

1487 Chapter 5

1488 Simple Bayesian Models

1489 In this chapter, we lay out the basic principles of Bayesian inference, building on the concepts
1490 of probability we developed earlier (Chapter 3). Our overarching purpose is to use the rules of
1491 probability to show how Bayes theorem works. We will make use of the conditional rule of probability
1492 and the law of total probability, so it might be useful to review these first principles (Section 3.2)
1493 before proceeding with this chapter.

1494 We begin with the central, underpinning tenet of the Bayesian view: the world can be divided
1495 into quantities that are observed and quantities that are unobserved. Unobserved quantities include
1496 parameters in models, latent states predicted by models, missing data, effect sizes, predictions of
1497 future states, and data before we observe them. We want to learn about these quantities using
1498 observations. Bayes provides a framework to achieve that understanding, a framework that is
1499 applied in exactly the same way regardless of the specifics of the research problem at hand or the
1500 nature of the unobserved quantities we want to understand.

1501 The feature of Bayesian analysis that most clearly sets it apart from all other types of statis-
1502 tical analysis is that Bayesians treat all unobserved quantities as random variables.¹ Because the

¹There is some argument among statisticians about whether states of ecological systems and parameters governing their behavior are truly random. Ecologists with traditional statistical training may object to viewing states and parameters as random variables. These objections might proceed like this. Consider the state, “the average biomass of trees in a hectare of Amazon rainforest.” It could be argued that there is nothing random about it, that at any instant in time there *is* an average biomass that is fixed and knowable at that instant – it is determined, not random. This is true, perhaps, but the practical fact is that if we were to attempt to know that biomass, which is changing by the minute, we would obtain different values depending on when and how we measured it. These values would follow a probability distribution. So, thinking of unknowns as random variables is a scientifically useful abstraction with enormous practical benefits, benefits we will demonstrate in later chapters. We will leave arguments about whether states and parameters are “truly random” to metaphysics. As an aside, Ben Bolker (personal communication) points out that “The same traditionally trained ecologists who object to treating states as random variables don’t mind using hypothesis tests that are grounded in the idea of a long-term frequency of observation in repeated observations, which don’t sensibly exist in many cases...”

1503 behavior of random variables is governed by probability distributions, it follows that unobserved
1504 quantities can be characterized by probability distributions like those we learned about in Section
1505 3.4. Bayesian analysis uses the rules of probability (Section 3.2) to discover the characteristics
1506 of the probability distributions of unobserved quantities. Understanding those distributions enables
1507 the ecological researcher to make statements about processes tempered by honest specifications of
1508 uncertainty.

1509 It is fundamental to Bayesian analysis to understand the distinctions among things that are
1510 known vs. unknown, observed vs. unobserved, and random variables vs. fixed quantities. The first
1511 distinction is this. Things that are *known* are not random variables but rather are treated as fixed.
1512 This might seem obvious, but it can be slippery. Numerical constants, for example π , are known.
1513 Things that are not observed, for example, parameters in a model, latent states, predictions, and
1514 missing data are unknown and are always modeled as random variables. But what about things we
1515 observe?

1516 Observations of responses (i.e., the y) are always modeled as random variables. How can this
1517 be? How can something that we observe be random? The key idea here is that the y are random
1518 variables *before they are observed*. After we observe them, we have quantities in hand that represent
1519 one instance of a stochastic process. So, this one instance of observations is fixed but if we repeated
1520 our observations of the response, we would not expect to always get identical values. The sources
1521 of stochasticity in responses will be treated in greater detail as we proceed.

1522 What about observed predictor variables (i.e., covariates, the x)? Are they random or fixed?
1523 Rightly or wrongly (usually wrongly), ecologists often treat predictor variables as being observed
1524 perfectly – they are observations but they are treated as if they were known, fixed quantities. They
1525 are not random variables if we assume they are measured without error, but they *are* random
1526 variables if we assume they have measurement or sampling errors that we seek to include in our
1527 model.

1528 5.1 Bayes theorem

1529 The basic problem in ecological research is to understand processes that we cannot observe based on
1530 quantities that we can observe. We represent unobserved processes as models made up of parameters
1531 and latent states, which we will notate here as θ . We make observations y to learn about θ . Before

1532 the data are observed, we treat them as random variables. The chance of observing the data
 1533 conditional on θ is given by a probability distribution, $[y|\theta]$. Because θ is also a random variable,
 1534 it is governed by a probability distribution, $[\theta]$. We want to discover the probability distribution of
 1535 the unobserved θ conditional on the observed data, that is $[\theta|y]$. Using the basic rules of conditional
 1536 probability for two random variables,

$$[\theta|y] = \frac{[\theta, y]}{[y]} \quad (5.1.1)$$

$$[y|\theta] = \frac{[\theta, y]}{[\theta]}. \quad (5.1.2)$$

1537 Solving 5.1.2 for $[\theta, y]$ we have

$$[\theta, y] = [y|\theta] [\theta]. \quad (5.1.3)$$

1538 Substituting the right hand side of 5.1.3 for $[\theta, y]$ in 5.1.1 we obtain,

$$[\theta|y] = \frac{[y|\theta] [\theta]}{[y]}. \quad (5.1.4)$$

1539 Because y is conditional on θ , the law of total probability (equations 3.2.13 and 3.2.14) for discrete
 1540 valued parameters shows that

$$[y] = \sum_{\theta} [y|\theta] [\theta] \quad (5.1.5)$$

1541 where we are summing over all possible values of θ . For parameters that are continuous,

$$[y] = \int [y|\theta] [\theta] d\theta. \quad (5.1.6)$$

1542 Substituting the right hand side of equation 5.1.5 for $[y]$ in 5.1.4 we obtain Bayes theorem for
 1543 discrete valued parameters,

$$[\theta|y] = \frac{[y|\theta] [\theta]}{\sum_{\theta} [y|\theta] [\theta]} \quad (5.1.7)$$

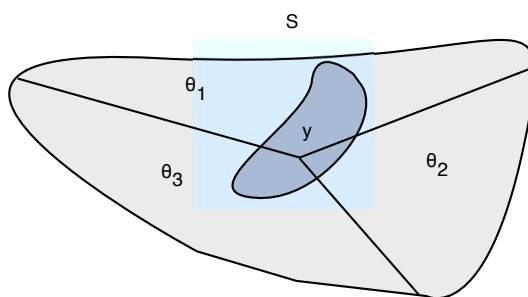
1544 and similarly substituting equation 5.1.6 for $[y]$ in 5.1.4 we find Bayes theorem for parameters that
 1545 are continuous,

$$[\theta|y] = \frac{[y|\theta] [\theta]}{\int [y|\theta] [\theta] d\theta}. \quad (5.1.8)$$

1546 Bayes theorem provides the basis for estimating the probability distribution of the unobserved
 1547 quantities θ informed by the data y . A simple example illustrates these ideas graphically (Box 5.1).

1548 **Box 5.1 Illustration of Bayes Theorem**

Imagine support for the parameter θ shown as the light colored polygon labeled S . Assume that θ can take on three values, θ_1, θ_2 , and θ_3 . We assume for simplicity that these are the *only* possible values – they are mutually exclusive and exhaustive, i.e., $\sum_i \text{area of wedge}_i = S$. The area of each θ_i wedge divided by the area of S reflects our prior knowledge of the parameter, $\frac{\text{area of wedge } \theta_i}{\text{area of } S} = \Pr(\theta_i)$. If we had no reason to favor one value of θ_i over another, $\Pr(\theta_1) = \Pr(\theta_2) = \Pr(\theta_3) = \frac{1}{3}$.



We now collect some data shown by the dark polygon y . The parameter θ controls how the data arise. So, for example, the data might be the number of survivors observed in a sample of n individuals during time Δt where θ is the probability that an individual survives over the time interval. We want to use the data to update our knowledge of θ .

Given that we have data in hand, we can limit attention to the wedge of the θ_i contained within the data polygon. The probability of θ_i is $\Pr(\theta_i|y) = \frac{\text{area of } \theta_i \text{ within } y}{\text{area of } y} = \frac{\text{area of } \theta_i \text{ within } y / \text{area of } S}{\text{area of } y / \text{area of } S} = \frac{\Pr(\theta_i \cap y)}{\Pr(y)} = \frac{\Pr(\theta_i, y)}{\Pr(y)}$. Using the conditional rule of probability to substitute for $\Pr(\theta_i, y)$ we have $\Pr(\theta_i|y) = \frac{\Pr(y|\theta_i)\Pr(\theta_i)}{\Pr(y)}$. Using $\Pr(y) = \frac{\text{area of } y}{\text{area of } S} = \sum_i \Pr(y|\theta_i)\Pr(\theta_i)$, we find Bayes theorem for discrete parameters,

$$\Pr(\theta_i|y) = \frac{\Pr(y|\theta_i)\Pr(\theta_i)}{\sum_j \Pr(y|\theta_j)\Pr(\theta_j)}. \quad (5.1.9)$$

The denominator is a normalizing constant assuring that $\sum_i \Pr(\theta_i|y) = 1$. As the number of wedges in S increases to infinity and their area decreases to 0, we have Bayes theorem for continuous parameters,

$$[\theta|y] = \frac{[y|\theta][\theta]}{\int [y|\theta][\theta] d\theta}. \quad (5.1.10)$$

1549 Understanding Bayesian inference and why it works requires that we understand each of its
 1550 components, which we now explain for continuous parameters. The *likelihood* $[y|\theta]$ (Figure 5.1.1)
 1551 plays a key role in Bayesian analysis by linking the unobserved θ to the observed y . It allows us
 1552 to answer a central question of science: “What is the probability that we would observe the data
 1553 if our deterministic model, $g(\theta)$, accurately portrays the process that gives rise to the data?” We
 1554 have seen the likelihood before (equation 4.1.4, Figure 4.2.1).

1555 The *prior distribution* of the unobserved quantities, $[\theta]$ represents our knowledge about θ before
 1556 we collect the data (Figure 5.1.1). The prior distribution can be informative, reflecting knowledge
 1557 gained in previous research, or it can be vague, reflecting a lack of information about θ before we
 1558 collected the data that are now in hand. We will treat priors in greater detail in the next section;
 1559 for now, we highlight prior distributions as one of the components of Bayes theorem.

1560 The product of the likelihood and the prior is the joint distribution² (Figure 5.1.1). We have seen
 1561 this product ($[y|\theta][\theta]$) before (equation 4.4.1), and we learned that it does not define a probability
 1562 distribution for θ because the area under the curve $[y|\theta][\theta]$ with respect to θ is not certain to equal
 1563 one.

1564 The marginal distribution of the data

$$[y] = \int [y|\theta][\theta] d\theta \quad (5.1.11)$$

1565 is the area under the joint distribution curve (Figure 5.1.1). Dividing each point on the joint
 1566 distribution $[y|\theta][\theta]$ by $\int [y|\theta][\theta] d\theta$ normalizes the curve with respect to θ , yielding the posterior
 1567 distribution $[\theta|y]$. The posterior distribution is a true probability density function that meets all
 1568 of the requirements for these functions (Section 3.4.1), including $\int [\theta|y] d\theta = 1$. Dividing the joint
 1569 distribution by $\int [y|\theta][\theta] d\theta$ assures that the posterior distribution integrates to 1, which is why $[y]$
 1570 is often referred to as a normalizing constant.

1571 Before the data are collected y is a random variable and the quantity $\int [y|\theta][\theta] d\theta$ is a marginal
 1572 distribution, a concept we will use frequently in later chapters (for review, see Section 3.4.2). It is
 1573 also called the prior predictive distribution – it tells us what we know about the *data* before they are
 1574 collected. However, after the data are collected, $\int [y|\theta][\theta] d\theta$ is a known, fixed quantity (a scalar).

²Recall that the joint distribution $[\theta, y] = [y|\theta][\theta]$.

$$[\theta|y] = \frac{[y|\theta] [\theta]}{\int_{\theta} [y|\theta] [\theta] d\theta}$$

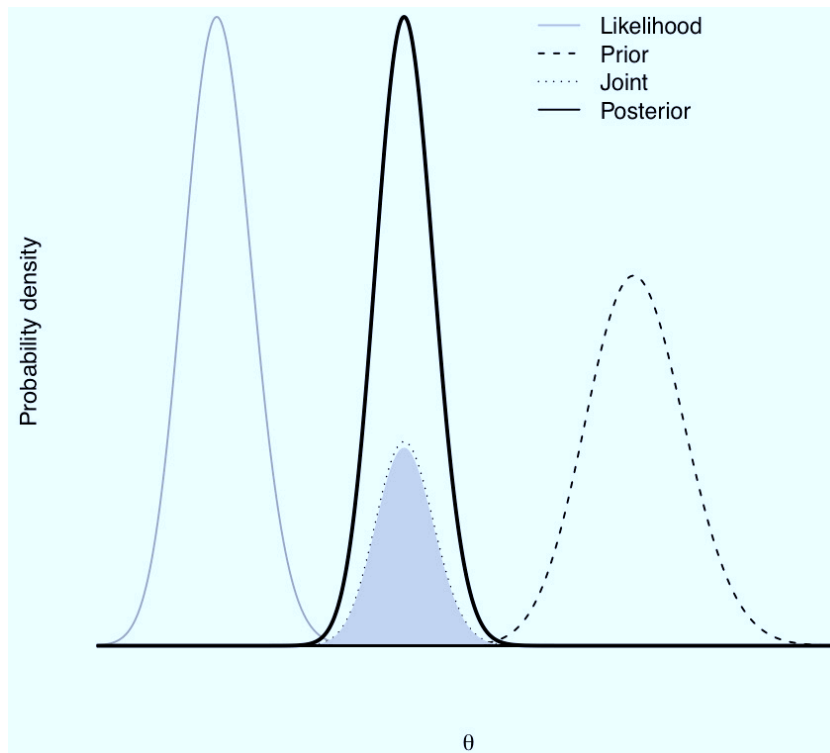


Figure 5.1.1: Illustration of Bayes theorem for data (y) and unobserved quantities (θ). The likelihood ($[y|\theta]$, grey solid line) gives the probability that we would observe the data conditional on the value of the parameter. The prior ($[\theta]$, dashed line) specifies the probability of θ based on our knowledge of θ before the data were collected. The joint distribution (dotted line) is the product of the prior and the likelihood. The marginal distribution of the data ($\int [y|\theta] [\theta] d\theta$) is the integral of the joint distribution, shown here as the shaded area. (See Section 3.4.2 for a review of the concept of marginal distributions.) The posterior is the distribution (black solid line) that results when we divide every point on the joint distribution curve by the area under the curve, effectively normalizing the joint distribution so that the area under the posterior distribution equals 1.

1575 This means that

$$[\theta|y] \propto [\theta, y] \tag{5.1.12}$$

$$\propto [y|\theta] [\theta]. \tag{5.1.13}$$

1576 We will make extensive use of this proportionality.³ We can use equation 5.1.13 to learn about the
1577 posterior distribution from the joint distribution even when we cannot directly calculate $[y]$, as will
1578 often be the case. We call equation 5.1.13 a simple Bayesian model because it represents the joint
1579 distribution of the observed and unobserved quantities as the product of the likelihood and the prior
1580 distributions.

1581 We could have developed the same ideas about discrete valued parameters using sums rather
1582 than integrals.

1583 5.2 The relationship between likelihood and Bayes

1584 The fundamental difference between inference based on maximum likelihood and inference based
1585 on Bayes theorem is that Bayes treats all unobserved quantities as random variables governed
1586 by probability distributions. This treatment is possible because dividing the joint distribution
1587 by the marginal distribution of the data assures that posterior distribution is a true probability
1588 distribution (Figure 5.2.1). This is a non-trivial result because it allows Bayesian inference to make
1589 probabilistic statements about unobserved quantities of interest. In contrast, the likelihood profile
1590 is not a probability distribution – there is nothing that assures that the area under the curve equals
1591 1 (Figure 5.2.1). Unknowns cannot be treated as random variables in the likelihood framework.
1592 Instead, likelihood depends on comparing the relative strength of evidence in data for one value
1593 of a parameter over another value. The use of prior information can be accomplished in Bayesian
1594 and likelihood analysis using $[y|\theta] [\theta]$. In likelihood, we find the values of θ that maximize $[y|\theta] [\theta]$.
1595 The normalization of this product by the marginal distribution of the data is what sets Bayesian
1596 inference apart from inference based on likelihood – it allows unobserved quantities to be treated
1597 as random variables.

³The constant of proportionality is the reciprocal of the marginal distribution of the data, which is a constant after the data are observed.

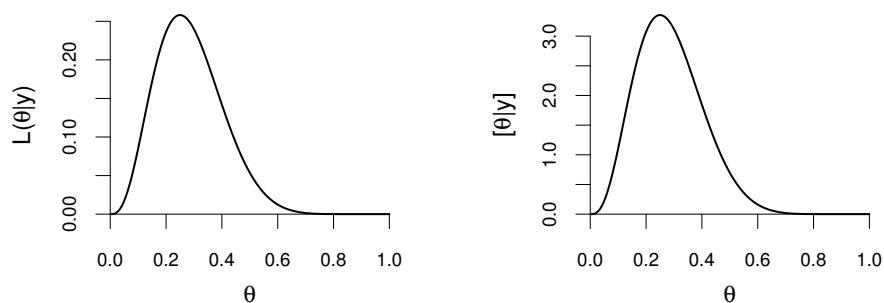


Figure 5.2.1: Likelihood profile (left panel) and posterior distribution (right panel) for the parameter probability of a success (θ) given the observation three successes on twelve trials with vague priors on θ . The shapes of the two curves are identical. The area under the likelihood profile does not equal 1. The area under the posterior distribution equals 1.

1598 5.3 Finding the posterior distribution in closed form

1599 A simple Bayesian model contains a joint distribution expressed as a likelihood multiplied by a
 1600 prior (or priors) $[y|\theta][\theta]$. There are special cases of this product where posterior distribution $[\theta|y]$
 1601 has the same form as the prior $[\theta]$. In these cases, the prior and the posterior are called *conjugate*
 1602 *distributions* (or simply conjugates) and the prior is called a conjugate of the likelihood. Conjugate
 1603 distributions are important for two reasons. For simple problems, they allow us to calculate the
 1604 parameters of posterior distributions on the back of a cocktail napkin.⁴ Moreover, the ease of
 1605 calculation of parameters of the posterior for simple problems becomes important for complicated
 1606 problems if we can break them down into parts that can be attacked one at time. We will learn
 1607 about the role of conjugates in this process in the chapter on Markov chain Monte Carlo (Chapter
 1608 7).

1609 It is perfectly possible to make use of conjugate priors effectively without knowing how each
 1610 one is derived. Seeing a single derivation (Box 5.3) is adequate background for most ecologists
 1611 who seek to use Bayesian methods. However, we will offer a couple of examples here to provide
 1612 intuition for conjugate relationships. More detailed treatment as well as tables showing the known
 1613 conjugate distributions can be found in Bayesian textbooks (e.g., Gelman, 2006). The ones we use
 1614 most frequently are shown in Appendix table A.3.

⁴It is embarrassing to do an elaborate numerical procedure to obtain results that can be obtained on a napkin.

Box 5.3 Derivation of the posterior distribution for a beta prior and binomial likelihood

We seek the posterior distribution of the parameter ϕ , the probability of a success conditional on n trials and y observed successes. The beta distribution is a conjugate prior for the binomial likelihood. Using Bayes theorem:

$$[\phi|y, n] \propto \underbrace{\binom{n}{y} \phi^y (1 - \phi)^{n-y}}_{\text{binomial likelihood}} \underbrace{\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \phi^{\alpha-1} (1 - \phi)^{\beta-1}}_{\text{beta prior}}, \quad (5.3.1)$$

where α and β are the parameters of the beta prior distribution. By dropping the normalizing constants $\left(\binom{n}{y}, \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\right)$ we obtain:

$$[\phi|y, n] \propto \underbrace{\phi^y (1 - \phi)^{n-y}}_{\text{binomial likelihood}} \underbrace{\phi^{\alpha-1} (1 - \phi)^{\beta-1}}_{\text{beta prior}}. \quad (5.3.2)$$

Simplifying:

$$[\phi|y, n] \propto \phi^{y+\alpha-1} (1 - \phi)^{\beta+n-y-1}. \quad (5.3.3)$$

Let $\alpha_{new} = y + \alpha$, $\beta_{new} = \beta + n - y$. Multiplying 5.3.3 by the normalizing constant, $\frac{\Gamma(\alpha_{new} + \beta_{new})}{\Gamma(\alpha_{new})\Gamma(\beta_{new})}$ we obtain the posterior distribution of ϕ , a beta distribution with parameters α_{new} and β_{new} :

$$[\phi|y, n] = \frac{\Gamma(\alpha_{new} + \beta_{new})}{\Gamma(\alpha_{new})\Gamma(\beta_{new})} \phi^{\alpha_{new}-1} (1 - \phi)^{\beta_{new}-1}. \quad (5.3.4)$$

1615 Here are a couple of examples to show how conjugate prior-likelihood relationships can be used
 1616 to estimate posterior distributions easily and quickly. Imagine that you are studying infection of
 1617 whitebark pine (*Pinus albicaulis*) with blister rust (*Cronartium ribicola*). You desire information on
 1618 the proportion of individuals in a stand that are infected, that is, the prevalence of the disease, ϕ .
 1619 You take a sample of 80 individuals and find 17 that are infected. What is the posterior distribution
 1620 of ϕ ? We will use the simple Bayesian model,

$$[\phi|\mathbf{y}] = \frac{[\mathbf{y}|\phi][\phi]}{[\mathbf{y}]} \quad (5.3.5)$$

1621 We have no prior knowledge of disease prevalence in the stand, so a reasonable choice for a prior
 1622 distribution of ϕ , a quantity that can take on continuous values between 0 and 1, is a beta distribu-
 1623 tion with parameters $\alpha_{prior} = 1, \beta_{prior} = 1$, i.e., $\phi \sim \text{beta}(1, 1)$ which defines a uniform distribution
 1624 over (0,1). A logical choice for the likelihood of ϕ is a binomial distribution with $y = 17$ successes
 1625 given $n = 80$ trials where we seek to know the probability of a “success,” i.e., $y \sim \text{binomial}(\phi, 80)$.
 1626 Thus,

$$\text{beta}(\phi|\alpha_{posterior}, \beta_{posterior}) = \frac{\text{binomial}(y|\phi, n) \text{beta}(\phi|\alpha_{prior}, \beta_{prior})}{[\mathbf{y}]} \quad (5.3.6)$$

1627 Using the beta-binomial conjugate prior relationship, we can calculate the parameters of the poste-
 1628 rior beta distribution using $\alpha_{posterior} = \alpha_{prior} + y$ and $\beta_{posterior} = \beta_{prior} + n - y$. So, in this example,
 1629 the posterior distribution of ϕ is $\text{beta}(1+17, 1+80-17)$ which has a mean of $\frac{\alpha_{posterior}}{\beta_{posterior} + \alpha_{posterior}} = .219$
 1630 and variance $\frac{\alpha_{posterior}\beta_{posterior}}{(\alpha_{posterior} + \beta_{posterior})^2(\alpha_{posterior} + \beta_{posterior} + 1)} = .00206$ (Section 3.4.4 and Appendix table
 1631 A.2). Using the quantile function for a beta distribution, we can calculate that the true value of ϕ
 1632 lies between .137 and .314 with probability 0.95.

1633 As a second example, suppose you are studying copepods in an arctic lake during summer. You
 1634 want to estimate the posterior distribution of the mean abundance per unit volume using

$$[\lambda|\mathbf{y}] = \frac{[\mathbf{y}|\lambda][\lambda]}{[\mathbf{y}]} \quad (5.3.7)$$

1635 Prior research has shown that lakes like the one you are studying have a mean abundance of λ_{prior}
 1636 = 52 individuals per liter with a standard deviation of 6.8. You take a sample of four scoops of
 1637 one liter of water and count the individuals they contain finding $\mathbf{y} = (64, 48, 59, 52)'$. What can

1638 we say about the abundance of copepods informed by the data and the prior estimate? A good
 1639 choice for the likelihood in this example (i.e., $[\mathbf{y}|\lambda]$) is the Poisson because the data are discrete
 1640 and because the variance is approximately the same as the mean. A gamma prior distribution (i.e.,
 1641 $[\lambda]$) is conjugate to the Poisson likelihood, so the posterior distribution for the mean of the Poisson
 1642 (i.e., $[\lambda|\mathbf{y}]$) is also gamma. Thus,

$$\text{gamma}(\lambda|\alpha_{\text{posterior}}, \beta_{\text{posterior}}) = \frac{\prod_{i=1}^4 \text{Poisson}(y_i|\lambda) \text{gamma}(\lambda|\alpha_{\text{prior}}, \beta_{\text{prior}})}{[\mathbf{y}]} \quad (5.3.8)$$

1643 The parameters of a gamma posterior are $\alpha_{\text{posterior}} = \alpha_{\text{prior}} + \sum_{i=1}^4 y_i$ and $\beta_{\text{posterior}} = \beta_{\text{prior}} + n$.
 1644 To use the prior information we must first convert the prior mean and standard deviation to prior
 1645 parameters using moment matching (Section 3.4.4), $\alpha_{\text{prior}} = \frac{\mu_{\text{prior}}^2}{\sigma_{\text{prior}}^2} = 58.5$ and $\beta_{\text{prior}} = \frac{\mu}{\sigma^2} = 1.12$.
 1646 It follows that the parameters of the gamma posterior distribution of the mean abundance are
 1647 $\alpha_{\text{posterior}} = 58.5 + 64 + 48 + 59 + 52 = 281.5$ and $\beta_{\text{posterior}} = 4 + 1.12 = 5.12$. The mean of the
 1648 posterior is $\frac{\alpha_{\text{posterior}}}{\beta_{\text{posterior}}} = 55.0$ with variance $\frac{\alpha_{\text{posterior}}}{\beta_{\text{posterior}}^2} = 10.7$ and standard deviation 3.3. The upper
 1649 0.975 quantile for a gamma distribution with parameters $\alpha = 281.5$ and $\beta = 5.12$ is 61.5 and the
 1650 lower 0.025 quantile is 48.7. Thus, the probability is 0.95 that the true mean number of individuals
 1651 per liter is between 48.7 and 61.5.

1652 5.4 More about prior distributions

1653 We devote an entire section in this chapter to prior distributions because ecologists who have not
 1654 received formal training in Bayesian methods will be especially unfamiliar with the use of priors, a
 1655 concept that, in contrast to likelihood, has no parallel in traditional statistical training. We also
 1656 include this section because ecologists often seek to minimize the influence of the prior on inference.
 1657 This is a place where it is easy to make errors. Finally, we want to advocate the thoughtful use of
 1658 informed priors in Bayesian modeling.

1659 Some view the choice of a prior in Bayesian models as a contentious topic because it is a decision
 1660 that can influence inference. However we will attempt to convince you that:

- 1661 1. There is no such thing as a noninformative prior, but certain priors influence the posterior
 1662 distribution more than others.
- 1663 2. Informative priors, when properly justified, can be tremendously useful in Bayesian modeling

1664 (and science, in general).

1665 It is important to remember that one of the objectives of Bayesian analysis is to provide informa-
1666 tion that can inform subsequent analyses; the posterior distribution obtained in one investigation
1667 becomes the prior in subsequent investigation. Thus we agree with the view of Gelman (2006) that
1668 “noninformative” priors are provisional. They are a starting point for analysis. As scientists, we
1669 should always prefer to use appropriate, well constructed, informative priors on θ .

1670 5.4.1 “Noninformative” Priors

1671 We use the double quotes in this section title because there is no such thing as a noninformative
1672 prior. By that we mean that all priors will have some influence on the posterior distribution of some
1673 transformation of the parameter you may be interested in learning about. With that said, let’s begin
1674 by studying potential priors for a very simple Bayesian model, a model for binary data. Consider
1675 the set of binary data (i.e., zeros and ones) denoted by y_i , for $i = 1, \dots, n$. If we are interested in
1676 inference concerning the probability that a given observation will be one, $p = \Pr(y = 1)$, then we
1677 could formulate the following parametric model

$$y_i \sim \text{Bern}(p), \quad (5.4.1)$$

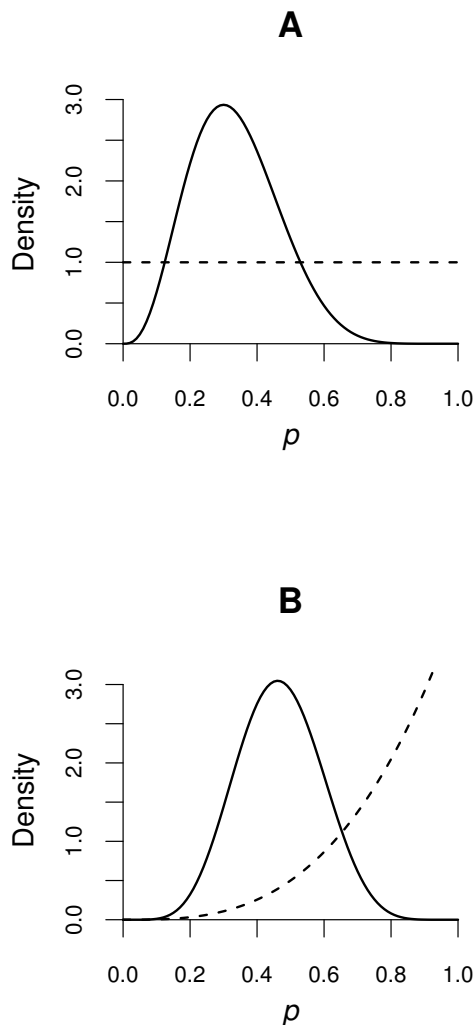
1678 where $i = 1, \dots, n$. In this case, a Bernoulli distribution is the “model” that we assume stochastically
1679 generated the data. The Bernoulli distribution contains the parameter p , thus a complete Bayesian
1680 model requires a prior distribution for p . Let’s examine a few priors for p as well as their influence
1681 on the posterior distribution for the following data set: $\mathbf{y} = (0, 0, 1, 0, 1, 0, 0, 0, 1, 0)'$.

1682 Perhaps the most commonly chosen prior for p is the uniform distribution such that $0 < p < 1$.
1683 The uniform is a specific case of the more flexible beta distribution, thus it is common to select the
1684 prior

$$p \sim \text{beta}(\alpha, \beta), \quad (5.4.2)$$

1685 where, if $\alpha = \beta = 1$, this distribution becomes a uniform. The uniform distribution is commonly
1686 thought to be “noninformative” in this setting because all possible values of p are equiprobable. The
1687 uniform can be contrasted with a prior where larger values of p are more probable, such as when
1688 $\alpha = 4, \beta = 1$. We compare the posterior distributions arising from these two choices for a prior in

1689 Figure 5.4.1. Notice how the prior in Figure 5.4.1 **B** “pulls” the posterior toward the larger values,
 1690 thus influencing it.



1690 *Figure 5.4.1:* Prior (dashed line) and resulting posterior distributions (solid line) for a model with a Bernoulli
 1691 likelihood and a beta prior with two prior specifications: **A** $\alpha = 1, \beta = 1$ and **B** $\alpha = 4, \beta = 1$.

1691 An alternative to the visual approach for assessing the influence of the prior on the posterior is
 1692 to inspect the closed form mathematical expression for the posterior (i.e., the result of conjugate
 1693 relationships, Section 5.3). For the Bernoulli-beta model⁵ we are using in this example, the posterior

⁵We showed the derivation of the expression for the posterior distribution when the prior is beta and the likelihood is binomial (Section 5.3). Recall that the Bernoulli is a special case of the binomial where the number of trials, $n = 1$.

1694 distribution for p is

$$[p|\mathbf{y}] = \text{beta} \left(\sum_{i=1}^n y_i + \alpha, \sum_{i=1}^n (1 - y_i) + \beta \right). \quad (5.4.3)$$

1695 In our simple example, using the data \mathbf{y} , the resulting beta posterior distribution has parameters
1696 $3 + \alpha$ and $7 + \beta$. Notice that larger values for α and β will have more of an effect on these parameters
1697 in the posterior. Similarly as α and β both get small, the posterior distribution appears to become
1698 less influenced by the prior (leaving only statistics related to the data in the posterior). Thus, a beta
1699 prior with $\alpha = 1, \beta = 1$ will be less influential on the posterior than a beta prior with $\alpha = 4, \beta = 1$.
1700 This is a seemingly sensible result, and one that is very commonly used to justify the specification of
1701 priors, especially for probabilities (i.e., p), regression coefficients (i.e., β), and variance components
1702 (i.e., σ^2). Perfect flatness can only be achieved in bounded priors like the beta; but, priors that
1703 *approach* flatness are often referred to as “flat” nonetheless. You will also see them called “diffuse,”
1704 “weak,” or “vague.”

1705 It is important to recognize that even the uniform prior for p technically has some influence
1706 on the posterior distribution because prior parameters $\alpha = 1, \beta = 1$ yield the posterior parameters
1707 $3 + 1, 7 + 1$, which are not the same as $3, 7$ as would be the case if only statistics related to the data
1708 appeared in the posterior. Using this argument, one might be tempted to use $\alpha = 0, \beta = 0$ as their
1709 prior parameters, but recall from the definition of the beta distribution that both parameters need
1710 to be greater than zero to ensure a valid probability distribution. Furthermore, the sensibility of
1711 using very small values for α and β in the beta prior breaks down because, as we see in Figure 5.4.2,
1712 a beta prior with $\alpha = 0.001, \beta = 0.001$ actually pulls the posterior distribution toward zero. A
1713 U-shaped prior distribution implied by the $\text{beta}(0.001, 0.001)$ has most of its mass near zero and
1714 one, thus suggesting that p is more likely to be large or small, but not moderate (i.e., close to 0.5).

1715 The take-home message is that all priors have an influence on the posterior distribution and what
1716 might seem like a good trick to minimize the prior influence may not always do what you think it
1717 should. With that said, one can always overwhelm any amount of prior influence with enough data.
1718 In our example, if n gets large, then any prior values for α and β become inconsequential in the
1719 posterior; they will be very minimal compared with the large values for $\sum_{i=1}^n y_i$ and $\sum_{i=1}^n (1 - y_i)$.
1720 Thus, to some extent, the simplest way to minimize prior influence is to collect a larger data set to
1721 begin with!

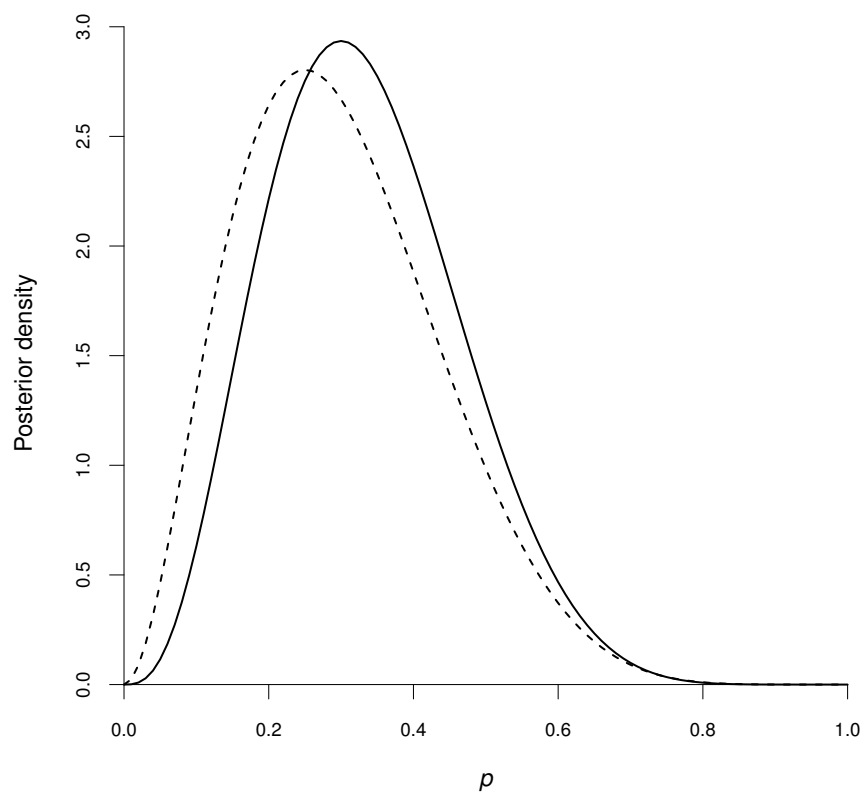


Figure 5.4.2: Resulting posterior distributions for the Bernoulli-beta model with prior specifications $\alpha = 1, \beta = 1$ (solid line) and $\alpha = 0.001, \beta = 0.001$ (dashed line).

1722 Another caution in specifying priors that appear to minimize the influence on the posterior
1723 distribution pertains to “propriety.” A proper probability distribution is a positive function that
1724 integrates to one over the support of its random variable (Section 3.4). If the function does not
1725 integrate to one, then it is termed “improper” and is not technically a valid probability distribution.
1726 That means that we can’t use it for statistical inference because all statistical theory depends on
1727 the basic axioms about probability distributions. For example, continuing the previous discussion
1728 about how to make the beta distribution less influential, we would be tempted to use $\alpha = 0, \beta = 0$.
1729 However, because both parameters need to be positive to guarantee a proper prior distribution, the
1730 $\text{beta}(0, 0)$ is not a valid probability density function and thus its use is not advised. Interestingly,
1731 the resulting posterior, which we can still work out analytically, ends up being a $\text{beta}(3, 7)$, which is
1732 proper in this specific case. Therefore, an improper prior can *sometimes* lead to a proper posterior,
1733 but that result has to be shown for the particular model being fit and almost always depends on
1734 the data. If you cannot mathematically show that your posterior is proper, then it’s best to avoid
1735 improper priors.

1736 Let’s consider another situation. Suppose you have the same data and Bayesian model but are
1737 interested in obtaining inference related to the quantity p^2 , rather than p . The seemingly benign
1738 uniform prior (i.e., $\text{beta}(1, 1)$) for p then becomes quite informative for p^2 . To illustrate this point,
1739 we can find the implied prior distribution for p^2 using a Jacobian transformation technique.⁶ In this
1740 case, if we use a uniform prior for p , the implied prior for p^2 (the quantity we desire inference about)
1741 is proportional to $1/p$. Therefore the values of p^2 under its implied prior are not equiprobable like
1742 they are for p . Specifically, the uniform prior for p says that smaller values for p^2 are more probable
1743 than larger values. That result may not be what we had in mind when we chose the $\text{beta}(1, 1)$
1744 prior for p . A prior whose information about a parameter does not change when you transform the
1745 parameter is called “invariant to transformation.” The Jeffreys prior was developed for this exact
1746 purpose, to help specify priors that are invariant to transformation.

1747 The Jeffreys prior depends on the form of the likelihood (also called the data model). More
1748 specifically, the Jeffreys prior is proportional to Fisher information raised to the half power.⁷ That

⁶The details of this technique are beyond the scope of this book, but can be found in any graduate level mathematical statistics book.

⁷This is the same Fisher information that is used to find asymptotic variance of an MLE.

1749 is, if we can calculate the negative expectation of the second derivative of the log likelihood ex

$$-E_y \left(\frac{d^2 \log[\mathbf{y}|p]}{dp^2} \right) \quad (5.4.4)$$

1750 then we have something proportional to the Jeffreys prior. The Jeffreys prior for our ongoing binary
1751 data example (5.4.1) is, perhaps surprisingly, a beta(0.5, 0.5) distribution. This Jeffreys prior will
1752 contain the same information for p as it will for p^2 , or any other transformation of p for that matter.
1753 Unfortunately, the Jeffreys prior is often called “noninformative,” but for those same reasons cited
1754 above, it is not noninformative. One might use a Jeffreys prior when they don’t know what else
1755 to use, in this case, because it happens to be invariant to transformation. For our example, the
1756 Jeffreys prior is U-shaped; not quite as extremely U-shaped as the beta(0.001, 0.001) prior for p , but
1757 it will still give more prior preference to those values close to zero and one than 1/2. The Jeffreys
1758 prior for this particular example turns out to be proper, but it is not guaranteed to be proper for
1759 all models.

1760 You will commonly see a normal prior with large variance used as a prior distribution for a
1761 variety of parameters. A normal distribution with large variance (i.e., $N(0, 1000)$) is often justified
1762 as an attempt to find a vague prior that is conjugate.⁸ Given that the normal distribution is not
1763 bounded, it will be impossible to make it perfectly flat, so the large variance serves as a mechanism
1764 to at least spread it out.⁹ A normal with infinite variance would be flat, but then would also
1765 not be proper (i.e., would not integrate to one). The use of a normal prior with large, but finite,
1766 variance seems to work fine without complications for parameters that are means, and where plenty
1767 of information exists in the data. However, for other types of parameters, say transformations of
1768 probabilities such as $\text{logit}(p)$, the normal prior with large variance can have a dubious influence on
1769 the posterior.

1770 To illustrate our point, suppose we have the same Bernoulli model for the binary data we’ve

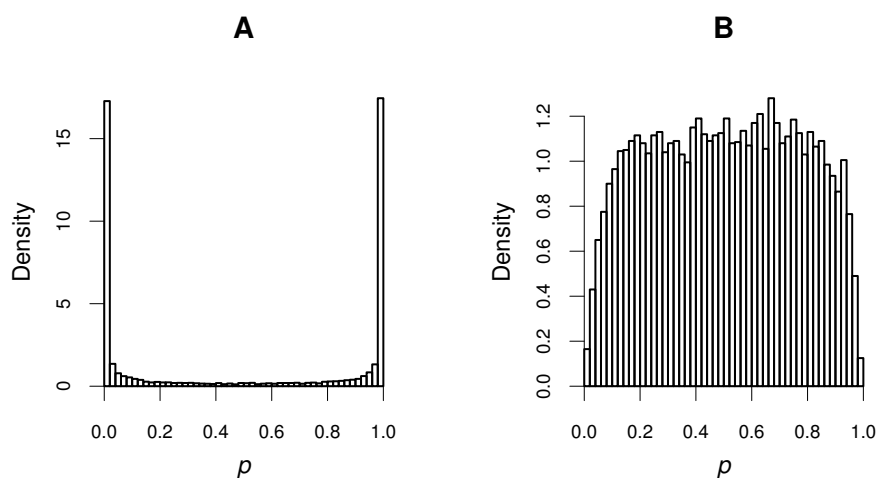
⁸Recall that conjugacy occurs when a prior and posterior have the same form. There can be many analytical and computational advantages to using conjugate priors, but they are not always the best choice.

⁹Keep in mind that the prior variance is always relative to the scale of the parameter. For example, if the data indicate that a parameter should be 1000, then a $N(0, 100)$ prior for that parameter will probably be informative unless the sample size is huge, because a variance of 100 is small relative to 1000, as is the prior mean of 0.

1771 been discussing in this Section and we use the prior for $\text{logit}(p)$ such that

$$\text{logit}(p) \sim \text{normal}(0, \sigma_p^2), \quad (5.4.5)$$

1772 where, σ_p^2 is set to be a large number. The question is: What prior does this imply for p (rather than
 1773 $\text{logit}(p)$)? Simulating 10,000 random draws from a normal distribution and taking the inverse logit
 1774 transformation, we can see in Figure 5.4.3 that a normal with $\sigma_p^2 = 100$ is much more informative
 than a normal with $\sigma_p^2 = 2$.



1775 *Figure 5.4.3:* Histograms of p based on samples drawn from prior distributions for A) $\text{logit}(p) \sim$
 1776 $\text{normal}(0, 100)$ and B) $\text{logit}(p) \sim \text{normal}(0, 2)$.

1777 Priors with large variance might seem vague or less informative, but they are not always, thus it
 1778 is a good idea to check the implied prior distribution in the transformation of the parameter for
 1779 which you desire inference. You can do this by varying the values of the parameters for the prior
 and examining how that variation effects the posterior.

1780 It's worth mentioning that the same methods are commonly used for choosing priors for vari-
 1781 ance components. In fact, we show models that contain such priors throughout this book. It is
 1782 important to realize that such priors are not truly noninformative, for the same reasons we have
 1783 described above. For example, suppose we have data that can be sufficiently modeled with a normal
 1784 distribution

$$y_i \sim \text{normal}(\mu, \sigma^2), \quad (5.4.6)$$

1785 for $i = 1, \dots, n$, and where the mean μ is assumed to be known (for now) and our interest lies
 1786 in obtaining inference about the variance σ^2 . A conjugate prior for the variance parameter is the
 1787 inverse gamma distribution

$$\sigma^2 \sim \text{inverse gamma}(\alpha, \beta), \quad (5.4.7)$$

1788 which yields the posterior for σ^2

$$[\sigma^2 | \mathbf{y}] = \text{inverse gamma} \left(\frac{n}{2} + \alpha, \frac{\sum_{i=1}^n (y_i - \mu)^2}{2} + \beta \right). \quad (5.4.8)$$

1789 Notice that, similar to the beta posterior we discussed previously, here in the inverse gamma
 1790 posterior for σ^2 , if α and β get small then the influence of the prior on the posterior is minimized.
 1791 Thus, it is common to see priors for variance components specified as inverse gamma (0.001, 0.001),
 1792 in an attempt to minimize prior influence (but see Gelman, 2006) However, these priors are not
 1793 “noninformative” and are not invariant to transformation. An example of where this sort of prior
 1794 can be misleading is if one was interested in obtaining inference about the standard deviation σ ,
 1795 rather than the variance σ^2 .¹⁰ As an alternative, the Jeffreys prior could be used for σ^2 . For this
 1796 model, the Jeffreys prior turns out to be proportional to $1/\sigma^2$, which has the form of an inverse
 1797 gamma with $\alpha = 0$ and $\beta = \infty$. This formulation for the inverse gamma does not yield a proper
 1798 prior because both parameters (α and β) must be positive. However, the Jeffreys prior, as always,
 1799 is invariant to transformation. Like in the case with the Bernoulli model previously discussed, the
 1800 Jeffreys prior for σ^2 yields a proper posterior as long as at least one observation is available (i.e.,
 1801 $n \geq 1$).

1802 Finally, there is another approach to finding priors that are minimal in their influence on the
 1803 posterior; these priors are called “reference” priors. A reference prior is found by maximizing the
 1804 Kullback-Leibler divergence between the posterior and prior distributions.¹¹ The heuristic concept
 1805 behind reference priors is that a prior which is as different as possible from the posterior may
 1806 be desirable if you has no actual prior information or expertise and just needs a default prior to

¹⁰This is more common than you might think, as it is easier for us to interpret the standard deviation σ than the variance σ^2 .

¹¹The development of this concept is beyond the scope of this book, but in short, the K-L divergence provides a way to measure discrepancy between two distributions; it involves explicit integration and can be difficult to compute in practice, making this approach quite technical.

1807 use so that you can still obtain Bayesian inference.¹² Interestingly, for univariate parameters, the
1808 reference prior approach yields the Jeffreys prior! However, in multivariate situations, the reference
1809 prior needs to be found for each individual model where it is being used. Actually calculating the
1810 correct reference prior can be quite analytically and numerically challenging. The field of objective
1811 Bayesian inference focuses on this task for various models.

1812 **5.4.2 Informative Priors**

1813 We have learned that all priors influence the posterior in some way, but that we can often assess the
1814 amount of influence and sometimes even control it. Philosophically, when formulating statistical
1815 models, we might ask ourselves why we're trying to limit the influence of the prior on the posterior
1816 in the first place. The illusion of objectivity has been put on a pedestal in science, almost to
1817 the extent that we are to believe that only new data can be used to reach scientific conclusions.
1818 Extrapolating this concept to Bayesian statistics would then imply that we *should* be looking for
1819 priors that have no influence on the posterior (hence the previously mentioned subfield of objective
1820 Bayesian inference). However, a not often recognized point is that all parametric statistical modeling
1821 approaches are subjective, including maximum likelihood. The very fact that we have to choose
1822 a likelihood function implies that we have made a strong assumption about the data generating
1823 mechanism. Nonparametric statistical approaches seek to minimize such assumptions, but they
1824 make their own set of strong assumptions based on their associated computational algorithms for
1825 providing inference. Any constraint you put on the data or parameters in order to obtain inference
1826 imparts subjectivity. As we discuss in Chapter 9, the various forms of regularization, including
1827 penalization methods and model selection, put extreme constraints on parameters, yet they are used
1828 throughout statistics and across all applied fields without much fanfare concerning their inherent
1829 subjectivity. More importantly, these approaches are recognized as being helpful in many ways!

1830 Our view is that we would be remiss if we ignored decades of important scientific learning in
1831 the field of ecology and that there should be a way to rigorously incorporate this learning into
1832 our statistical approaches. Fortunately, the Bayesian framework provides such a mechanism. The
1833 posterior distribution itself is a formal, mathematically valid way to combine information from
1834 current as well as previous scientific studies. In that light, it is not difficult to see that the posterior

¹²Some argue that this very concept seems contrary to the Bayesian spirit by trying to avoid its biggest utility, the ability to properly account for previous research efforts in making scientific conclusions.

1835 distribution and Bayesian framework are a mathematical representation of the scientific method
1836 itself.

1837 In the scientific method, we use existing data and expertise to formulate hypotheses about how
1838 the world works, then we make conclusions and update hypotheses using new data. In Bayesian
1839 statistics, we summarize our understanding of how the world works in a prior distribution and then
1840 “update” (i.e., compute the posterior distribution) our understanding using new data. Science would
1841 be completely haphazard if we threw out everything we knew about the world every time we began a
1842 new study. Haphazard is not even a strong enough word to describe science performed in a manner
1843 where we pretend to be completely ignorant about our study system; perhaps lazy or irresponsible
1844 would be a better descriptor! In all seriousness, we challenge the reader (and ourselves) to provide
1845 an example of a parameter in a statistical model they wish to fit where they know absolutely nothing
1846 about it. Nothing at all. At a minimum, we should all know at least the support (i.e., the values the
1847 parameter can assume) for any parameter, but we often know quite a bit more than that. Ignoring
1848 prior information when you actually have it, is like selectively throwing away data before an analysis.

1849 Instead, we argue that science would be better off if we all took the time to carefully collect
1850 and represent our prior understanding of parameters in Bayesian models. Doing so can be hard
1851 work, as it sometimes requires a mathematical transformation of moments into natural parameters
1852 occurring in the distribution we, as experts, value as best representing the data and parameters.
1853 It also could include being more responsible in our knowledge of preexisting scientific findings, for
1854 example, by more carefully reading the literature and translating those findings into quantitative
1855 information we can use in our prior. Formulating honest and responsible priors may also involve
1856 communicating with other experts on the topic under study, probing them for details that can be
1857 represented in probability distributions to serve as priors. Yes, this is beyond what we normally do
1858 in statistical analyses, but Bayesian methods provide the tools to incorporate such information and
1859 we should be obligated to use them responsibly.

1860 Aside from our inherent obligation to be responsible in accounting for the body of accumulated
1861 scientific knowledge when we make new inference, so-called informative priors can be quite helpful.
1862 For example, strong priors can be beneficial in the following ways:

- 1863 • They allow us to borrow strength across several sources of information, including different

1864 data sources and expert knowledge. Given that priors are most influential when paired with
1865 small data sets, it can be incredibly helpful to obtain meaningful inference by having a formal
1866 mechanism to combine several smaller, but independently collected data sets into a single
1867 modeling framework. An additional likelihood involving a separate data source can often be
1868 written as a prior in the original model containing the primary data source. We cover use of
1869 multiple likelihoods in the joint distribution in Section 6.2.5.

1870 • Informative priors stabilize computational algorithms. This benefit is not an inferential one,
1871 but definitely a practical one. When statistical models accumulate parameters in such a way
1872 that the ratio of unknowns to knowns grows, the probability surfaces we need to explore during
1873 the fitting process can acquire pathological problems such as lack of identifiable parameters,
1874 multicollinearity, and flat likelihood or posterior surfaces. Such issues can cause numerical
1875 approaches to become unstable (e.g., failure to converge). Stronger priors add definition to
1876 the surfaces that are being explored by the statistical fitting algorithms and thus improve
1877 computational stability.

1878 • Stronger priors offer a formal way to place constraints on the unknowns in statistical models.
1879 A seldom recognized fact is that such constraints are the basis for important inferential tools
1880 such as model selection. We cover this topic in great detail throughout Chapter 9, but as a
1881 preview we note here that important concepts such as information criteria, penalized likelihood
1882 methods, ridge regression, Lasso, and cross-validation are used regularly in many fields and can
1883 all be effectively considered as different ways to improve inference through the use of stronger
1884 priors. Most statisticians are now recognizing that imposing a constraint on an optimization
1885 problem (e.g., like maximizing a likelihood) is the same concept as specifying a prior in a
1886 Bayesian model and both can be helpful for the same reasons.

1887 Excellent examples of the benefits of using informative priors can be found in Crome et al. (1996);
1888 McCarthy and Masters (2005); Elder et al. (2006) and McCarthy et al. (2008).

Up to now, we have considered informative priors as single distributions as if they were obtained from a single, previously conducted investigation. What do we do if we have multiple sources of prior knowledge informing a parameter, θ ? Recall the idea of a mixture distribution (Section 3.4.5). We can compose a prior from multiple previous studies by mixing their estimates of θ . A prior on

θ using information from L different studies can be written as

$$[\theta] = \sum_{l=1}^L w_l [\theta]_l \quad (5.4.9)$$

$$\sum_{l=1}^L w_l = 1 \quad (5.4.10)$$

1889 where the w_l are weights. If we believe that all studies were conducted equally well, then the w_l
 1890 would be chosen to be equal. As an example, assume we had three studies of the intercept (β_0)
 1891 in a regression with an associated variance and we wanted to combine them in a prior. We might
 1892 reasonably use

$$[\beta_0] = \frac{1}{3} \text{normal}(\beta_{0,1}, \sigma_1^2) + \frac{1}{3} \text{normal}(\beta_{0,2}, \sigma_2^2) + \frac{1}{3} \text{normal}(\beta_{0,3}, \sigma_3^2). \quad (5.4.11)$$

1893 Now that we can see the potential value of priors informed by single or multiple studies, we need
 1894 to know how to represent existing scientific knowledge in the form of a probability distribution.
 1895 Several different approaches exist for manifesting expert knowledge about a parameter into a prior
 1896 distribution, but rather than cover each one generically, instead consider the following example.

1897 5.4.3 Example: Priors for Moth Predation¹³

1898 A particular species of moth rests during the day on tree trunks (it is active at night), and their
 1899 coloration acts as camouflage to protect them against predatory birds. A study was conducted to
 1900 evaluate predation of a common moth species. Suppose that n sites (for $i = 1, \dots, n$) were selected,
 1901 and a varying number of dead moths (N_i) were glued to tree trunks at each site. After 24 hours,
 1902 the number of moths that had been removed (y_i), presumably by predators, were recorded. A
 1903 reasonable data model for the moth counts would be a binomial with N_i “trials” per site such that

$$y_i \sim \text{binomial}(N_i, p), \quad (5.4.12)$$

1904 where the parameter p corresponds to the probability of predation and is the unknown about which
 1905 we desire inference. Consider the following three scenarios in formulating an appropriate prior

¹³This example is gratefully modified from the excellent text of Ramsey and Schafer (2012) using ideas from Kiona Ogle.

1906 distribution for this model:

- 1907 1. We desire a relatively vague prior that contributes information equivalent to two additional
1908 “placed” moths and an expected prior probability of predation of 0.5.
- 1909 2. We desire an informative prior based on a previous observational study that reported an
1910 average of 10% (standard deviation of 2.5%) of the moths in a population were eaten by
1911 predators in a 24 hour time period.
- 1912 3. We desire an informative prior based on a pilot study that suggests the proportion removed
1913 in any given 24 hour period is unlikely to exceed 0.5 or be less than 0.1.

1914 In scenario 1, the fact that we do not feel we have much prior information pertaining to p means
1915 that we want to spread out the probability mass in the prior between zero and one such that our
1916 prior has a mean of 0.5, but no real strong preference for any range of values. A beta distribution
1917 could work well here such that $p \sim \text{beta}(\alpha, \beta)$. But how do we assess the information content of the
1918 prior in terms of an effective increase in sample size? The answer comes from looking at the form
1919 of the posterior distribution for this model:

$$[p|\mathbf{y}] = \text{beta} \left(\sum_{i=1}^n y_i + \alpha, \sum_{i=1}^n (N_i - y_i) + \beta \right). \quad (5.4.13)$$

1920 In (5.4.13) we can see a very similar form for the posterior as we had in the Bernoulli model
1921 we discussed previously where each of the updated posterior parameters contain a sum of two
1922 components, one coming from the data and one coming from the prior. In the first parameter
1923 $\sum_{i=1}^n y_i + \alpha$, we see that the sum of y_i over all sites contains the total number of moths placed that
1924 were preyed upon plus the prior parameter α . In the second posterior parameter $\sum_{i=1}^n (N_i - y_i) + \beta$
1925 we see that it is the number of total moths not predated plus prior parameter β . Thus if we set
1926 $\alpha = 1$ and $\beta = 1$, it is kind of like adding two moths to the sample size in such a way that it
1927 does not impose any preference for predation. In this case, the implied prior is a $\text{beta}(1, 1)$ or a
1928 uniform distribution. Of course, we could have said that we wanted a uniform to begin with, but
1929 it is instructive to see that the prior parameters α and β can be thought of as augmenting the
1930 sample size if that helps specify prior information. Knowing that, what prior would be induced if
1931 we had the equivalent of 10 extra moths worth of prior information such that 60% were in favor

1932 of predation and the other 40% were against predation? The answer could be easily visualized by
 1933 plotting a beta probability density function with parameters $\alpha = 4$ and $\beta = 6$.

1934 In scenario 2, we have information from a former study on moth predation. That study provides
 1935 inference pertaining to the mean and standard deviation of the proportion of moths predated upon.
 1936 In translating this information into our beta prior, we can consider the mean and variance equations
 1937 associated with a beta random variable:

$$E(p) = \frac{\alpha}{\alpha + \beta} \quad (5.4.14)$$

$$\text{Var}(p) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \quad (5.4.15)$$

1938 Setting $E(p) = 0.1$ and $\sqrt{\text{Var}(p)} = 0.025$, we can back solve for α and β to find the appropriate
 1939 prior (as we discussed in Section 3.4.4). Letting the reader check our algebra, we arrive at $\alpha = 14.3$
 1940 and $\beta = 128.7$ as parameters in our prior.

1941 Scenario 3 is slightly more involved, but entirely realistic, in that it is common for prior informa-
 1942 tion to arise as bounds on likely values for a parameter. In this scenario, if we assume that the term
 1943 “unlikely” implies that p should fall between a lower bound and upper bound with high probability
 1944 (e.g., 95%), then we need to take a similar approach to the moment matching technique, but instead
 1945 of relating moments to the results of a pilot study, we relate quantiles of the distribution to the
 1946 results of a pilot study. That is, we need to solve the system of equations

$$\int_0^{0.1} \text{beta}(p|\alpha, \beta) dp = 0.025 \quad (5.4.16)$$

$$\int_{0.5}^1 \text{beta}(p|\alpha, \beta) dp = 0.025, \quad (5.4.17)$$

1947 for α and β . The system of equations are really just representing the quantile function associated
 1948 with the beta distribution. This calculation would be quite difficult to perform analytically (i.e.,
 1949 pencil and paper) but could be approximated numerically using an optimization algorithm in a
 1950 mathematical or statistical software package. We used the function `optim()` in R (R Core Team,
 1951 2013) to find the appropriate prior parameter values $\alpha = 4.8$ and $\beta = 12.7$ by minimizing the
 1952 difference between the output of beta quantiles function (i.e., `qbeta()` in R) and .025.

1953 So, as we can see, there are a variety of ways to convert preexisting scientific information and
1954 expertise into probability distributions for use as priors in Bayesian models. These informative
1955 priors can be very useful in many ways, but only when care is taken to appropriately specify them.
1956 It is a common concern that if Bayesian models fall into the wrong hands, they could be misused by
1957 those seeking to mislead science or policy. The fact is, even under such dubious intentions, the priors
1958 would have to be clearly written out in any scientific communication and would be scrutinized just as
1959 any other scientific finding is scrutinized during peer review. Furthermore, for those with villainous
1960 intentions, there are much easier ways to mislead science or the general public, for example, by
1961 outright fabrication of scientific studies. We feel that carelessness by well intending scientists (in
1962 the field, lab, or in specifying inappropriate likelihoods or priors) is probably a much more common
1963 cause of erroneous inference than mischief.

1964 **5.4.4 Guidance**

1965 We admit that the cautionary statements in this Section could make the choice of priors seem
1966 complicated and difficult; however that is not our aim. We feel that priors can be an important
1967 component of science and can be helpful in obtaining models that are useful for inference. Our
1968 goal in this discussion of priors is to instill a sense of awareness about the decisions being made in
1969 the model building process. If you are more thoughtful about the specification of priors and the
1970 associated consequences after reading this Section, then we have done our job.

1971 The fact is, not many of these details are made clear in other texts on applied Bayesian statistics
1972 and we wrote this Section, at least in part, as a reminder to ourselves to think deeply about how we
1973 can manifest prior scientific knowledge into the form of a probability distribution for use as a prior.
1974 You'll notice that we commonly use default priors in examples throughout this book. It would seem
1975 that by doing so, we encourage this practice, but in reality we don't claim to be experts in all of the
1976 applied subjects in the diverse examples we offer. Thus, it is with a touch of "do as we say, not as
1977 we do" that we suggest that our model specifications throughout are only placeholders for a model
1978 that might actually be used by an expert in the relevant field. This Section also serves as a prelude
1979 to Chapter 8 where we show a concrete example of the value of prior information and to Chapter
1980 9 where we describe ways that priors are an example of regularization, an approach widely used in
1981 statistics to improve model fit.

1982 Although we have provided you with several approaches for specifying priors for specific models
1983 in this Section, it would be too lengthy to list all possible options for all possible models. Thus, in
1984 the big picture, we echo the guidance provided by Seaman III et al. (2012) and leave you with a few
1985 further general diagnostics and remedies to consider when specifying priors in Bayesian models:

1986 • Bear in mind that one of the objectives of Bayesian analysis is to provide information that can
1987 inform subsequent analyses; the posterior distribution obtained in one investigation becomes
1988 the prior in subsequent investigation. Thus we agree with the view of Gelman (2006) that
1989 vague priors are provisional – they are a starting point for analysis. As scientists, we should
1990 always prefer to use appropriate, well constructed informative priors.

1991 • Visualize the prior you choose in terms of the parameters for which you desire inference.
1992 We did this above for the $\text{logit}(p)$ (i.e., Figure 5.4.3). Sometimes this can be accomplished
1993 analytically (i.e., with pencil and paper, using calculus), but it’s often easier to just simulate
1994 values from your prior, then transform them to represent the quantity you want inference on
1995 and plot a histogram.

1996 • Perform a prior sensitivity analysis. Try several different priors, maybe by simply choosing
1997 different prior variances, and see how much the posterior distribution moves around as a result.
1998 You’ll often see little posterior sensitivity to priors when there is a high data to parameter
1999 ratio. However, if the posterior is sensitive to the prior and you truly desire a prior that is only
2000 weakly informative, you will need to rethink your prior by changing its form or parameters .
2001 Alternatively, you must carefully justify your choice of prior in relation to the inference you
2002 seek.

2003 • An influential prior could be indicated if the posterior inference differs greatly from maximum
2004 likelihood inference. Of course, this can only be checked in models where both approaches
2005 can be implemented easily, so it may not be practical for more complicated Bayesian models.
2006 Still, in some Bayesian models, inference will approach what would be obtained with inference
2007 based on maximum likelihood if certain priors are used.

2008 • Parameters that are dependent should probably have priors that acknowledge the dependence.
2009 We are often lured into thinking that we can get away with specifying independent priors for

2010 parameters, but a prior that is a joint multivariate distribution or conditional distribution for
2011 one parameter given another is often more appropriate. An example is in regression where
2012 $y_i \sim \text{normal}(\beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i}, \sigma^2)$ for $i = 1, \dots, n$. In this case, it is common to use
2013 the independent normal priors, $\beta_0 \sim \text{normal}(0, \sigma_\beta^2)$ and $\beta_2 \sim \text{normal}(0, \sigma_\beta^2)$, but it can
2014 be helpful to use a multivariate normal prior for both regression coefficients simultaneously:
2015 $\boldsymbol{\beta} \sim \text{normal}(\mathbf{0}, \boldsymbol{\Sigma})$, where $\boldsymbol{\beta}$ is the vector containing β_1 and β_2 , $\mathbf{0}$ is a vector of zeros, and $\boldsymbol{\Sigma}$
2016 is a covariance matrix.

2017 • Keep in mind that even in large sample size situations there may be not be enough information
2018 in the data to tease apart different parameters, regardless of their priors. This is more of
2019 an identifiability problem¹⁴ rather than a problem with the prior and the form of model
2020 itself should be reconsidered. An example of where this can happen is with binomial data
2021 $y_i \sim \text{binomial}(N, p)$ for $i = 1, \dots, n$ where N and p are unobserved. There is not enough data
2022 in the world to learn about both N and p individually in this case, but a strong prior on one of
2023 the two parameters (if warranted) can help focus the inference on the other. However, without
2024 sufficient prior information this is not a really useful model in an inverse (i.e., statistical)
2025 setting.

¹⁴Parameters are identifiable if they can be estimated given a large amount of data. They are unidentifiable if they cannot. See section 6.3 for a more complete definition of identifiability.