &= \frac{1}{d_2^{c_2} \Gamma(c_2)} b_2^{c_2-1} e^{-b_2/d_2}\end{aligned}$$

Since

$$g(x; \alpha, \beta) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} & \text{for } x > 0 \\ 0 & \text{elsewhere} \end{cases}$$

where $\alpha > 0$ and $\beta > 0$ for Gamma Distribution (α, β) .

Posterior Distribution

$$\pi(\theta, \lambda, k, b_1, b_2 | y) \propto f(y_i | \theta) f(y_j | \lambda) \pi(\lambda | a_2, b_2) \pi(\theta | a_1, b_1) \pi(b_1) \pi(b_2) \pi(k)$$

Explicitly,

$$\begin{aligned} \pi(\theta, \lambda, k, b_1, b_2 | y) &\propto (e^{-k\theta} \theta^{\sum_{i=1}^k y_i}) (e^{-(n-k)\theta} \lambda^{\sum_{i=k+1}^n y_i}) (d_1^{-c_1} b_1^{c_1-1} e^{-b_1/d_1}) \\ &\quad (d_2^{-c_2} b_2^{c_2-1} e^{-b_2/d_2}) I[k \in \{1, 2, \dots, n\}] (b_2^{-a_2} \theta^{a_2-1} e^{-\theta/b_2}) (b_1^{-a_1} \theta^{a_1-1} e^{-\theta/b_1}) \end{aligned}$$

So the conditional distributions are:

(a)

$$\begin{aligned} \pi(\theta | y, \lambda, k, b_1, b_2) &\propto \int \pi(\theta, \lambda, k, b_1, b_2 | y) d\lambda db_1 db_2 \\ &\propto \int (e^{-k\theta} \theta^{\sum_{i=1}^k y_i} \theta^{a_1-1} e^{-\theta/b_1}) \text{constant} d\theta \\ &\propto \int (e^{-(k+1/b_1)\theta} \theta^{\sum_{i=1}^k y_i + a_1 - 1}) d\theta \\ &= \text{Gamma}(a_1 + \sum_{i=1}^k y_i, b_1 + k) \end{aligned}$$

(b)

$$\begin{aligned} \pi(\lambda | y, \theta, k, b_1, b_2) &\propto \int \pi(\theta, \lambda, k, b_1, b_2 | y) d\theta db_1 db_2 d\theta \\ &\propto \int (e^{-(n-k)\lambda} \theta^{\sum_{i=k+1}^n y_i} \theta^{a_2-1} e^{-\theta/b_2}) \text{constant} d\lambda \\ &\propto \int (e^{-((n-k)+b_2)\lambda} \theta^{\sum_{i=1}^k y_i + a_2 - 1}) d\lambda \\ &= \text{Gamma}(a_2 + \sum_{i=k+1}^n y_i, b_2 + (n-k)) \end{aligned}$$

(c)

$$\begin{aligned}
\pi(b_1 | y, \lambda, \theta, k, b_1) &\propto \int \pi(\theta, \lambda, k, b_1, b_2 | y) dy d\theta db_2 d\theta \\
&\propto \int b_1^{c_1-1} e^{-b_1/d_1} \theta^{a_1-1} e^{-\theta/b_1} \text{constant} db_1 \\
&\propto \int b_1^{c_1+a_1-1} e^{-b_1/(d_1+\theta)} db_1 \\
&= \text{Gamma}(a_1 + c_1, d_1 + \theta)
\end{aligned}$$

(d)

$$\begin{aligned}
\pi(b_2 | y, \lambda, \theta, k, b_1) &\propto \int \pi(\theta, \lambda, k, b_1, b_2 | y) dy d\theta db_1 dk \\
&\propto \int b_2^{c_2-1} e^{-b_2/d_2} \lambda^{a_2-1} e^{-\lambda/b_2} \text{constant} db_2 \\
&\propto \int b_2^{c_2+a_2-1} e^{-b_2/(d_2+\lambda)} db_2 \\
&= \text{Gamma}(a_2 + c_2, d_2 + \lambda)
\end{aligned}$$

(e)

$$\begin{aligned}
\pi(k | y, \lambda, \theta, b_1, b_2) &\propto \int \pi(\theta, \lambda, k, b_1, b_2 | y) dy d\theta db_1 db_2 \\
&= \frac{e^{(\lambda-\theta)k} (\theta/\lambda)^{\sum_{i=1}^k y_i}}{\sum_{j=1}^n \left\{ e^{(\lambda-\theta)j} (\theta/\lambda)^{\sum_{i=1}^j y_i} \right\}} I [k \in \{1, 2, \dots, n\}]
\end{aligned}$$

As the conditional distribution of k is discrete, thereby characterized by a probability mass function.

DATA AUGMENTATION

Suppose Y_0, Y_1, \dots, Y_n is a time series of random defined by $Y_0 = 0$, $i = 1, \dots, n$, $Y_i = Y_{i-1} + S_i$ where $S_i \sim \text{Beta}(\theta, \theta)$, $\theta > 0$. Therefore

$$Y_i | Y_0, Y_1, \dots, Y_{i-1} \sim Y_{i-1} + S_i$$

Likelihood for the observations (y_0, y_1, \dots, y_n) of θ :

$$\begin{aligned} f(y_0, \dots, y_n | \theta) &= f(y_0 | \theta) \prod_{i=1}^n f(y_i | y_0, \dots, y_{i-1}, \theta) = \prod_{i=1}^n f(y_i | y_{i-1}, \theta) \\ &= \prod_{i=1}^n \frac{\Gamma(2\theta)}{\{\Gamma(\theta)\}^2} (y_i - y_{i-1})^{\theta-1} \{1 - (y_i - y_{i-1})\}^{\theta-1} I[0 < y_i - y_{i-1} < 1]. \end{aligned}$$

However suppose that an observation y_i^* is missing. So since the likelihood is no longer in closed form, we assume that Y_1, \dots, Y_n are iid data from the mixture density:

$$f(y_i | \theta) = \frac{1}{2} \left(\frac{1}{(2\pi)^{1/2}} e^{-y_i^2/2} + \frac{1}{(2\pi)^{1/2}} e^{-(y_i - \theta)^2/2} \right).$$

and

$$L(y_i | \theta) = \prod_{i=1}^n f(y_i | \theta) \left(e^{-y_i^2/2} + e^{-(y_i - \theta)^2/2} \right).$$

Let

- z = the additional variables included in the model (z may be just a single variable or a vector containing several variables).
- θ = original parameters in the model with prior $\pi(\theta)$.
- y = vector of observations.

So the posterior distribution of (θ, z) is proportional to

$$\pi(\theta, z | y) \propto f(y, z | \theta) \pi(\theta)$$

Data augmentation proceeds by carrying out Gibbs sampling to sample successively from θ and z to produce a sample from this joint distribution. The marginal distribution of θ therefore the posterior distribution of interest.

- $y = (y_0, \dots, y_n)$ excluding y_i^* , the missing observation $z = y_i^*$

$$f(y_i^*, y \mid \theta) \propto f(y, y_i^* \mid \theta)$$

$$\propto [(y_i^* - y_{i-1}^*) \{1 - (y_i^* - y_{i-1}^*)\} (y_{i+1}^* - y_i^*) \{1 - (y_{i+1}^* - y_i^*)\}]^{\theta-1}$$

Therefore, the posterior density of θ given (y, y_i^*) is

$$\pi(\theta, y_i^* \mid y) \propto f(y, y_i^* \mid \theta)\pi(\theta)$$

the conditional posterior distribution of y_i^* is:

$$f(y_i^*, y \mid \theta) \propto f(y, y_i^* \mid \theta)$$

$$\propto [(y_i^* - y_{i-1}^*) \{1 - (y_i^* - y_{i-1}^*)\} (y_{i+1}^* - y_i^*) \{1 - (y_{i+1}^* - y_i^*)\}]^{\theta-1}$$

here $y_i^* \in (y_{i-1}^*, y_{i-1}^* + 1) \cap (y_{i+1}^* - 1, y_{i+1}^*)$.

- Here z is a sequence of n heads on tails, with one element per iteration. $z = (z_1, \dots, z_n)$ and

$$z_i = \begin{cases} 1 & \text{if i-th observation is head} \\ 2 & \text{if i-th observation is fail} \end{cases}$$

Suppose the prior of $\theta \sim \mathcal{N}(0, 1)$. Then

$$f(y_i, z_i \mid \theta) = f(y_i \mid z_i, \theta)P(z_i).$$

where

$$f(y_i \mid z_i, \theta) = \begin{cases} e^{-y_i^2/2} & \text{if } z_i = 1 \\ e^{-(y_i - \theta)^2/2} & \text{if } z_i = 2 \end{cases}$$

and $P(z_i = 1) = P(z_i = 2) = 1/2$

Using that

$$\pi(\theta, z \mid y) \propto \left\{ \prod_{i=1}^n f(y_i \mid z_i, \theta)P(z_i) \right\} \pi(\theta)$$

$$\propto \left(e^{-\frac{1}{2} \sum_{i:z_i=1} (y_i - \theta)^2} \right) \left(e^{-\theta^2} \right) \sim \mathcal{N}\left(\frac{\sum_{i:z_i=2} y_i}{1 + n_2}, \frac{1}{1 + n_2} \right)$$

where $n_2 = \#$ of observations for which $z_i = 2$

$$\begin{aligned} P(z_i = 1 \mid \theta, y) &= \frac{P(z_i = 1)}{P(z_i = 1) + P(z_i = 2)} \\ &= \frac{e^{-y_i^2/2}}{e^{-(y_i - \theta)^2/2} + e^{-y_i^2/2}} \end{aligned}$$

$$P(z_i = 2 \mid \theta, y) = \frac{P(z_i = 1)}{P(z_i = 1) + P(z_i = 2)} = \frac{e^{-(y_i - \theta)^2/2}}{e^{-(y_i - \theta)^2/2} + e^{-y_i^2/2}}$$

Then a Gibbs sampler is used to simulate the posterior distribution of $(\theta, z_1, \dots, z_n)$.

Example .0.2. A genetic linkage $y = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$ with cell probabilities $(\frac{2+\theta}{4}, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{\theta}{4})$; $0 \leq \theta \leq 1$. Prior of $\theta \sim Uniform(0, 1)$, so the posterior density of θ :

$$\pi(\theta | y) \propto f(y | \theta)\pi(\theta) \propto (2 + \theta)^{y_1}(1 - \theta)^{y_2+y_3}\theta^{y_4} I[\theta \in (0, 1)]$$

Then

- (a) Sample the posterior distribution of θ directly (e.g. via rejection sampling)
- (b) Use data augmentation

Augment the observed data (y_1, y_2, y_3, y_4) by dividing the first cell into two partitions, with respective probabilities proportional to θ and 2. That is,

$$z | y, \theta \sim Binomial(y_1, \frac{\theta}{2 + \theta})$$

Then the likelihood function is:

$$\begin{aligned} f(y, z | \theta) &= f(y | \theta)\pi(z | y, \theta) \\ &\propto (2 + \theta)^{y_1}(1 - \theta)^{y_2+y_3}\theta^{y_4} \binom{y_1}{z} \left(\frac{\theta}{2 + \theta}\right)^z \left(\frac{2}{2 + \theta}\right)^{y_1-z} \\ &= (1 - \theta)^{y_2+y_3}\theta^{y_4}\theta^2 z^{y_1-z} \binom{y_1}{z} \end{aligned}$$

The conditional posterior of θ is:

$$\begin{aligned} (\theta | y, z) &= \int f(y, z | \theta) dz dy \propto \theta^{y_4+z}(1 - \theta)^{y_2+y_3} I[\theta \in (0, 1)] \\ &= Beta(y_4 + z + 1, y_2 + y_3 + 1) \end{aligned}$$

To complete the Gibbs sampler, we also generate draws from the conditional posterior distribution of z which is

$$z | y, \theta \sim Binomial(y_1, \frac{\theta}{2 + \theta})$$

R Codes

Listing 1. Triplot Code in R

```
1 #####
2 #
3 #     A Sample Triplot by Anil Aksu #
4 #     It is developed to show some basics of R #
5 #                                     #
6 #####
7
8 ## the range of sampling
9 x=seq(-4,4,length=101)
10 ## this function gets numbers from console
11 prior=dnorm(x, mean = 0.5, sd = 0.7, log = FALSE)
12 likelihood=dnorm(x, mean = 0.49, sd = 0.65, log = FALSE)
13 posterior=dnorm(x, mean = 0.52, sd = 0.5, log = FALSE)
14
15
16 ## let's plot them
17 plot(range(x), range(c(likelihood,prior,posterior)), ...
18       type='n', xlab=expression(paste(theta)), ...
19       ylab=expression(paste("f(", theta, " )")))
20 lines(x, prior, type='l', col='blue')
21 lines(x, likelihood, type='l', col='red')
22 lines(x, posterior, type='l', col='green')
23
24 title("Prior, Likelihood and Posterior Distribution")
25 legend(
26   "topright",
27   lty=c(1,1,1),
28   col=c("blue", "red", "green"),
29   legend = c("prior", "likelihood", "posterior")
30 )
```

Listing 2. Inference Plots Code in R

```
1 #####
```

```

2 # #
3 # Posterior, Perspective and Contour plots #
4 # by Anil Aksu #
5 # #
6 #####
7
8 ## the range of sampling
9 x=seq(0,20,length=101)
10 ## this function gets numbers from console
11 posterior=dnorm(x, mean = 7, sd = 1.5, log = FALSE)
12
13
14 ## let's plot them
15 plot(range(x), range(c(posterior)), type='n', ...
      xlab=expression(paste(theta)), ...
      ylab=expression(paste("f(", theta, "|x )")))
16
17 lines(x, posterior, type='l', col='blue',lwd=5)
18
19 title("Posterior Distribution")
20 legend = c("posterior")
21
22 ## perspective plot
23 x <- seq(-10, 10, length= 30)
24 y <- x
25 f <- function(x, y) { r <- sqrt(x^2+y^2); 10 * ...
      sin(r)/r }
26 z <- outer(x, y, f)
27 z[is.na(z)] <- 1
28 op <- par(bg = "white")
29 persp(x, y, z, theta = 30, phi = 30, expand = 0.5, ...
      col = "lightblue")
30 persp(x, y, z, theta = 30, phi = 30, expand = 0.5, ...
      col = "lightblue",
31       ltheta = 120, shade = 0.75, ticktype = "detailed",
32       xlab = "X", ylab = "Y", zlab = "Sinc( r )"
33 ) -> res
34 round(res, 3)
35
36 # contour plot
37 a <- expand.grid(1:20, 1:20)
38 b <- matrix(a[,1]^2 + a[,2]^2, 20)
39 filled.contour(x = 1:20, y = 1:20, z = b,
40               plot.axes = { axis(1); axis(2); ...
41                             points(10, 10) })
42
43 ## bivariate posterior sampling

```

```

44
45 ## the range of sampling
46 x=seq(-4,6,length=101)
47 ## this function gets numbers from console
48 posterior=0.8*dnorm(x, mean = 0, sd = 1, log = ...
      FALSE)+0.2*dnorm(x, mean = 4, sd = 1, log = FALSE)
49
50 ## let's plot them
51 plot(range(x), range(c(posterior)), type='n', ...
      xlab=expression(paste(theta)), ...
      ylab=expression(paste("f(", theta, "|x )")))
52
53 lines(x, posterior, type='l', col='blue',lwd=5)
54
55 # title("Bivariate Posterior Distribution")
56 legend = c("posterior")
57
58 ## credible interval posterior plot
59
60 ## the range of sampling
61 x=seq(0,20,length=101)
62 ## this function gets numbers from console
63 posterior=dnorm(x, mean = 7, sd = 2, log = FALSE)
64
65 ## let's plot them
66 plot(range(x), range(c(posterior)), type='n', ...
      xlab=expression(paste(theta)), ...
      ylab=expression(paste("f(", theta, "|x )")))
67
68 lines(x, posterior, type='l', col='blue',lwd=5)
69
70 # title("Bivariate Posterior Distribution")
71 legend = c("posterior")

```

Listing 3. Rejection Sampling Plots Code in R

```

1 #####
2 #
3 #           Rejection sampling plots           #
4 #           by Anil Aksu                       #
5 #                                           #
6 #####
7
8
9 require(SMPracticals)
10 ## rejection sampling
11

```

```
12 ## the range of sampling
13 x=seq(-10,10,length=101)
14 ## this function gets numbers from console
15 posterior=0.6*dnorm(x, mean = 0, sd = 4, log = ...
      FALSE)+0.4*dnorm(x, mean = 6, sd = 2, log = FALSE)
16 envelope=2*dnorm(x, mean = 2, sd = 5, log = FALSE)
17 ## let's plot them
18 plot(range(x), range(c(posterior,envelope)), ...
      type='n', xlab=expression(paste(theta)), ylab="")
19 lines(x, posterior, type='l', col='blue',lwd=5)
20 lines(x, envelope, type='l', col='red',lwd=5)
21
22 title("Rejection Sampling Plot")
23 legend("topright", legend=c(expression(paste("L(", ...
      theta, ")",pi,"(", theta, ...
      ")")),expression(paste(mu,"g(", theta, "))),
24 lty=1, col=c('blue', 'red'),inset = .02)
25
26 ## British Coal Mining accidents
27 data(coal)
28 # years of coal mining accidents
29 years <- unique(as.integer(coal$date))
30 # the number of accidents in each year
31 accident <- integer(length(years))
32 for (i in 1:length(years)){
33 accident[i]<-sum(as.integer(coal$date) == years[i])
34 }
35
36 plot(years ,accident, col='blue',lwd=2, xlab="year", ...
      ylab="# of disasters")
37 #rug(coal$date)
```

BIBLIOGRAPHY

1. Allen B. Dawney. *Think Bayes: Bayesian Statistics in Python*. O'REILLY, 2013.
2. D. Gamerman and F. L. Herbert. *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference*. Chapman and Hall/CRC, 2006.
3. Peter D. Hoff. *A First Course in Bayesian Statistical Methods*. Springer Publishing Company, Incorporated, 1st edition, 2009.
4. Sheldon Ross. *Introduction to Probability Models*. Academic Press, Boston, 2014.