# MVE550: Stochastic Processes and Bayesian Inference

Petter Mostad

31st December 2020

## Foreword

The course "MVE550 Stochastic processes and Bayesian inference" is mostly about stochastic processes. When introducing the course in 2018, we wanted it to also have an element of inference, i.e., theory about how one can find stochastic process models appropriate for given data. Further, we wanted this inference to use a Bayesian framework. However, we could not find existing teaching material that perfectly fitted our plans. The solution was to use [Dob16] as the main textbook, but to supplement it with some additional material, contained in this compendium. Further information about Bayesian inference for stochastic inference can be found in for example [IRW12].

For the 2020 version of the course, each chapter will be successively revised, and will appear on Canvas before the corresponding lecture.

# Contents

For	ewor	rd	iii	
Cor	ntent	s	v	
List	t of I	Figures	vii	
1	Basi	ics of Bayesian inference	1	
	1.1	Tossing a coin	2	
	1.2	The Beta and Binomial distributions	5	
	1.3	The Poisson Gamma conjugacy	8	
	1.4	Summary so far	13	
	1.5	Bayesian inference using discretization	14	
	1.6	Bayesian inference using numerical integration	16	
	1.7	Exercises	17	
<b>2</b>	Hide	den Markov Models	19	
	2.1	Examples of Hidden Markov Models	20	
	2.2	Inference questions for HMMs	22	
3	Som	e basic inference for Markov chains and HMMs	<b>25</b>	
	3.1	The Multinomial Dirichlet conjugacy	25	
	3.2	Inference for time-homogeneous Markov chains with finite state		
	~ ~	space	26	
	3.3	Interence for HMMs, known model parameters <sup>*</sup>	30	
	3.4	Interence for HMMs, unknown model parameters	32	
	3.5	Exercises	34	
4	Som	e basic inference for Branching processes	39	
	4.1	Using a Binomial likelihood	39	
	4.2	Using a Multinomial likelihood	41	
	4.3	Exercise	42	
5	Markov chain Monte Carlo (MCMC)			
	5.1	Bayesian inference with the Metropolis Hastings algorithm	44	
	5.2	Using a random walk proposal	45	
	5.3	Using independent proposals	46	
	5.4	Gibbs sampling	46	

	5.5 Advantages and disadvantages with MCMC for Bayesian inference	48	
Ap	pendices	<b>49</b>	
Α	Some solutions to some ExercisesA.1 Exercises from Chapter 1A.2 Exercises from Chapter 3A.3 Exercises from Chapter 4	<b>51</b> 51 55 59	
в	List of some probability distributions	61	
С	List of some conjugacies	65	
Bib	Bibliography		

# List of Figures

1.1	The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The	
	prior used is $\theta = 0.7^x \cdot 0.3^{1-x}$ where $x \sim \text{Bernoulli}(0.5)$	3
1.2	The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The	
	prior used is $Uniform(0, 1)$	7
1.3	The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The	
	prior used is $Beta(33.4, 33.4)$ .	8
1.4	The the prior and posterior probability for $\theta$ when the prior	
	Beta(33.4, 33.4) is used. The prior is the dotted line	9
1.5	The three posterior distributions $Gamma(20, 1)$ , $Gamma(44, 2)$ , and $Gamma(67, 3)$ .	11
1.6	Two possible predictions for the fourth observation $k_4$ : One derived in our Bayesian computations (shown with pluses) and one derived	
	with a classical approach (shown in circles)	12
2.1	A hidden Markov model. Note that this is not a transition graph relating states of a Markov chain, but instead a graph showing the dependencies between the random variables $x_0, x_1, \ldots, x_T$ and	
2.2	$y_0, y_1, \ldots, y_T$	20
	observed actual sales of the medicine	21

## CHAPTER 1

## **Basics of Bayesian inference**

This course is mostly about stochastic processes. Such processes can function as models for many real phenomena where some uncertainty is involved. In the simplest cases, we can set up a precise stochastic model based only on reasonable assumptions, and then go on to make predictions from these models. For example, when throwing a dice, it is reasonable to assume that the probability of obtaining each of the outcomes 1 through 6 is 1/6. From this we can compute such things as the probability of obtaining a total of 9 in the first three throws, or the expected wait until we get 3 consecutive sixes. Similarly, using a deck of 52 playing cards, it is reasonable to assume that each draw from it is independent, and we can compute such things as the probability of being dealt a straight flush.

However, for most potential applications of stochastic processes, and of mathematical statistics in general, the situation is more complex. We cannot make predictions of future observations based only on reasonable assumptions, we must also use earlier observations, *data*, to find a reasonable stochastic model. Then we can make predictions from this model. For example, if we want to predict the range of an electric car on full batteries, we could use data for the ranges of similar cars. Using this data, we would build a model for the range of the car in question, and we could then use the model to make predictions. A simple model in this situation could be a normal distribution, with *parameters*  $\mu$  and  $\sigma$ , representing the expectation and standard deviation of the range.

Building a stochastic model using data is called *inference*. There are two quite common ways of thinking, or *paradigms*. One is the *classical* or *frequentist* approach. In this approach, we start with building a stochastic model which can be used for making the predictions we want, defining the model in terms of some *parameters* which are regarded as *unknown*. We then use the data to *estimate* these parameters. Finally, plugging in these estimates in the model, we can use it to make the predictions we want.

The alternative approach is *Bayesian inference*. In this approach, we build a stochastic model using only general reasonable assumptions, but we include in the model random variables representing both the data we have observed and the future observations we want to predict. We then compute the conditional distribution for the future observations given that the variables representing data are fixed to the observed values. This conditional distribution is used for prediction. So we start with a stochastic model, we *update* it using the observations from the data (in a way we *learn* from the data), and then we use the updated model for prediction.

An example may make this clearer. Consider repeated throws of a six-sided dice. If you know that the dice is fair, all sequences of equal length of outcomes will have the same probability: Observing 1,1,1 will have the same probability as observing 2,4,1, namely  $(1/6)^3 = 0.00463$ . But what if you suspect the dice is not fair? In practice, you would throw the dice a number of times, and if one outcome appears more often than other outcomes, you might start to suspect that the dice is loaded in favour of this outcome. This would increase your belief that this outcome would appear again, in your next throw. So if you have thrown 1,2,3,1,1,4,1,2,1,1,5, continuing this sequence with a 1 would seem more probable than continuing it with a 6. Even a short sequence like 1,1,1 would be slightly more probable than a sequence like 2,4,1. In other words, if you suspect a loaded dice, the outcomes are no longer independent.

Using a Bayesian approach, one could set up a model for how the dice could be loaded, which would yield predicted probabilities for any sequence of observed outcomes. Let  $p_1, p_2, \ldots, p_6$  denote the probabilities of observing six sequences that are identical except for the last throw, which would be  $1, 2, \ldots, 6$ , respectively. If you have observed the throws these sequences have in common, the conditional probabilities for the next throw to be either of  $1, 2, \ldots, 6$  given the observed outcomes would be  $p_1/(p_1 + \cdots + p_6), \ldots, p_6/(p_1 + \cdots + p_6)$ , respectively.

## 1.1 Tossing a coin

Let us explore the idea above in greater detail in a slightly simpler setting: A loaded coin has probability 0.7 for either H (heads) or T (tails), but you do not know which. In fact, you think there is an equal probability that it is loaded in favour of heads or tails. The probability for observing k heads in n throws is now

$$\Pr(k) = 0.5 \cdot \text{Binomial}(k; n, 0.7) + 0.5 \cdot \text{Binomial}(k; n, 0.3)$$
(1.1)

Here we write, for example, Binomial(k; n, 0.7) for the value at k of the Binomial probability mass function with parmeters n and 0.7. We obtain the formula above by conditioning on whether the coin is loaded in favour of either H or T: The probability of each of these possibilities is 0.5, and, given each choice, the probabilities of k heads in n throws is either Binomial(k; n, 0.7) or Binomial(k; n, 0.3), respectively.

Figure 1.1 illustrates a particular sequence of observations of heads (represented as crosses at 1) and tails (represented as crosses at 0). Together with each observation, we also plot the probability, *before* this observation is made but *given* all the previous observations, of observing heads. Before any observations are made, the probabilities of observing either heads of tails is 0.5, because of the symmetry of the situation. Once we start making observations, the prediction for the next observation will jump up and down a bit, depending on those observations. However, after a while, it settles close to 0.3: we then have such a substantial overweight of tails in the data that it seems reasonable to believe that the coin is loaded towards tails. In a way, we have *learned* from the data that the coin is loaded this way.

The probabilities displayed in Figure 1.1 may be computed as follows: The conditional probability of heads after observing a specific sequence of heads



Figure 1.1: The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The prior used is  $\theta = 0.7^x \cdot 0.3^{1-x}$  where  $x \sim \text{Bernoulli}(0.5)$ .

and tails is equal to the probability of observing this sequence, followed by observing heads, divided by the probability of observing the sequence. Thus, if the sequence contains  $y_H$  heads and  $y_T$  tails, the probability is given by

$$\Pr(\text{heads} \mid \text{sequence}) = \frac{0.5 \cdot 0.7^{y_H + 1} 0.3^{y_T} + 0.5 \cdot 0.3^{y_H + 1} 0.7^{y_T}}{0.5 \cdot 0.7^{y_H} 0.3^{y_T} + 0.5 \cdot 0.3^{y_H} 0.7^{y_T}}$$
(1.2)

Here, we have used that the probability of observing a specific sequence with  $y_H$  heads and  $y_T$  tails is  $\theta^{y_H}(1-\theta)^{y_T}$ , where  $\theta$  is the probability of observing heads each time. So we have implicitly used the order of the observations of heads and tails as part of the data. However, we get the same results by assuming that the data only contains the counts  $y_H$  and  $y_T$ , see Exercise 4 of Section 1.7.

To sum up, Bayesian inference is done with the following steps:

- 1. Based on reasonable assumptions, create a stochastic model containing random variables representing observed data and observations you would like to predict.
- 2. Use for prediction the conditional distribution of prediction variables given that the the data variables are fixed to their observed values.

A frequentist approach in the situation above might focus on estimating whether the coin is loaded towards heads or towards tails. With n low, one might conclude that there is not enough information to reliably make an estimate, and continue to predict heads in the next throw with probability 0.5. Then, as n increases, one would reach a point where one would conclude that one could reliably estimate the direction of the bias. From that point on, the predicted probability for heads would be exactly 0.7 or 0.3, depending on the estimate.

#### 1. Basics of Bayesian inference

Although the Bayesian approach above is not formulated in terms of estimating the direction of the bias, we may reformulate it as stepwise learning about a stochastic variable  $\theta$  representing the true probability of heads. Thus in our current setup  $\theta$  has the two possible values 0.7 and 0.3, with some probability for each. We may reformulate the model as encompassing two dependent random variables,  $\theta$  and k (the count of heads after n trials):

$$\begin{array}{rcl} x & \sim & \mathrm{Bernoulli}(0.5) \\ \theta & = & 0.7^x \cdot 0.3^{1-x} \\ k \mid \theta & \sim & \mathrm{Binomial}(n, \theta). \end{array}$$

Here, we write Bernoulli(0.5) for a random variable which has value 1 with probability 0.5 and otherwise value 0. We also write Binomial $(n, \theta)$  for a random variable which is Binomially distributed with parameters n and  $\theta$ . We write  $x \sim \text{Bernoulli}(0.5)$  and  $k \mid \theta \sim \text{Binomial}(n, \theta)$  to indicate that the random variables x and  $k \mid \theta$  have the given distributions.

Now, let y be 1 or 0 depending on whether the n + 1'st throw is heads or not. Using standard probability theory formulas, the conditional distribution of y given the count k of heads in the n first throws can be written

$$\Pr(y \mid k) = \sum_{\theta=0.3,0.7} \Pr(y, \theta \mid k) = \sum_{\theta=0.3,0.7} \Pr(y \mid \theta, k) \Pr(\theta \mid k)$$
$$= \sum_{\theta=0.3,0.7} \Pr(y \mid \theta) \Pr(\theta \mid k).$$
(1.3)

The key step above is that we can write  $\Pr(y \mid \theta, k) = \Pr(y \mid \theta)$ . This is because once we know the value of  $\theta$ , the probability that y is 1 does not depend on k: In fact, the probability that y = 1 is exactly  $\theta$ . According to the computation above, we can compute the probability  $\Pr(y = 1 \mid k)$  by first computing the two probabilities  $\Pr(\theta = 0.3 \mid k)$  and  $\Pr(\theta = 0.7 \mid k)$ , and then multiplying these as weights with the probabilities  $\Pr(y = 1 \mid \theta = 0.3) = 0.3$ and  $\Pr(y = 1 \mid \theta = 0.7) = 0.7$ .

At this point, we introduce a generic notation for probability mass functions which may also be used for probability density functions: For example, we write  $\pi(k)$  instead of  $\Pr(k)$  and  $\pi(y \mid \theta, k)$  instead of  $\Pr(y \mid \theta, k)$ . If  $z \sim \text{Exponential}(\lambda)$  so that z is a continuous random variable with an Exponential distribution with parameter  $\lambda$ , we also write  $\pi(z) = \lambda \exp(-\lambda z)$  for the density function. This generic notation is helpful, as so many probability computations are the same whether the underlying functions are probability mass functions or probability density functions.

The distribution  $\pi(\theta \mid k) = \Pr(\theta \mid k)$  is called the *posterior* for  $\theta$ . The unconditional distribution for  $\theta$ , in which  $\theta$  has probability 0.5 for both 0.7 and 0.3 in our case, is called the *prior*. We may compute the posterior using *Bayes* formula:

$$\pi(\theta \mid k) = \frac{\pi(k \mid \theta)\pi(\theta)}{\pi(k)} = \frac{\pi(k \mid \theta)\pi(\theta)}{\sum_{\theta} \pi(k \mid \theta)\pi(\theta)}$$

That Bayes formula appears in our approach is the reason why we call it Bayesian inference. In our case, we get

$$\pi(\theta \mid k) = \frac{\pi(k \mid \theta)\pi(\theta)}{\sum_{\theta} \pi(k \mid \theta)\pi(\theta)}$$

$$= \frac{\text{Binomial}(k; n, \theta) \cdot 0.5}{\text{Binomial}(k; n, 0.3) \cdot 0.5 + \text{Binomial}(k; n, 0.7) \cdot 0.5}$$
$$= \frac{\binom{n}{k} \theta^k (1 - \theta)^{n-k}}{\binom{n}{k} 0.3^k 0.7^{n-k} + \binom{n}{k} 0.7^k 0.3^{n-k}} = \frac{\theta^k (1 - \theta)^{n-k}}{0.3^k 0.7^{n-k} + 0.7^k 0.3^{n-k}}$$

Using Equation 1.3 we now get

$$\pi(y \mid k) = 0.3 \frac{0.3^{k} 0.7^{n-k}}{0.3^{k} 0.7^{n-k} + 0.7^{k} 0.3^{n-k}} + 0.7 \frac{0.7^{k} 0.3^{n-k}}{0.3^{k} 0.7^{n-k} + 0.7^{k} 0.3^{n-k}}$$

Comparing with Equation 1.2, we see that we have arrived at exactly the same result as we got there.

## 1.2 The Beta and Binomial distributions

Above, we made the rather curious assumption that  $\theta$  was either equal to 0.3 or to 0.7. A more realistic assumption is that  $\theta$  is just some real number between 0 and 1. Specifically, let us now assume that  $\theta$  has as prior the uniform distribution on the interval [0, 1], so that  $\pi(\theta) = 1$ .

Assume further that k out of n observations are heads. Even if  $\theta$  is now a continuous variable instead of a discrete one, we can still compute the posterior using Bayes formula, we just need to use an integral instead of a sum:

$$\pi(\theta \mid k) = \frac{\pi(k \mid \theta)\pi(\theta)}{\pi(k)}$$
(1.4)  

$$= \frac{\pi(k \mid \theta)\pi(\theta)}{\int_{\theta}\pi(k,\theta) d\theta}$$
  

$$= \frac{\pi(k \mid \theta)\pi(\theta)}{\int_{\theta}\pi(k \mid \theta)\pi(\theta) d\theta}$$
  

$$= \frac{\text{Binomial}(k; n, \theta)}{\int_{\theta}\text{Binomial}(k; n, \theta) d\theta}$$
  

$$= \frac{\binom{n}{k}\theta^{k}(1-\theta)^{n-k}}{\int_{\theta}\binom{n}{k}(1-\theta)^{n-k} d\theta}$$
  

$$= \frac{\theta^{k}(1-\theta)^{n-k}}{\int_{\theta}\theta^{k}(1-\theta)^{n-k} d\theta}$$

To go on, we might compute the integral in the denominator. To do so, we may use a shortcut, looking up the density of for the Beta distribution. In fact  $x \in [0, 1]$  has a Beta distribution with parameters  $\alpha > 0$  and  $\beta > 0$  if its density is

$$\pi(x \mid \alpha, \beta) = \frac{1}{\mathcal{B}(\alpha, \beta)} x^{\alpha - 1} (1 - x)^{\beta - 1}$$

where  $B(\alpha, \beta)$  is the *Beta function*, defined by

$$B(\alpha,\beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

where the Gamma function  $\Gamma(t)$  in turn is defined for t > 0 by

$$\Gamma(t) = \int_0^\infty x^{t-1} \exp(-x) \, dx.$$

Right now, the important thing for us is that the Beta density integrates to 1, so that, for all  $\alpha > 0$  and  $\beta > 0$ ,

$$\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} \, d\theta = \mathcal{B}(\alpha,\beta).$$

Plugging this into the computations above by setting  $\alpha - 1 = k$  and  $\beta - 1 = n - k$ , we get

$$\pi(\theta \mid k) = \frac{\theta^k (1-\theta)^{n-k}}{B(k+1, n-k+1)}$$

However, we can now use the definition of the Beta density again, to recognize that the posterior density  $\pi(\theta \mid k)$  is in fact a Beta density, specifically a Beta density with parameters k + 1 and n - k + 1.

Our goal is to find the probability of heads in the n + 1'st throw assuming that k out of the n first throws were heads. Similar to above, we may compute

$$\pi(y = 1 \mid k) = \int_{\theta} \pi(y = 1, \theta \mid k) d\theta \qquad (1.5)$$
$$= \int_{\theta} \pi(y = 1 \mid \theta, k) \pi(\theta \mid k) d\theta$$
$$= \int_{\theta} \pi(y = 1 \mid \theta) \pi(\theta \mid k) d\theta$$
$$= \int_{\theta} \theta \pi(\theta \mid k) d\theta$$

We could compute this integral. But we may also notice that the integral is the expectation of the posterior distribution  $\theta \mid k$ . Having derived that this posterior is the Beta(k+1, n-k+1) distribution, we can look up its expectation, finding that it is (k+1)/(k+1+n-k+1) = (k+1)/(n+2), so that

$$\pi(y = 1 \mid k) = \frac{k+1}{n+2}.$$

In Figure 1.2 we show, as in the previous section, a sequence of observed heads and tails, together with the probability of observing heads given all *previous* observations. The difference compared to Figure 1.1 is that we now use the prior Uniform(0, 1) for  $\theta$ , while in Figure 1.1 we use a prior where  $\theta$  is equal to 0.3 or 0.7 with equal probability. Thus the result does not stabilize as easily. But it does seem to stabilize eventually.

Notice how we repeatedely took advantage of knowledge about the Beta distribution in the computations above. This simplified our computations, but, in fact, we may go one step further in simplification: As we are computing a posterior density for  $\theta$ , we know that whatever we compute will always integrate to 1 when we integrate it as a function of  $\theta$  over the possible values for  $\theta$ , i.e., the interval [0, 1]. Thus, there is no loss of information if we multiply or divide by factors that do not depend on  $\theta$ . These factors can always be reinstated in the end, by using the requirement that our density must integrate to 1. To take advantage of this idea, we define the symbol  $\propto_{\theta}$  (read "proportional to") to mean that two expressions are identical up to a factor not involving  $\theta$ . Thus we can write for example  $\theta \propto_{\theta} 3\theta$  and  $\theta/(1+\theta) \propto_{\theta} \alpha\theta/(1+\theta)$ . Using this notation, the computations of Equation 1.4 can be written

$$\pi(\theta \mid k) \propto_{\theta} \pi(k \mid \theta) \pi(\theta) = \text{Binomial}(k; n, \theta) \propto_{\theta} \theta^{k} (1 - \theta)^{n-k}.$$
 (1.6)



Figure 1.2: The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The prior used is Uniform(0, 1).

By comparing with the density for a Beta distribution, we see that the posterior  $\pi(\theta \mid k)$  must necessarily be a Beta(k + 1, n - k + 1) density. This trick of removing uninteresting factors until we need them is going to be used repeatedly in the rest of this text.

Above, we assume a Uniform(0, 1) prior for  $\theta$ . However, the computation of Equation 1.6 suggests that, as long as the prior has the form  $\theta^{\text{something}}(1 - \theta)^{\text{something else}}$ , we will get a posterior that has the form of a Beta density. So, specifically, we now assume that  $\theta$  has a Beta( $\alpha, \beta$ ) prior for some parameters  $\alpha > 0$  and  $\beta > 0$ . The computation of the posterior density now becomes

$$\pi(\theta \mid k) \propto_{\theta} \pi(k \mid \theta) \pi(\theta) \propto_{\theta} \theta^{k} (1-\theta)^{n-k} \theta^{\alpha-1} (1-\theta)^{\beta-1} = \theta^{\alpha+k-1} (1-\theta)^{\beta+n-k-1}$$

from which we can read off that the posterior  $\theta \mid k$  is a Beta $(\alpha + k, \beta + n - k)$  distribution. The computations of Equation 1.5 apply unchanged to our more general situation. Referring again to what we know about the expectation of the Beta distribution, we get

$$\pi(y=1 \mid k) = \frac{\alpha+k}{\alpha+\beta+n}$$

We now have the possibility to do slightly more realistic learning about the biasedness of the coin. Neither guessing that  $\theta$  is either 0.3 or 0.7 or assuming that it can be any number between 0 and 1 with equal probability seems very realistic. Rather, one might guess that the coin is not too far away from fair, but it might be slightly unfair, i.e.,  $\theta$  might for example most likely be in the interval from 0.4 to 0.6. Selecting a Beta density that is symmetric and has approximately 90% of its probability mass in this interval (see Exercise 5 of

Section 1.7), we find that we may use the prior Beta(33.4, 33.4). Figure 1.3 is comparable to Figures 1.1 and 1.2, but now using this new prior. Because we have now put much more information into the prior, the prediction probabilities are much more stable from the start. Another way to illustrate what is going on is with Figure 1.4. It shows the prior density Beta(33.4, 33.4) and the posterior density Beta(33.4 + 11, 33.4 + 19) after considering all of the 30 observations illustrated in Figure 1.3. Notice how the posterior is slightly narrower than the prior, as it is based on more information and thus represents less uncertainty. It is also slightly shifted to the left, as there are 11 heads and 19 tails in the data.



Figure 1.3: The probability of heads at each point in a sequence of observations, conditioning on the previous observations of heads and tails. The prior used is Beta(33.4, 33.4).

## 1.3 The Poisson Gamma conjugacy

In the previous section, we saw that when considering data with a Binomial $(n, \theta)$  likelihood and using a prior for  $\theta$  with a Beta distribution, the posterior for  $\theta$  also became a Beta distribution. This kind of situation is in fact quite common in basic Bayesian inference. We say that a family of distributions is *conjugate* to a likelihood if selecting the prior in the family leads to a posterior in the same family. Thus the Beta family of distributions is conjugate to the Binomial likelihood, when considering the probability  $\theta$  as the parameter. In this section, we will look at the Poisson Gamma conjugacy. Note that Appendix C contains an overview of several examples of conjugacy.

Assume you are monitoring the number of incoming requests for data to some internet database. Assuming these requests happen independently, we see in Chapter 6 of Dobrow how it may be reasonable to model the number of requests per time unit as Poisson distributed. In other words, if k is the



Figure 1.4: The the prior and posterior probability for  $\theta$  when the prior Beta(33.4, 33.4) is used. The prior is the dotted line.

number of requests per time unit, the probability mass function is

$$\pi(k \mid \theta) = e^{-\theta} \frac{\theta^k}{k!},$$

where  $\theta$  is the expected number of requests for this time unit. Assume you count the number of such requests for successive time units. After a number of such counts, you want to predict the count for the coming time unit.

The situation is quite similar to the one in the previous section, and indeed, we can use much of the same thinking. As successive counts are independent if we know the true value of  $\theta$ , we can handle the situation as follows: We set up a prior density  $\pi(\theta)$  for  $\theta$ , representing our knowledge about the expected count before we have made any actual counts. Then, we find a way to compute the posterior density  $\pi(\theta \mid k_1)$  given an observed count  $k_1$ . Notice that, if we then make another count,  $k_2$ , we can update our knowledge about  $\theta$  again, but now using  $\pi(\theta \mid k_1)$  as the prior and obtaining a posterior density  $\pi(\theta \mid k_1, k_2)$ . Continuing like this for some counts, we obtain a final posterior  $\pi(\theta \mid k_1, \ldots, k_n)$ . We can then use this posterior to make a prediction for the next count  $k_{n+1}$ , using computations like those in Equation 1.5:

$$\pi(k_{n+1} \mid k_1, \dots, k_n) = \int_{\theta} \pi(k_{n+1} \mid \theta) \pi(\theta \mid k_1, \dots, k_n) \, d\theta.$$
(1.7)

To make computations in practice, we need to decide on a prior  $\pi(\theta)$  to start with. In the previous example, we used the uniform distribution on the interval [0, 1] to indicate starting out with "no knowledge". In our current situation,  $\theta$  can be any positive number. One possibility might be to try to use a uniform distribution on the interval  $[0, \infty)$ . Notice that there is no such thing, as the integral of any positive function constant on this interval is infinite. However, without going into the technical arguments here, it turns out that in Bayesian statistics, we may use "densities" that integrate to infinity: We call these densities "improper". This will work fine as long as the posterior density we compute is an ordinary "proper" density.

So we might use a "constant density" on  $[0, \infty)$ ; we would denote this as  $\pi(\theta) \propto_{\theta} 1$ , for  $\theta \geq 0$ . However, this improper density may actually not correspond very well to "having no knowledge" about  $\theta$ . In fact, it would for example appear to assign the same probability to the intervals [1, 2] [10, 11], and [1000000, 1000001]. A better representation of "no knowledge" about the parameter  $\theta$  might be that the intervals [1, 2], [10, 20] and [1000000, 2000000]have the same prior probability. It would then be reasonable to use the prior

$$\pi(\theta) \propto_{\theta} \frac{1}{\theta}$$

Notice that this prior is also improper, as the integral of  $1/\theta$  over  $[0,\infty)$  is infinite.

Assuming we have observed a count k, we can now get, using Bayes formula,

$$\pi(\theta \mid k) \propto_{\theta} \pi(k \mid \theta) \pi(\theta) \propto_{\theta} e^{-\theta} \theta^{k} \cdot \frac{1}{\theta} = e^{-\theta} \theta^{k-1}$$

Using a similar trick as in the previous section, we look up the density for a Gamma distribution:

Gamma
$$(\theta; \alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} \theta^{\alpha-1} \exp(-\beta\theta)$$

Thus we see that we must have

$$\pi(\theta \mid k) = \text{Gamma}(\theta; k, 1).$$

Continuing with the script from the previous section, we can now try out using the prior

$$\pi(\theta) = \text{Gamma}(\theta; \alpha, \beta)$$

Computations with Bayes formula give

$$\pi(\theta \mid k) \propto \pi(k \mid \theta) \pi(\theta) \propto_{\theta} e^{-\theta} \theta^k \theta^{\alpha - 1} \exp(-\beta \theta) = \theta^{\alpha + k - 1} \exp(-(\beta + 1)\theta)$$

so the posterior is a  $\text{Gamma}(\theta; \alpha + k, \beta + 1)$  distribution, and we have proved that the Gamma family of distributions is conjugate to the Poisson likelihood.

It is now easy to determine the effect of observing counts  $k_1, k_2, \ldots, k_n$ . Starting with the distribution  $\operatorname{Gamma}(\alpha, \beta)$  for  $\theta$ , each time we observe a count  $k_i$  we add  $k_i$  to the first parameter and 1 to the second parameter. Thus, after n repeated updates, we get a  $\operatorname{Gamma}(\alpha + \sum_{i=1}^{n} k_i, \beta + n)$  distribution. It is worth noticing that the two improper densities for  $\theta$  that we considered at the start can fit into this framework: The prior  $\pi(\theta) \propto_{\theta} 1$  could be named a "Gamma(1,0)" density, resulting in a Gamma  $(1 + \sum_{i=1}^{n} k_i, n)$  posterior. Similarly, the prior  $\pi(\theta) \propto_{\theta} 1/\theta$  could be named a "Gamma(0,0)" density, resulting in a Gamma  $(\sum_{i=1}^{n} k_i, n)$  posterior.

We can also make the computation of Equation 1.7 explicit:

$$\pi(k_{n+1} \mid k_1, \dots, k_n) \tag{1.8}$$

$$= \int_{\theta} \pi(k_{n+1} \mid \theta) \pi(\theta \mid k_1, \dots, k_n) d\theta$$
  

$$= \int_{\theta} e^{-\theta} \frac{\theta^{k_{n+1}}}{k_{n+1}!} \frac{(\beta+n)^{\alpha+\sum_{i=1}^{n}k_i}}{\Gamma(\alpha+\sum_{i=1}^{n}k_i)} \theta^{\alpha+\sum_{i=1}^{n}k_i-1} \exp\left(-(\beta+n)\theta\right) d\theta$$
  

$$= \frac{(\beta+n)^{\alpha+\sum_{i=1}^{n}k_i}}{\Gamma(\alpha+\sum_{i=1}^{n}k_i)k_{n+1}!} \int_{\theta} \theta^{\alpha+\sum_{i=1}^{n+1}k_i-1} \exp\left(-(\beta+n+1)\theta\right) d\theta$$
  

$$= \frac{(\beta+n)^{\alpha+\sum_{i=1}^{n}k_i}}{\Gamma(\alpha+\sum_{i=1}^{n}k_i)k_{n+1}!} \cdot \frac{\Gamma(\alpha+\sum_{i=1}^{n+1}k_i)}{(\beta+n+1)^{\alpha+\sum_{i=1}^{n+1}k_i}}.$$

In the last step, we have again compared with the density for a Gamma distribution to compute the integral.



Figure 1.5: The three posterior distributions Gamma(20, 1), Gamma(44, 2), and Gamma(67, 3).

Let us illustrate our results so far in a concrete example. We start with the prior  $\pi(\theta) \propto_{\theta} 1/\theta$ , and we then make the consequtive observations  $k_1 = 20$ ,  $k_2 = 24$ , and  $k_3 = 23$ . The posteriors after one, two, and three observations are Gamma(20, 1), Gamma(44, 2), and Gamma(67, 3), respectively. These posteriors are illustrated in Figure 1.5. We see that our knowledge about  $\theta$  is increasing in each step, as the variances of the distributions are decreasing. Note that the expectations of the Gamma distributions are 20/1 = 20, 44/2 = 22, and 67/3 = 22.33, respectively.

Figure 1.6 illustrates predictions we may make for the fourth observation  $k_4$  after observing  $k_1, k_2, k_3$ . The pluses represent probabilities for various values of  $k_1$  computed according to the formula for  $\pi(k_4 \mid k_1, k_2, k_3)$  above. The circles represents a more classical prediction: Using the three observations  $k_1, k_2, k_3$ ,



Figure 1.6: Two possible predictions for the fourth observation  $k_4$ : One derived in our Bayesian computations (shown with pluses) and one derived with a classical approach (shown in circles).

the maximum likelihood estimate for the parameter  $\theta$  is 67/3 = 22.33. The circles plots a Poisson mass density function with parameter 22.33. We see that the predicted distribution is then narrower than the one using Equation 1.8. The reason is that with the classical prediction, we have "thrown away" our remaining uncertainty about  $\theta$ : We believe that its value is exactly 22.33 instead of believing that its value is given by the posterior Gamma(67,3) depicted in Figure 1.5. In this sense, we can say that the frequentist model is *overfitted*.

The computations shown in Equation 1.8 may seem a bit messy. Let us close this section with showing how we can do such computations in a more structured way that make them simpler to follow. In general, if we are in a situation with conjugacy, so that all the densities  $\pi(k \mid \theta)$ ,  $\pi(\theta)$ , and  $\pi(\theta \mid k)$  are expressed in nice analytic formulas, the last relevant density,  $\pi(k)$ , can also be expressed in a nice analytic formula using the following formula, which can be derived immediately from Bayes rule:

$$\pi(k) = \frac{\pi(k \mid \theta)\pi(\theta)}{\pi(\theta \mid k)}$$
(1.9)

If all the densities on the right have nice formulas, we get a nice formula also for  $\pi(k)$  on the left. Moreover, as the right-hand side contains  $\theta$  while the left-hand side does not, we know that if we put in explicit formulas on the left-hand side, the  $\theta$  must somehow disappear from the formula after simplifications.

Let us assume that  $\pi(k \mid \theta) = \text{Poisson}(k; \theta), \ \pi(\theta) = \text{Gamma}(\alpha, \beta), \ \text{and} \ \pi(\theta \mid k) = \text{Gamma}(\alpha + k, \beta + 1).$  We then get

$$\pi(k) = \frac{\pi(k \mid \theta)\pi(\theta)}{\pi(\theta \mid k)}$$

$$= \frac{\text{Poisson}(k;\theta) \operatorname{Gamma}(\theta;\alpha,\beta)}{\operatorname{Gamma}(\theta;\alpha+k,\beta+1)}$$

$$= \frac{e^{-\theta}\frac{\theta^{k}}{k!} \cdot \frac{\beta^{\alpha}}{\Gamma(\alpha)} \theta^{\alpha-1} \exp(-\beta\theta)}{\frac{(\beta+1)^{\alpha+k}}{\Gamma(\alpha+k)} \theta^{\alpha+k-1} \exp(-\beta\theta - \theta)}$$
(1.10)

12

$$= \frac{\beta^{\alpha}\Gamma(\alpha+k)}{(\beta+1)^{\alpha+k}\Gamma(\alpha)k!}$$

As we knew it should, the  $\theta$ 's disappeared from our computations, and we are left with what is called the *prior predictive* distribution for k, given the prior  $\theta \sim \text{Gamma}(\alpha, \beta)$ . If we instead start with a posterior, for example the posterior  $\theta \sim \text{Gamma}(\alpha + \sum_{i=1}^{n} k_i, \beta + n)$ , the formula above gives the *posterior predictive*. Replacing  $\alpha$  with  $\alpha + \sum_{i=1}^{n} \text{ and } \beta$  with  $\beta + n$  in Equation 1.10, we see that we get exactly the result of Equation 1.8.

The trick we just used to compute the predictive distribution can be used in all situations where you have conjugacy. Often the resulting predictive distribution turns out to be in a well-known family of distributions. In our case, the probability mass function for  $\pi(k)$  found in Equation 1.10 is actually a Negative Binomial distribution: A stochastic variable x taking on as possible values any positive integer has a Negative Binomial distribution if its probability mass function is given by

$$\pi(x \mid r, p) = \binom{x+r-1}{x} \cdot (1-p)^x p^r = \frac{\Gamma(x+r)}{\Gamma(x+1)\Gamma(r)} (1-p)^x p^r \qquad (1.11)$$

where r > 0 and  $p \in (0, 1)$  are parameters. Thus we see that the prior predictive density computed in Equation 1.10 is a Negative-Binomial $(\alpha, \beta/(1 + \beta))$  distribution.

## 1.4 Summary so far

Bayesian inference can be summarized in the following way: Let y represent a vector of random variables which you have observed, and let  $y_{new}$  represent a vector of random variables you would like to predict. There are then two steps:

- 1. Based on reasonable assumptions, create a stochastic model relating y and  $y_{new}$ .
- 2. Make predictions for  $y_{new}$  using the conditional distribution  $y_{new} \mid y$ , where y is fixed to its observed values.

Most often, the stochastic model is formulated using an additional stochastic variable  $\theta$ , a parameter or vector of parameters, so that y and  $y_{new}$  are conditionally independent given  $\theta$ , i.e., for the densities,

$$\pi(y_{new} \mid \theta, y) = \pi(y_{new} \mid \theta).$$

Then,

$$\pi(y_{new} \mid y) = \int_{\theta} \pi(y_{new}, \theta \mid y) \, d\theta = \int_{\theta} \pi(y_{new} \mid \theta) \pi(\theta \mid y) \, d\theta \tag{1.12}$$

and computation of  $\pi(y_{new} \mid y)$  is done by first computing  $\pi(\theta \mid y)$  and then computing the integral above.

To find  $\pi(\theta \mid y)$  one can generally formulate the joint distribution  $\pi(\theta, y) = \pi(y \mid \theta)\pi(\theta)$  and then use Bayes theorem:

$$\pi(\theta \mid y) = \frac{\pi(y \mid \theta)\pi(\theta)}{\pi(y)} \propto_{\theta} \pi(y \mid \theta)\pi(\theta)$$

The distributions  $\pi(\theta)$  and  $\pi(\theta \mid y)$  are called the *prior* and *posterior*, respectively.

In many cases, the data y is a random sample  $y = (y_1, \ldots, y_n)$ , where the  $y_i$  are conditionally independent given the parameter  $\theta$ , in other words,

$$\pi(y \mid \theta) = \prod_{i=1}^{n} \pi(y_i \mid \theta).$$

We get for the posterior

$$\pi(\theta \mid y) \propto_{\theta} \prod_{i=1}^{n} \pi(y_i \mid \theta) \pi(\theta),$$

and it can be obtained by stepwise updating the prior  $\pi(\theta)$ , using the data values  $y_i$  in any order, and using the posterior from one update as the prior for the next update.

A parametric family of probability distributions for a parameter  $\theta$  is *conjugate* to a likelihood  $\pi(y \mid \theta)$  if, when the prior is in the family, the posterior is also in the family. We have so far looked at two examples of conjugacy: The Beta-Binomial conjugacy and the Gamma-Poisson conjugacy. Whenever we have conjugacy, the *prior predictive density* 

$$\pi(y) = \int_{\theta} \pi(y,\theta) \, d\theta = \int_{\theta} \pi(y \mid \theta) \pi(\theta) \, d\theta$$

has a simple closed form, which may be computed using Equation 1.9. The *posterior predictive density* of Equation 1.12 can also be computed from the same equation, replacing the prior  $\pi(\theta)$  with the posterior  $\pi(\theta \mid y)$ . Explicitly,

$$\pi(y_{new} \mid y) = \frac{\pi(y_{new} \mid \theta)\pi(\theta \mid y)}{\pi(\theta \mid y_{new}, y)}.$$

## 1.5 Bayesian inference using discretization

We saw in the previous section that Bayesian inference requires computation of the posterior  $\pi(\theta \mid y)$  and the predictive distribution  $\pi(y_{new} \mid y)$  using (most often) Equation 1.12. However, in many practical applications, these computations cannot be done using conjugacy; there simply does not exist a conjugate prior to the likelihood  $\pi(y \mid \theta)$  one would like to use. In this and the next section we look at some simple alternative computational approaches.

We first turn to discretization, which can be a very good alternative if  $\theta$  has only one dimension (i.e., component) and  $\pi(\theta)$  is positive only within some bounded interval. Assume  $\theta_1, \theta_2, \ldots, \theta_k$  are equally-spaced values within this interval, so that the prior density  $\pi(\theta)$  can reasonably be approximated by a discrete distribution on  $\{\theta_1, \ldots, \theta_k\}$  specified by

$$a_i = \Pr(\theta = \theta_i) = \frac{\pi(\theta_i)}{\sum_{j=1}^k \pi(\theta_j)}$$

Writing  $b_i = \pi(y \mid \theta_i)$  for i = 1, ..., k we can then approximate the posterior with a discrete distribution on  $\{\theta_1, \ldots, \theta_k\}$  specified by

$$c_i = \Pr\left(\theta = \theta_i \mid y\right) = \frac{\pi(y \mid \theta = \theta_i) \Pr\left(\theta = \theta_i\right)}{\sum_{j=1}^k \pi(y \mid \theta = \theta_j) \Pr\left(\theta = \theta_j\right)} = \frac{a_i b_i}{\sum_{j=1}^k a_j b_j}$$

Finally, for a specific value of  $y_{new}$  we may approximate the predictive distribution as

$$\pi(y_{new} \mid y) = \int_{\theta} \pi(y_{new} \mid \theta) \pi(\theta \mid y) \, d\theta$$
$$\approx \sum_{i=1}^{k} \pi(y_{new} \mid \theta_i) \Pr(\theta = \theta_i \mid y) \, d\theta = \sum_{i=1}^{k} \pi(y_{new} \mid \theta_i) c_i.$$

How good this approximation is depends of course on how large k is, as well as the regularity of the functions involved.

#### Example

Assume

$$p \sim \text{Beta}(2.3, 4.1)$$
  
 $y \mid p \sim \text{Binomial}(17, p)$   
 $y_{new} \mid p \sim \text{Binomial}(3, p).$ 

Assume we would like to compute the probability  $\pi(y_{new} = 1)$  given that y = 4. Using the theory developed above, we get the posterior

$$p \mid (y = 4) \sim \text{Beta}(2.3 + 4, 4.1 + 17 - 4) = \text{Beta}(6.3, 17.1).$$

According to the results of Exercise 8 of Section 1.7, the predictive distribution is Beta-Binomial, so we get

$$\pi(y_{new} = 1 \mid y = 4) = \frac{B(6.3 + 1, 17.1 + 3 - 1)}{B(6.3, 17.1)} \binom{3}{1}$$
$$= \frac{\Gamma(7.3)\Gamma(20.1)\Gamma(6.3 + 17.1)}{\Gamma(7.3 + 20.1)\Gamma(6.3)\Gamma(17.1)} \cdot \frac{3!}{1!2!}$$
$$= 0.403364$$

The following R code approximates this result using discretization:

```
p <- seq(0, 1, length.out=20)
a <- dbeta(p, 2.3, 4.1)
b <- dbinom(4, 17, p)
c <- a*b/sum(a*b) #No need to divide a by its sum before this step
d <- dbinom(1, 3, p)
sum(c*d)</pre>
```

The code results in 0.4033704, which we see is a good approximation even if k is only 20. The advantage with the R code is of course that one may use any prior density on [0, 1], not just a Beta density.

## 1.6 Bayesian inference using numerical integration

Instead of discretizing one may apply numerical integration. After all, the answers we seek can be expressed as integrals:

$$\begin{aligned} \pi(y_{new} \mid y) &= \int_{\theta} \pi(y_{new} \mid \theta) \pi(\theta \mid y) \, d\theta = \int_{\theta} \pi(y_{new} \mid \theta) \frac{\pi(y \mid \theta) \pi(\theta)}{\int_{\theta} \pi(y \mid \theta) \pi(\theta) \, d\theta} \\ &= \frac{\int_{\theta} \pi(y_{new} \mid \theta) \pi(y \mid \theta) \pi(\theta) \, d\theta}{\int_{\theta} \pi(y \mid \theta) \pi(\theta) \, d\theta} \end{aligned}$$

## Example

For simplicity we continue with the example from the previous section. So assume we would like to compute, as above, the probability that  $y_{new} = 1$  given that y = 4. We get

$$\pi(y_{new} = 1 \mid y = 4) = \frac{\int_0^1 \text{Binomial}(1; 3, \theta) \text{Binomial}(4; 17, \theta) \text{Beta}(\theta; 2.3, 4.1) \, d\theta}{\int_0^1 \text{Binomial}(4; 17, \theta) \text{Beta}(\theta; 2.3, 4.1) \, d\theta}$$

and the R code

which produces the answer 0.403364, i.e., an even better approximation than with the discretization. As with discretization, the computations above can be done in principle with any densities, one is not limited to using conjugate priors.

Note the use of the Vectorize function above: The integrate function in R requires as input a function that not only computes the functional value, but given a vector of input values it should compute the vector of the corresponding functional values. The function Vectorize is a trick in R to convert any function into a vectorized one. In our particular example, the functions f1 and f2 are already vectorized, as you may check. However, this issue is easy to forget, so it is safest to include Vectorize as done in the code above.

Discretization or numerical integration works well in the simple example above, yielding accurate results. When variables are defined on unbounded intervals, one may need to make transformations before doing discretization. However, the biggest limitation to these methods is the dimension of the  $\theta$  vector.

Let's imagine that, to get some kind of reasonable accuracy when discretizing a real variable, you cannot use fewer than 10 values to represent it. A density in n dimensions will then need  $10^n$  points to represent it. For many problems nwill be higher than, say, 10, giving at least  $10^{10}$  gridpoints, which is unfeasible to handle. In fact we may very well want to study problems with millions of dimensions. Cleary, neither discretization nor numerical integration are then useful tools.

## 1.7 Exercises

1. A survey has been made about the type of living conditions and the political opinions of people in a city. Probabilities for observing each combination have been estimated and is listed in the following table

	Party A	Party B	Party C	Party D
Rental flat	0.11	0.03	0.08	0.01
Self-owned flat	0.09	0.01	0.14	0.03
House	0.13	0.04	0.09	0.24

If you observe that a person lives in a rental flat, what is the probability that the person votes for party B?

- 2. A disease is affecting 0.7% of the population. Initial diagnosis is done with a somewhat unreliable test. If a person is affected, the test will be positive with a 95% probability. However, if the person is not affected, there is still a 5% chance that the test is positive. Given that the test is positive, what is the probability that the person is affected?
- 3. Assume you are making repeated independent experiments with a probability of success  $\theta$  in each experiment. Initially, you make 12 experiments, of which 9 are successful.
  - a) Using a prior for  $\theta$  that is uniform on the interval [0, 1], what is the posterior for  $\theta$  given the results of the 12 experiments?
  - b) Assume now that you continue with doing 19 more experiments, of which 11 are successful. Given the combined information from all your 31 experiments, what is the posterior for  $\theta$ ?
  - c) Given all the information above, what is the probability for success in your thirty second experiment?
- 4. Refer to the coin-flipping example in the beginning of the chapter.
  - a) Write down the probability of observing  $y_H$  heads and  $y_T$  tails during a sequence with  $y_H + y_T$  coinflips.
  - b) Write down the probability of observing  $y_H$  heads and  $y_T$  tails followed by observing heads, during  $y_H + y_T + 1$  coinflips.
  - c) Write down the conditional probability of observing heads after having observed  $y_H$  heads and  $y_T$  tails. Verify that your result is the same as that in Equation 1.2.
- 5. Write an R program that computes the parameter  $\alpha$  for a Beta $(\alpha, \alpha)$  distribution which has 90% of its density in the interval [0.4, 0.6].
- 6. Assume  $y \sim \text{Negative-Binomial}(\alpha, p)$ , where  $\alpha$  is fixed and known and p is the unknown parameter. Prove that the Beta family is a conjugate prior family.
- 7. Assume  $y \sim \text{Normal}(\mu, \tau^{-1})$  where  $\mu$  is fixed and  $\tau$  is the unknown parameter. Prove that the Gamma family of distributions is a conjugate prior family.

#### 1. Basics of Bayesian inference

- 8. Find the formula for the probability mass function of the prior predictive distribution for the Beta Binomial conjugacy. Try to look up if this distribution has a standard name.
- 9. In the beginning of this chapter, we discussed an example with a loaded dice. We will now make some explicit computations for this example. We will use a parameter vector  $\theta = (\theta_1, \theta_2, \theta_3)$ , where each  $\theta_i \in (0, 1)$ . Given this parameter, and based on the geometry of standard dice, we model the probability of obtaining the value k with the dice as

$$\pi(k \mid \theta) = \begin{cases} \frac{1}{3}\theta_1 & k = 1\\ \frac{1}{3}\theta_2 & k = 2\\ \frac{1}{3}\theta_3 & k = 3\\ \frac{1}{3}(1 - \theta_3) & k = 4\\ \frac{1}{3}(1 - \theta_2) & k = 5\\ \frac{1}{3}(1 - \theta_1) & k = 6 \end{cases}$$

For the prior on  $\theta$ , we use

 $\pi(\theta) = \text{Beta}(\theta_1; 20, 20) \cdot \text{Beta}(\theta_2; 20, 20) \cdot \text{Beta}(\theta_3; 20, 20)$ 

- a) Find a formula for the posterior for  $\theta$  given a specific observed sequence of throws  $k_1, k_2, \ldots, k_n$ .
- b) Find the probabilities for each of the outcomes 1, 2, ..., 6 conditional on having observed the following sequence: 2, 4, 1, 6, 3, 6, 6, 3, 4, 2, 2.
- 10. Assume you have defined a likelihood function  $\pi(y \mid \theta)$  and are given a family of priors  $q_{\gamma}(\theta)$ , parametrized by a vector  $\gamma$  of parameters, with  $\gamma \in \Omega$ . Assume now that this family is conjugate, so that, if the prior is  $q_{\gamma}(\theta)$  for some  $\gamma \in \Omega$ , then the posterior  $\theta \mid y$  has density  $q_{f(\gamma)}(\theta)$  for some other  $f(\gamma) \in \Omega$ .
  - a) Fix an integer k > 1 and describe a new family of priors as consisting of all densities

$$r(\theta) = \sum_{i=1}^{k} \lambda_i q_{\gamma_i}(\theta)$$

where  $\lambda_1, \ldots, \lambda_k$  are nonnegative real numbers summing to 1, and for all  $i, \gamma_i \in \Omega$ . Prove that this family is a conjugate family. Derive explicit formulas for the posterior given a prior like the one above.

- b) Compute an explicit formula for the prior predictive distribution in this case.
- c) Can you imagine an application where using this kind of *mixture prior* as a model could be advantageous?
- 11. An example of conjucacy is the Normal Normal conjugacy; see the Appendix in Chapter C. Assume that  $x \mid \mu \sim \text{Normal}(\mu, \tau^{-1})$  and  $\mu \sim \text{Normal}(\mu_0, \tau_0^{-1})$  for fixed and known  $\tau$  and  $\tau_0$ . Then it can be shown that the prior predictive distribution for x is normal. Use this fact, together with what is called in Dobrow the Law of Total Expectation and the Law of Total Variance, to find the parameters of this normal distribution.

## CHAPTER 2

## Hidden Markov Models

Let us start with reviewing Exercise 2.20 of Dobrow: We let  $X_0, X_1, \ldots$  be a Markov chain with transition matrix

$$P = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ p & 1 - p & 0 \end{pmatrix}$$

for some 0 . Let g be a function defined by

$$g(x) = \begin{cases} 0, & \text{if } x = 1\\ 1, & \text{if } x = 2, 3 \end{cases}$$

If we let  $Y_n = g(X_n)$  for  $n \ge 0$ , we might ask if also  $Y_0, Y_1, \ldots$  is a Markov chain. However, it is not: If  $Y_i = 1$  and  $X_i = 2$ , the probability that  $Y_{i+1} = 0$  is zero, as we necessarily have  $X_{i+1} = 3$ . However, if  $Y_i = 1$  and  $X_i = 3$ , then the probability that  $Y_{i+1} = 0$  is p. Thus, given  $Y_i$ , the value of  $Y_{i+1}$  would depend on  $X_i$ , and thus on  $Y_j$  for j < i, as these  $Y_j$  and  $X_i$  are not independent. In other words,  $Y_0, Y_1, \ldots$  is not a Markov chain.

This very simple example illustrates a situation which can be found in many types of data indexed with an integer: The sequence  $Y_0, Y_1, \ldots$  of the data itself is not a Markov chain, as each  $Y_i$  does not contain all the underlying information that we imagine is used to "generate" a new state  $Y_{i+1}$  in the chain. However, there may exist, as in the example above, some underlying variables  $X_i$  that encode the underlying information, in such a way that  $X_0, X_1, \ldots$  is a Markov chain.

The situation can be illustrated with Figure 2.1. Note that the graph does not illustrate transitions between states of a Markov chain, but instead the dependencies between the random variables  $x_0, x_1, \ldots, x_T$  and  $y_0, y_1, \ldots, y_T$ . We assume  $x_0, x_1, \ldots$  is a Markov chain. The figure also illustrates the assumption that  $y_i$  depends on  $x_i$ , but given  $x_i$ , it is conditionally independent of all the other variables. In some situations,  $y_i$  may also depend on  $y_{i-1}$ . There is then an arrow from  $y_{i-1}$  to  $y_i$  for all *i*. This will be the case in Example 2 below. However, for simplicity, we will focus on the theory for the case where  $y_i$  does not depend on  $y_{i-1}$  when  $x_i$  is given. The extension to the case where there is such a dependency is generally not so difficult.

The model described above is called a *Hidden Markov Model*, or HMM. Below, we will look at some examples of situations where such models may be useful. Afterwards, we will summarize the inference questions we might have



Figure 2.1: A hidden Markov model. Note that this is not a transition graph relating states of a Markov chain, but instead a graph showing the dependencies between the random variables  $x_0, x_1, \ldots, x_T$  and  $y_0, y_1, \ldots, y_T$ .

in such situations. In the next chapter, we will look at how these inference questions can be answered using Bayesian inference.

## 2.1 Examples of Hidden Markov Models

## Example 1

Each day i, a pharmacy sells on-line  $Y_i$  bottles of cough-medicine. We assume the number of bottles sold is Poisson distributed based on an underlying expected demand  $X_i$  in the area of the pharmacy, so that

$$Y_i \sim \text{Poisson}(X_i)$$

The underlying demand  $X_i$  depends on whether there is a flu-infection in the area or not: We simplify and say that  $X_i = 30$  if there is an infection present and  $X_i = 10$  if not. We then model the presence of the infection with a Markov chain. Specifically, we assume  $X_0, X_1, \ldots$  is a two-state Markov chain with possible values 10 and 30 and transition matrix

$$P = \begin{bmatrix} 0.95 & 0.05\\ 0.2 & 0.8 \end{bmatrix}.$$

Figure 2.2 shows a simulation of this process over 60 days. The full line represents the underlying demand for cough-medicine, which is assumed to be 30 if there is a flu-infection in the area and 10 if not. What is actually observed are the daily sales, which are assumed to be Poisson distributed based on the expected demand. From the figure, we see that the sales go up when there is an infection in the area. Thus, observing the sales might be used as a way to detect such infection. However, as the actual sales is a noisy observation of the underlying demand, it is not clear-cut to conclude from this data exactly when the infection is present. And, of course, in a more realistic model, this would be much more difficult.

## Example 2

As our second and more realistic example, we look at the problem of finding promotor regions for genes in vertebrate genomes. The example is based on



Figure 2.2: A simulation from the flu model. The full line represents the underlying expected demand for cough-medicine, based on whether there is a flu-infection in the area or not. The dots represent the observed actual sales of the medicine.

[Axe10], which is an excellent reference for both theory and examples of Hidden Markov Models, in particular when used in genetics.

DNA consists of the four nucleotides, or bases, A, C, G, and T. One may try to model DNA strands consisting of sequences of such bases as a Markov chain, with the bases being the states of the Markov chain and each position on the DNA corresponding to a "time" step i. Of course, the bases are not randomly distributed along the DNA, so one might expect to capture some (small part) of the structure using a Markov chain model.

In fact, the two matrices below represent transition matrices for two types of genetic strands:

$$P_{+} = \begin{bmatrix} 0.180 & 0.274 & 0.426 & 0.120 \\ 0.171 & 0.368 & 0.274 & 0.188 \\ 0.161 & 0.339 & 0.375 & 0.125 \\ 0.079 & 0.355 & 0.384 & 0.182 \end{bmatrix}, P_{-} = \begin{bmatrix} 0.300 & 0.205 & 0.285 & 0.210 \\ 0.322 & 0.298 & 0.078 & 0.302 \\ 0.248 & 0.246 & 0.298 & 0.208 \\ 0.177 & 0.239 & 0.292 & 0.292 \end{bmatrix}$$

The transition matrix  $P_+$  has been found by studying DNA from so-called CpG islands: These are places on the genome where the usual tendency for CG-pairs to mutate to different pairs is dampened by evolution, as the actual

sequences are important for the functioning of the gene regulation processes. The transition matrix  $P_{-}$  is based on data from other parts of the genome.

One may now try to use the differences between these two transition matrices to detect whether some sequence is a CpG island or not. For a given sequence, one might try out the two models to see which one fits best. However, a more realistic application is that one is given the entire DNA sequence of a chromosome, and one would like to detect which parts of that sequence are CpG islands. Just like in Example 1, it might then not be entirely obvious how best to determine the boundaries between CpG islands and non-CpG islands.

This is where the theory of Hidden Markov Models can be used: Let  $X_i$  be 1 or 0 depending on whether position i of the genome is in a CpG island or not. Then let  $Y_i$  be A, C, G, or T, depending on the observed nucleotide at i. The probabilities of observing possible value for  $Y_i$  would depend on  $Y_{i-1}$ , where one would use the  $P_+$  transition matrix if  $X_i = 1$  and  $P_-$  otherwise. In order to complete the model, we would need transition probabilities for a Markov chain modelling the  $X_i$ . Such probabilities could be found by considering data where the CpG islands are known.

Note how in this example  $Y_i$  depends on both  $X_i$  and  $Y_{i-1}$ . Thus, a proper illustration of this model would look like Figure 1.1, but with arrows added between the y variables.

## 2.2 Inference questions for HMMs

Let us try to get an overview of the types of inference questions that might be considered for HMMs. We can divide them into two types: Cases where the model parameters are known, and cases where they are not.

### Known model parameters

Imagine that the transition matrix for the Markov chain  $X_0, X_1, \ldots$  is known, as well as its initial distribution. Also, imagine that all the conditional probability distributions  $\pi(Y_i \mid X_i)$  (or, in the more general case, the probability distributions  $\pi(Y_i \mid X_i, Y_{i-1})$ ) are known. Finally, assume that we have observed as data  $Y_0, Y_1, \ldots, Y_T$ . What is unknown is then the values of the hidden variables  $X_0, X_1, \ldots, X_T$ .

From a Bayesian perspective, we would like to known the posterior distribution for  $X_0, \ldots, X_T$  given the data  $Y_0, \ldots, Y_T$ . But as this is a (potentially long) sequence of random variables, we need in practice to be more specific about what we want to compute. Three options are common:

- 1. We find a specific sequence of values  $x_0, x_1, \ldots, x_T$  which maximizes the posterior distribution. In example 1 above, it would mean to find our best possible guess for the time points when the flu infection appears and disappears in the area. In example 2, it would produce our best possible guess the locations of the CpG islands.
- 2. For each single step i, we find the marginal posterior distribution for  $X_i$  given the data. In other words, for every position, we find the most likely value for  $X_i$ , but we can also consider which values for  $X_i$  are almost as likely. For example, in example 2 above, we could compute the probability that a specific location is in a CpG island.

3. We might generate sample sequences  $x_0, x_1, \ldots, x_T$  from the posterior distribution given the data. For example, in example 2 above, we might use such samples to estimate the average lengths of the CpG islands that we predict.

In the next chapter, we briefly discuss some methods for the types of inference above.

## Unknown model parameters

In our examples, we have simply presented fixed transition matrices. But model parameters, such as transition matrices, ultimately come from data. In HMMs, the model parameters consist of the transition matrix, the initial distribution on  $X_0$ , and the parameters of the distributions  $\pi(Y_i \mid X_i)$  (or  $\pi(Y_i \mid X_i, Y_{i-1})$ ). We need to learn about these from data.

We must then either make some very strong assumptions about the hidden Markov chain, its states, and their relationship with observed values (as in Example 1 above), or we need data that contain observations of both the hidden chain values  $X_i$  and the values  $Y_i$  (as in Example 2 above).

We will consider some methods for Bayesian inference for HMM model parameters in the next Chapter.

## CHAPTER 3

# Some basic inference for Markov chains and HMMs

Assume you would like to use a Markov chain  $X_0, X_1, \ldots$ , as your model in an applied setting where some data is available. The parameters of the Markov chain are the transition matrix P and the probability distribution p on the initial state  $X_0$ . You would then like to learn about these parameters from the data. This data could take many forms, for example, you could have observed only some specific selection of the variables  $X_i$ . In this chapter, we will assume that you have observed the whole sequence  $X_0, X_1, \ldots, X_n$  up to some number n. After learning about the parameters, you can use this to predict further steps  $X_{n+1}, X_{n+2}, \ldots$  of the chain.

We will also consider inference for Hidden Markov Models (HMMs). There is then a wider range of what your data could consist of. In this chapter, we will consider a situation where the model parameters of the HMM are known and only the values of the hidden states  $X_0, X_1, \ldots, X_n$  are unknown. But we will also consider a situation where model parameters are learned from data.

## 3.1 The Multinomial Dirichlet conjugacy

The Multinomial distribution counts the number of outcomes in each of k possible classes when n independent trials are performed and the probability of ending up in each of the classes is given by the probability vector  $p = (p_1, \ldots, p_k)$ . (Recall that a probability vector p of length k is a vector of non-negative real numbers such that  $\sum_{i=1}^{k} p_i = 1$ .) In other words, a vector  $x = (x_1, \ldots, x_k)$  of non-negative integers has a Multinomial distribution with parameters n and p if  $\sum_{i=1}^{k} x_i = n$  and the probability mass function is given by

$$\pi(x \mid n, p) = \binom{n}{x_1 \ x_2 \ x_3 \ \dots \ x_k} p_1^{x_1} p_2^{x_2} \cdots p_k^{x_k}.$$

Recall that the *multinomial coefficient* above is given by

$$\binom{n}{x_1 \ x_2 \ x_3 \ \dots \ x_k} = \frac{n!}{x_1! x_2! \dots x_k!}.$$

Note that the Multinomial distribution with k = 2 can be identified with the Binomial distribution.

#### Some basic inference for Markov chains and HMMs

A vector  $\theta = (\theta_1, \dots, \theta_k)$  of non-negative numbers satisfying  $\sum_{i=1}^k \theta_i = 1$  has a *Dirichlet* distribution with parameter vector  $\alpha = (\alpha_1, \dots, \alpha_k)$ , with each  $\alpha_i > 0$ , if it has probability density function

$$\pi(\theta \mid \alpha) = \frac{\Gamma(\alpha_1 + \alpha_2 + \dots + \alpha_k)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\dots\Gamma(\alpha_k)} \theta_1^{\alpha_1 - 1} \theta_2^{\alpha_2 - 1} \dots \theta_k^{\alpha_k - 1}$$

Note that the Dirichlet distribution with k = 2 can be identified with the Beta distribution.

As the Beta family is conjugate to the Binomal likelihood, it is natural to check if the Dirichlet family is conjugate to the Multinomial likelihood. So assume  $\theta$  has the prior  $\theta \mid \alpha \sim \text{Dirichlet}(\alpha)$  for some  $\alpha$ , and assume we have the Multinomial likelihood  $x \mid n, \theta \sim \text{Multinomial}(n, \theta)$ . Bayes formula gives

$$\pi(\theta \mid x) \propto_{\theta} \pi(x \mid \theta)\pi(\theta)$$
  
$$\propto_{\theta} \theta_{1}^{x_{1}}\theta_{2}^{x_{2}}\dots\theta_{k}^{x_{k}}\theta_{1}^{\alpha_{1}-1}\theta_{2}^{\alpha_{2}-1}\dots\theta_{k}^{\alpha_{k}-1}$$
  
$$= \theta_{1}^{\alpha_{1}+x_{1}-1}\theta_{2}^{\alpha_{2}+x_{2}-1}\dots\theta_{k}^{\alpha_{k}+x_{k}-1}$$

from which we deduce that

$$\theta \mid x \sim \text{Dirichlet}(\alpha_1 + x_1, \alpha_2 + x_2, \dots, \alpha_k + x_k)$$

and we have shown conjugacy. For the predictive distribution we get

$$\pi(x) = \frac{\pi(x \mid \theta)\pi(\theta)}{\pi(\theta \mid x)}$$

$$= \frac{\text{Multinomial}(x; n, \theta) \text{Dirichlet}(\theta; \alpha)}{\text{Dirichlet}(\theta; \alpha + x)}$$

$$= \frac{\frac{n!}{x_1! \dots x_k!} \theta_1^{x_1} \dots \theta_k^{x_k} \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \theta_1^{\alpha_1 - 1} \dots \theta_k^{\alpha_k - 1}}{\frac{\Gamma(\alpha_1 + \dots + \alpha_k + x_1 + \dots + x_k)}{\Gamma(\alpha_1 + x_1) \dots \Gamma(\alpha_k + x_k)}}$$

$$= \frac{n!}{x_1! \dots x_k!} \cdot \frac{\Gamma(\alpha_1 + x_1)}{\Gamma(\alpha_1)} \dots \frac{\Gamma(\alpha_k + x_k)}{\Gamma(\alpha_k)} \cdot \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1 + \dots + \alpha_k + x_1 + \dots + x_k)}$$
(3.1)

which is a generalization of the Beta Binomial distribution.

See Exercises 1 and 2 for more about these distributions.

## 3.2 Inference for time-homogeneous Markov chains with finite state space

Consider a time-homogeneous markov chain consisting of random variables  $X_0, X_1, \ldots$ , with finite state space S with k elements. The parameters of this model are  $p_0$ , the probability vector describing the distribution on  $X_0$ , and P, the transition matrix of the chain. Assume first that  $p_0$  is known while P is unknown, and assume we would like to learn about P using a sequence of observations  $x_0, x_1, \ldots, x_n$  for the variables  $X_0, \ldots, X_n$ . Let us write  $P_i$  for the *i*'th row of P: We now consider it a random variable, in fact, a probability vector with non-negative entries summing to 1. The probability of the data  $x_0, \ldots, x_n$  for a fixed P is

$$\pi(x_0, \dots, x_n \mid P) = \pi(x_0) \prod_{r=1}^n \pi(x_r \mid x_{r-1}, P)$$

$$= \pi(x_0) \prod_{r=1}^{n} P_{x_{r-1},x_r}$$
$$= \pi(x_0) \prod_{i=1}^{k} \prod_{j=1}^{k} (P_{ij})^{c_{ij}}$$

where  $c_{ij}$  is the count of times the chain  $x_0, x_1, \ldots, x_n$  transits from state *i* to state *j*. We now define a prior on *P* with

$$\pi(P) = \prod_{i=1}^{k} \text{Dirichlet}(P_i; \alpha_i)$$

where  $\alpha_i = (\alpha_{i1}, \ldots, \alpha_{ik})$  is a vector of parameters for  $i = 1, \ldots, k$ . For the posterior, we get

$$\pi(P \mid x_0, \dots, x_n) \propto_P \pi(x_1, \dots, x_n \mid P) \pi(P)$$
$$\propto_P \prod_{i=1}^k \prod_{j=1}^k (P_{ij})^{c_{ij}} \prod_{i=1}^k \prod_{j=1}^k (P_{ij})^{\alpha_{ij}-1}$$
$$\propto_P \prod_{i=1}^k \prod_{j=1}^k (P_{ij})^{\alpha_{ij}+c_{ij}-1}$$

from which we read off that

$$\pi(P \mid x_0, \dots, x_n) = \prod_{i=1}^n \text{Dirichlet}(P_i; \alpha_i + c_i)$$

where  $c_i = (c_{i1}, ..., c_{ik}).$ 

#### Example

Assume you have a Markov chain with three possible states, and that a sequence of values  $x_0, \ldots, x_{20}$  have been observed. Assume the counts  $c_{ij}$  of transitions from state *i* to state *j* are given by the following table:

	1	2	3
1	3	4	1
2	3	3	0
3	2	0	4

A classical inference approach might try to use the observed frequencies of transitions as the values in the transition matrix P. But we see that some values of  $P_{ij}$  would then become zero. We have seen that whether the entries of the transition matrix are zero or positive can have a decisive influence on the properties of the Markov chain, and it may seem rash to conclude that some transitions have probability zero simply because they have not been observed in a short sequence of the chain. Even more fundamentally, with other data, some states might not have been visited at all. All counts in the vector  $c_i$  would then be zero, and it would be impossible to compute frequencies summing to one from such a vector. Note that even if such problems might be overcome with more data, the amount of data needed for frequency estimates increases dramatically with the number s of states in the state space.

In practice, to get a useful result for P, one may need to use more information than that available in the counts  $c_{ij}$ . In the Bayesian inference above, such information is provided in the parameters  $\alpha_i$  of the prior for P. Note that if  $x \sim \text{Dirichlet}(\lambda)$ , then  $\mathbf{E}[x] = \lambda / \sum_{i=1}^k \lambda_i$ . Thus, in our case<sup>1</sup>,

$$E(P_i \mid x_0, \dots, x_n) = \frac{\alpha_i + c_i}{\alpha_{i1} + \dots + \alpha_{ik} + c_{i1} + \dots + c_{ik}}$$

Thus, as long as we use parameters  $\alpha_{ij} > 0$ , all posterior expectations of values in *P* will be nonzero. In many situations, a reasonable choice may be  $\alpha_{ij} = 1$ for all i, j, leading to

$$E(P_i \mid x_0, \dots, x_n) = \frac{(1, 1, \dots, 1) + c_i}{k + c_{i1} + \dots + c_{ik}}$$

In the example above, we get for example

$$E(P_1 \mid x_0, \dots, x_{20}) = \frac{(1, 1, 1) + (3, 4, 1)}{3 + 3 + 4 + 1}$$

and the posterior expectation for P would be the matrix

$$E(P \mid x_0, \dots, x_{20}) = \begin{bmatrix} 4/11 & 5/11 & 2/11 \\ 4/9 & 4/9 & 1/9 \\ 3/9 & 1/9 & 5/9 \end{bmatrix}.$$
 (3.2)

The values  $\alpha_{ij}$  are sometimes called *pseudocounts*; however, they do not need to be integers.

The situation where the distribution  $p_0$  for  $X_0$  is unknown can be handled in a similar way. Note however that if we have observed only one sequence  $x_0, \ldots, x_n$ , only  $x_0$  informs us about  $p_0$ , so unless we have observed a number of sequences from the chain, the distribution for  $p_0$  will be more or less determined by the prior.

## Prediction

Let us assume we would like to predict the observation  $x_{n+1}$  of  $X_{n+1}$  based on the sequence  $x_0, \ldots, x_n$ . We can write

$$\pi(x_{n+1} \mid x_0, \dots, x_n) = \int \pi(x_{n+1} \mid x_n, P) \pi(P \mid x_0 \dots, x_n) dP$$
  
=  $\int P_{x_n, x_{n+1}} \pi(P_{x_n} \mid x_0 \dots, x_n) dP_{x_n}$  (3.3)

If we use a prior and compute the posterior for  $P_{x_k}$  as in the last subsection, we get

 $P_{x_n} \mid x_0, \ldots, x_n \sim \text{Dirchlet}(\alpha_{x_n} + c_{x_n}).$ 

 $<sup>^1\</sup>mathrm{Make}$  sure you interpret this Equation correctly: The numerator is a vector and the denominator is a number
According to Equation 3.3, the predictive distribution is given as the Expectation vector of the posterior Dirichlet distribution:

$$\pi(x_{n+1} \mid x_0, \dots, x_n) = \frac{\alpha_{x_n} + c_{x_n}}{\alpha_{x_n, 1} + \dots + \alpha_{x_n, k} + c_{x_n, 1} + \dots + c_{x_n, k}}.$$

Consider the example of the previous section, where  $x_0, \ldots, x_{20}$  were observed; assume that  $x_{20} = 2$ . We showed in Equation 3.2 that the expectation of  $P_2$ , the second row of the transition matrix, was (4/9, 4/9, 1/9). Thus we get that

$$\pi(x_{21} \mid x_0, \dots, x_{20}) = (4/9, 4/9, 1/9).$$
(3.4)

Let us also consider the prediction of a whole chain of observations  $x_{n+1}, \ldots, x_{n+r}$  based on the sequence  $x_0, \ldots, x_n$ . In the same way as above, we get

$$\pi(x_{n+1}, \dots, x_{n+r} \mid x_0, \dots, x_n) = \int \left[\prod_{s=1}^r \pi(x_{n+s} \mid x_{n+s-1}, P)\right] \pi(P \mid x_0 \dots, x_n) dP$$
$$= \int \left[\prod_{s=1}^r P_{x_{n+s}, x_{n+s-1}}\right] \pi(P \mid x_0 \dots, x_n) dP.$$
(3.5)

When the posterior for P is a product of Dirichlet distributions, it is in fact possible to compute the value of this integral, in a similar way as in computations for the predictive distribution for the Multinomial Dirichlet conjugate pair given in Equation 3.1. See Exercise 4 for concrete computations.

## Extensions

The Dirichlet priors we have considered here assumes that the transition matrix is positive. However, there may be situations where certain transitions may be ruled out apriori. In such cases, an alternative is to use Dirichlet distributions on the parameters in each line that could be non-zero. See Exercise 3 for an example.

In other situations, the assumption used above that the lines of the transition matrix are a priori independent may be unreasonable. In such cases, a prior reflecting this situation could be used.

In yet other situations, it may be known that the Markov chain is time reversible, so that inference about its parameters should be done under this restriction. A possibility is then to represent the Markov chain as a random walk on a weighted undirected graph, and infer the weights from data. It is even possible to use a conjugate analysis in this case.

The above discussion on predictions can also be extended in many directions. For example, making predictions for long stretches of a Markov chain may best be done by first simulating its transition matrix from the posterior and then continuing the Markov chain simulating with this transition matrix. Finally, one may study how a stationary distribution derived from a transition matrix changes when taking into account the posterior uncertainty in this transition matrix.

## 3.3 Inference for HMMs, known model parameters\*

THIS SECTION IS CURRENTLY NOT PART OF THE CURRICULUM FOR MVE550.

In Section 2.2 we discussed of possible goals for inference for HMMs. We will in this section focus on the following situation: The distributions  $\pi(X_{i+1} | X_i)$ and  $\pi(Y_i | X_i)$  are known for all *i*, and where the initial distribution  $\pi(X_0)$  is known. We also assume that we have observed the sequence  $y_0, y_1, \ldots, y_T$  for the variables  $Y_0, \ldots, Y_T$ , and that our main objective is to find the marginal posterior distribution for each of the variables  $X_0, \ldots, X_T$ . These posterior distributions may be of interest in themselves or one may use the posterior distributions to make other predictions of interest.

## The Forward algorithm

The objective is to compute, and store, for i = 0, 1, ..., T, the distributions  $\pi(X_i \mid Y_0, ..., Y_i)$ . This is done recursively, starting with i = 0, and at each step using the results of the previous step.

We first compute  $\pi(X_0 \mid Y_0)$  using Bayes formula:

$$\pi(X_0 \mid Y_0) = \frac{\pi(Y_0 \mid X_0)\pi(X_0)}{\pi(Y_0)} \propto_{X_0} \pi(Y_0 \mid X_0)\pi(X_0)$$

Then, assuming that we have computed and (somehow) stored  $\pi(X_i | Y_0 \dots, Y_i)$ we compute  $\pi(X_{i+1} | Y_0 \dots, Y_{i+1})$  again using Bayes formula:

$$\pi(X_{i+1} \mid Y_0, \dots, Y_{i+1}) \propto_{X_{i+1}} \pi(Y_{i+1} \mid X_{i+1}, Y_0, \dots, Y_i) \pi(X_{i+1} \mid Y_0, \dots, Y_i)$$
  
=  $\pi(Y_{i+1} \mid X_{i+1}) \int \pi(X_{i+1} \mid X_i) \pi(X_i \mid Y_0, \dots, Y_i) dX_i$ 

The details of how these computations are done, and how the results are stored, depend on the particular types of distributions involved. In an alternative description of the Forward algorithm, one recursively computes  $\pi(X_i | Y_0, \ldots, Y_{i-1})$  instead of  $\pi(X_i | Y_0, \ldots, Y_i)$ . The same ideas as above are used and the same computations are done; they are simply subdivided in a slightly different manner.

#### The Backward algorithm

The objective now is to compute and store, for  $i = T, \ldots, 0$ , the probabilities  $\pi(Y_{i+1}, \ldots, Y_T \mid X_i)$ . Note that when i = T, this expression is not really meaningful; however we will interpret  $\pi(Y_{i+1}, \ldots, Y_T \mid X_i)$  for i = T as a function equal to 1 for all values of  $X_T$ . Starting with this function, we go backwards, stepwise decreasing the index i, and computing  $\pi(Y_{i+1}, \ldots, Y_T \mid X_i)$  in terms of  $\pi(Y_{i+2}, \ldots, Y_T \mid X_{i+1})$ . We can to this by averaging out over  $X_{i+1}$ :

$$\pi(Y_{i+1}, \dots, Y_T \mid X_i) = \int \pi(Y_{i+1}, \dots, Y_T, X_{i+1} \mid X_i) \, dX_{i+1}$$
  
= 
$$\int \pi(Y_{i+1} \mid X_{i+1}) \pi(Y_{i+2}, \dots, Y_T \mid X_{i+1}) \pi(X_{i+1} \mid X_i) \, dX_{i+1}$$

The details of how these computations are done, and how the results are stored, depend on the particular types of distributions involved. In an alternative description of the Backward algorithm, one recursively computes  $\pi(Y_i \ldots, Y_T \mid X_i)$  instead of  $\pi(Y_{i+1}, \ldots, Y_T \mid X_i)$ . the same ideas as above are used and the same computations are done; they are simply subdivided in a slightly different manner.

### The Forward Backward algorithm

There are several interesting ways of putting together the two algorithms above. Let us for example assume that we would like to compute, for all i = 0, ..., T, the marginal posterior distributions  $\pi(X_i | Y_0, ..., Y_T)$ . We can do this by using (surprise!) Bayes formula:

$$\pi(X_i \mid Y_0 \dots, Y_T) \propto_{X_i} \pi(Y_{i+1}, \dots, Y_T \mid X_i, Y_0, \dots, Y_i) \pi(X_i \mid Y_0, \dots, Y_i)$$
  
=  $\pi(Y_{i+1}, \dots, Y_T \mid X_i) \pi(X_i \mid Y_0, \dots, Y_i)$ 

The distributions in the last line can be computed with the Backward and Forward algorithms, respectively.

Instead of finding the marginal distributions above, one might be interested in the joint distribution of all  $X_0, \ldots, X_T$  given the observed values for  $Y_0, \ldots, Y_T$ . As this is a high-dimensional distribution, it is easier to focus on obtaining a sequence  $x_0, x_1, \ldots, x_T$  which is a sample from this distribution. This can be done, for example, as follows: First, draw  $x_0$  from  $\pi(X_0 | Y_0, \ldots, Y_T)$  found as above. Then, for  $i = 1, 2, \ldots, T$ , draw  $x_i$  according to the density

$$\pi(X_i \mid Y_0, \dots, Y_T, X_0 = x_0, \dots, X_{i-1} = x_{i-1})$$
  

$$\propto_{X_i} \quad \pi(Y_i, \dots, Y_T \mid X_i, Y_0, \dots, Y_{i-1}, X_0 = x_0, \dots, X_{i-1} = x_{i-1})$$
  

$$\cdot \pi(X_i \mid Y_0, \dots, Y_{i-1}, X_0 = x_0, \dots, X_{i-1} = x_{i-1})$$
  

$$= \quad \pi(Y_i \mid X_i) \pi(Y_{i+1}, \dots, Y_T \mid X_i) \pi(X_i \mid X_{i-1} = x_{i-1})$$

As always, the details depend on the types of distributions involved.

#### Implementation when the state space is finite

Let us now assume the state space for X is finite, with the possible values  $1, \ldots, s$ . The Forward algorithm can be implemented as follows: To compute  $\pi(X_0 \mid Y_0) \propto_{X_0} \pi(Y_0 \mid X_0)\pi(X_0)$ , compute  $\pi(X_0)$  and  $\pi(Y_0 \mid X_0)$  for all s possible values of  $X_0$ , to obtain two vectors of length s. Multiply these two vectors termwise to obtain a new vector of length s, and divide by its sum to obtain a vector of length s representing  $\pi(X_0 \mid Y_0)$ .

For the recursive part of the computation, let  $v_i$  denote the probability vector representing  $\pi(X_i | Y_0 \dots, Y_i)$ . Then  $v_i P$  is the probability vector representing  $\int \pi(X_i + 1 | X_i) \pi(X_i | Y_0, \dots, Y_i) dX_i$  as a function of  $X_{i+1}$ . Multiplying termwise with the probability vector representing  $\pi(Y_{i+1} | X_{i+1})$  as a function of  $X_{i+1}$  and normalizing so that the sum becomes 1, we get the probability vector representing  $\pi(X_{i+1} | Y_0, \dots, Y_{i+1})$ .

The Backward algorithm is implemented similarly. See the exercise which concerns writing an R implementation of these algorithms.

## 3.4 Inference for HMMs, unknown model parameters

Assume  $X_0, X_1, \ldots, X_s, \ldots$ , is a Markov chain with a discrete state space, transition matrix P, and probability distribution p on the initial state  $X_0$ . Assume also  $Y_0, Y_1, \ldots, Y_s, \ldots$ , are discrete random variables so that

 $\pi(Y_s \mid X_0, X_1, \dots, Y_0, \dots, Y_{s-1}, Y_{s+1}, \dots) = \pi(Y_s \mid X_s).$ 

Assume further that  $\Pr(Y_s = j \mid X_s = i) = Q_{ij}$  is independent of s, so there is a single matrix Q describing the dependence of  $Y_s$  on  $X_s$ . Finally, assume one has observed  $X_0, \ldots, X_n$  and  $Y_0, \ldots, Y_n$ . How can we learn about the parameters P, p and Q of our model?

This will depend on what prior distribution we use for P, p, and Q. Generally, one will use independent priors for these. Then, the learning for P and p will be done as for any Markov chain  $X_0, \ldots, X_n$ , see Section 3.2. The learning for Q will be based on the n + 1 observed pairs  $(X_0, Y_0), (X_1, Y_1), \ldots, (X_n, Y_n)$ . In the example below, we try out two different priors for Q, to illustrate how this choice influences results.

#### Example

Assume the following values have been observed for  $X_0, \ldots, X_{20}$  and  $Y_0, \ldots, Y_{20}$ :

X	3	3	3	1	2	1	1	2	2	1	1	3	3	3	1	1	2	2	1	2	2
Y	1	4	3	2	3	2	1	1	4	1	1	3	3	4	0	0	3	0	0	2	2

We assume the X variable has possible values 1, 2, 3. The counts of transitions are exactly the same as those of the Example of Section 3.2, and we can learn about the transition matrix P in exactly the same way as in that section.

Let us first assume that Y can only take on the values 0, 1, 2, 3, 4, and use as a prior for Q a product of Dirichlet distributions. Specifically,

$$\pi(Q) = \prod_{i=1}^{3} \text{Dirichlet}(Q_i; \beta_i)$$

where  $Q_i$  is the i'th row of the Q matrix, and the vector  $\beta_i$  is the corresponding set of *pseudocounts* for the transitions from state *i*. The counts of transitions from the possible values 1, 2, 3 for X to the possible values 0, 1, 2, 3, 4 for Y are given in the following table:

	0	1	2	3	4
1	3	3	2	0	0
2	1	1	2	2	1
3	0	1	0	3	2

Just like in Section 3.2 we get that the posterior for Q is also a product of Dirichlet distributions:

$$\pi(Q \mid \text{data}) = \prod_{i=1}^{3} \text{Dirichlet}(Q_i; \beta_i + d_i)$$

where  $d_i$  is the vector of counts of transitions from state *i*. So, for example,  $d_1 = (3, 3, 2, 0, 0)$  and  $d_2 = (1, 1, 2, 2, 1)$ . Setting all the pseudocounts equal to 1, we get, explicitly,

 $\pi(Q \mid \text{data}) = \text{Dirchlet}(Q_1; 4, 4, 3, 1, 1) \text{Dirchlet}(Q_2; 2, 2, 3, 3, 2) \text{Dirichlet}(Q_3; 1, 2, 1, 4, 3).$ 

Computing expectations as in Section 3.2, if a state  $X_i$  for i > 20 has value 2, the probabilities for  $Y_i$  are given by the vector (2/12, 2/12, 3/12, 3/12, 2/12), i.e.,

### [1] 0.1666667 0.1666667 0.2500000 0.2500000 0.1666667

More generally we get the posterior expectation

$$\mathbf{E}(Q) = \begin{bmatrix} 4/13 & 4/13 & 3/13 & 1/13 & 1/13 \\ 2/12 & 2/12 & 3/12 & 3/12 & 2/12 \\ 1/11 & 2/11 & 1/11 & 4/11 & 3/11 \end{bmatrix}.$$

We can make a prediction for  $Y_{21}$  by conditioning on the posterior probabilities for  $X_{21}$  found in Equation 3.4. We get, for example, the posterior probability

$$\Pr(Y_{21} = 4) = \operatorname{E}[\Pr(Y_{21} = 4) | X_{21}]$$
  
=  $\Pr(X_{21} = 1) \operatorname{E}(Q_{14}) + \Pr(X_{21} = 2) \operatorname{E}(Q_{24}) + \Pr(X_{21} = 3) \operatorname{E}(Q_{34})$   
=  $\frac{4}{9} \cdot \frac{1}{13} + \frac{4}{9} \cdot \frac{2}{12} + \frac{1}{9} \cdot \frac{3}{11}$   
= 0.1385651.

For illustration, we try out a second, more structured prior:  $Y_i \sim$ Poisson $(\lambda X_i)$ , where we use a Gamma(2, 2) prior for  $\lambda$ . Such a prior might be chosen if there is a reason to believe that each  $Y_i$  is Poisson distributed with some underlying parameter that might be proportional to the value  $X_i$  of the "hidden chain". When  $X_i$  and  $\lambda$  are known all  $Y_i$  are independent and we can update our knowledge about  $\lambda$  stepwise, using information from one observed  $Y_i$  at the time. Specifically, if  $\lambda \sim \text{Gamma}(\alpha, \beta)$  and  $y \sim \text{Poisson}(x\lambda)$ , then

$$\pi(\lambda \mid y) \quad \propto_{\lambda} \quad \pi(y \mid \lambda)\pi(\lambda)$$
$$\propto_{\lambda} \quad e^{-\lambda x} \frac{(\lambda x)^{y}}{y!} \lambda^{\alpha-1} \exp(-\lambda\beta)$$
$$\propto_{\lambda} \quad \lambda^{\alpha+y-1} \exp(-\lambda(\beta+x))$$

so  $\lambda \mid y \sim \text{Gamma}(\alpha + y, \beta + x)$ . Thus, to find the posterior for  $\lambda$  given all the data, we take the first parameter and add the sum of all the  $Y_i$  observed, obtaining 2 + 40 = 42. To the second parameter we add 1 times the number of observations with  $X_i = 1, 2$  times the number of observations with  $X_i = 2$ , and 3 times the number of observations with  $X_i = 3$ , obtaining  $2+1\cdot8+2\cdot7+3\cdot6=42$ . Thus we have the posterior  $\lambda \mid \text{data} \sim \text{Gamma}(42, 42)$ .

If a state  $X_i$  for i > 20 has value 2, the probabilities for  $Y_i$  will be given by the Poisson distribution with parameter  $2\lambda$ , where  $\lambda \sim \text{Gamma}(42, 42)$ . Writing  $\lambda' = 2\lambda$ , it is fairly direct to show that  $\lambda' \sim \text{Gamma}(42, 42/2) = \text{Gamma}(42, 21)$ . From the general conjugacy theory of Chapter 1 we have that if  $y \sim \text{Poisson}(\lambda)$  and  $\lambda \sim \text{Gamma}(\alpha, \beta)$ , then  $y \sim \text{Negative-Binomial}(\alpha, \beta/(\beta+1))$ . Thus, in our case  $Y_i \mid \text{data} \sim \text{Negative-Binomial}(41, 21/(21+1))$ , and we can for example use the R command dnbinom(0:4, 42, 21/22) to compute some probabilities for  $Y_i$ , resulting in

### [1] 0.1417287 0.2705730 0.2644236 0.1762824 0.0901444

The probability for  $Y_{21} = 4$  using the distribution for  $X_{21}$  can be found with

sum(dnbinom(4, 42, 42/(1:3)/(42/(1:3) + 1))\*c(4, 4, 1)/9)

producing 0.06530846.

## 3.5 Exercises

- 1. Assume an experiment can have one of three outcomes; let us name the outcomes 1, 2, and 3. Assume the probabilities for these outcomes are  $p_1, p_2$ , and  $p_3$ , respectively, but that these probabilities are unknown. Assume 13 independent experiments are performed, of which 3 have outcome 1, 9 have outcome 2, and 1 has outcome 3.
  - a) Using a Dirichlet( $\alpha$ ) prior for  $p = (p_1, p_2, p_3)$ , where  $\alpha = (\alpha_1, \alpha_2, \alpha_3) = (1, 1, 1)$ , find the posterior for p.
  - b) Find the expected posterior value for p. (See the Appendix on probability distributions).
  - c) Still using the same prior, compute the probability that, among the next 4 experiments, there will be 1 with outcome 1, 2 with outcome 2, and 1 with outcome 3.
- 2.
- a) Convince yourself that a Dirichlet distribution with k = 2 is the same as a Beta distribution, just using different notation.
- b) Define a function on the set of non-negative vectors  $\theta = (\theta_1, \dots, \theta_k)$ with  $\sum_{i=1}^k \theta_i = K$  by

$$\pi(\theta \mid \alpha, K) = \frac{1}{K^{\alpha_1 + \dots + \alpha_k - 1}} \cdot \frac{\Gamma(\alpha_1 + \alpha_2 + \dots + \alpha_k)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\dots\Gamma(\alpha_k)} \theta_1^{\alpha_1 - 1} \theta_2^{\alpha_2 - 1} \dots \theta_k^{\alpha_k - 1}$$

Show that  $\pi$  is a density on this set, i.e., that it integrates to 1. We will use the notation  $\theta \mid \alpha, K \sim \text{Dirichlet}_K(\alpha)$ .

c) Assuming that  $(\theta_1, \ldots, \theta_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k)$ , use a proportionality argument to show that for any *i* with 2 < i < k,

 $\theta_1 \ldots, \theta_{i-1} \mid \theta_i, \ldots, \theta_k \sim \text{Dirichlet}_{1-\theta_i-\cdots-\theta_k}(\alpha_1, \ldots, \alpha_{i-1}).$ 

(Note also that the ordering of the indexes does not matter in our context).

d) For *i* with 1 < i < k, use the identity

$$\pi(\theta_1,\ldots,\theta_i) = \frac{\pi(\theta_1,\ldots,\theta_k)}{\pi(\theta_{i+1},\ldots,\theta_k \mid \theta_1,\ldots,\theta_i)}$$

to compute the marginal density  $\pi(\theta_1, \ldots, \theta_i)$  up to a constant not depending on  $\theta_1, \ldots, \theta_i$ .

- e) Assuming that  $(\theta_1, \ldots, \theta_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k)$ , find the distribution of the random vector  $(\theta_1, \ldots, \theta_i, \theta_{i+1} + \cdots + \theta_k)$ .
- f) Assuming that  $(\theta_1, \ldots, \theta_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k)$ , show that  $\theta_1 + \cdots + \theta_i$  has a Beta distribution, and find the parameters of this distribution.
- 3. Assume a Markov chain with state space containing the numbers 1, 2, 3, 4 has been observed for 26 steps. The values in these steps are 1, 2, 3, 2, 3, 4, 4, 3, 2, 3, 2, 1, 1, 2, 1, 2, 3, 4, 3, 4, 3, 3, 2, 1, 1.
  - a) Write down an estimate for the transition matrix P based only on frequencies of observed transitions.
  - b) Using a prior for the transition matrix consisting of a product of Dirchlet distributions with all pseudo-counts equal to 1, find the expectation of the posterior for the transition matrix given the observed sequence above.
  - c) Given the same prior, compute the posterior distribution for  $P_{11} \mid data$ . (Hint: You may need to look up, or solve Exercise 2 above, for the marginal distribution for the components of a Dirchlet distribution). In particular compute  $P(P_{11} > 0.3 \mid data)$ .
  - d) Assume now that you have prior information that transitions in the Markov chain cannot happen to states whose value differs more than one compared to the current state. Reformulate a new prior for P incorporating this information. Then, recompute the results from questions (b) and (c) above using this new prior.
- 4. We will now compute the explicit value for the predictive distribution in Equation 3.5.
  - a) Assume your data is a specific sequence  $z = (z_1, z_2, ..., z_n)$  of outcomes from n independent trials where each trial can have one of k outcomes, with probabilities of the outcomes given by a probability vector  $p = (p_1, ..., p_k)$ . The probability mass function for z is given by

$$\pi(z \mid n, p) = p_1^{c_1} p_2^{c_2} \cdots p_k^{c_k}$$

where  $c_i$  is the count of values in the sequence z equal to the *i*'th outcome. Go through the discussion in this chapter about the Multinomial Dirichlet conjugacy and show what needs to be changed when the Multinomial density is replaced with the density above. In particular, prove that the Dirichlet family is a conjugate family, and compute the predictive distribution corresponding to Equation 3.1.

#### 3. Some basic inference for Markov chains and HMMs

b) Assume you are in the situation of Section 3.2 where you have observed a sequence  $x_0, \ldots, x_n$  of a Markov chain and you want to predict the posterior probability of observing a specific continuation  $x_{n+1}, \ldots, x_{n+r}$ . You assume each line  $P_i$  of your transition matrix has a Dirichlet distribution with parameter vector  $\alpha_i = (\alpha_{i1}, \ldots, \alpha_{ik})$ . Further, you denote by  $c_{ij}$  and  $d_{ij}$  the counts of observed transitions from state *i* to state *j* in the sequences  $x_0, \ldots, x_n$  and  $x_{n+1}, \ldots, x_{n+r}$ , respectively. You also write, for  $i = 1, \ldots, k, c_i = (c_{i1}, \ldots, c_{ik})$ . Show that the probability  $\pi(x_{n+1}, \ldots, x_{n+r} \mid x_0, \ldots, x_n)$  can be written as a product over the values

$$Q_i = \int \left[\prod_{j=1}^k P_{ij}^{d_{ij}}\right] \text{Dirichlet}(P_i; \alpha_i + c_i) \, dP_i$$

for i = 1, ..., k.

c) Prove that

$$Q_i = \frac{\Gamma(\alpha_{i1} + c_{i1} + d_{i1})}{\Gamma(\alpha_{i1} + c_{i1})} \cdots \frac{\Gamma(\alpha_{ik} + c_{ik} + d_{ik})}{\Gamma(\alpha_{ik} + c_{ik})}$$
$$\cdot \frac{\Gamma(\alpha_{i1} + \dots + \alpha_{ik} + c_{i1} + \dots + c_{ik})}{\Gamma(\alpha_{i1} + \dots + \alpha_{ik} + c_{i1} + \dots + c_{ik} + d_{i1} + \dots + d_{ik})}$$

- d) Consider the data of Exercise 3 above, and the prior used in (3b) and (3c). Given this prior, compute the probability of observing the sequence 1,2,2,3 after the sequence given in 3.
- 5. Assume we have a Hidden Markov Model where the Markov chain has a state space consisting of 1, 2, 3, 4, a transition matrix

$$P = \begin{bmatrix} 0.7 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.7 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.7 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.7 \end{bmatrix}$$

and a distribution for the initial state  $X_0$  given by the probability vector (0.2, 0.4, 0.1, 0.3). Assume also that the possible values for  $Y_i$  are also 1, 2, 3, 4, and that the probability matrix for transitions from  $X_i$  to  $Y_i$  is given by

	0.9	0.06	0.03	0.01
0 -	0.04	0.9	0.04	0.02
Q =	0.02	0.04	0.9	0.04
	0.01	0.03	0.06	0.9

Assume we have observed the sequence 3, 4, 1, 1, 4, 3, 4, 3, 2, 2, 1 for  $Y_0, \ldots, Y_{10}$ .

- a) Implement in R the Forward algorithm for this situation. Store the computed distributions.
- b) Implement in R the Backward algorithm for this situation. Store the computed vectors.

- c) For each i = 0, ..., 10, compute in R the vector representing the marginal posterior  $X_i \mid Y_0, ..., Y_{10}$ , when the  $Y_i$  have the values above.
- d) Implement in R a function generating a sequence  $x_0, \ldots, x_{10}$  representing a sample from  $\pi(X_0, \ldots, X_{10} \mid Y_0, \ldots, Y_{10})$ , when the  $Y_i$  have the values above.
- 6. Assume you have the data:

Х	0	1	2	2	1	2	1	0	0	1	2	1	1	0
Y	0	1	1	2	2	1	1	0	1	2	1	1	1	0

We use a HMM model with the three possible states 0, 1, and 2 for the  $X_n$  and  $Y_n$ . For the transition matrix P use a prior that is a product of Dirichlet distributions with all pseudocounts set to 0.5. For the matrix Q defined by

$$Q_{ij} = \Pr\left(Y_n = j \mid X_n = i\right)$$

use a prior that is a product of Dirichlet distributions with all pseudocounts set to 2.

- a) Find the posterior expectation of P given the data.
- b) Find the posterior expectation of Q given the data.
- c) Designating the observations above as  $X_0, \ldots, X_{13}$  and  $Y_0, \ldots, Y_{13}$ , compute the posterior probability distribution for  $X_{14}$ .
- d) Compute the posterior probability distribution for  $Y_{14}$ .

## CHAPTER 4

## Some basic inference for Branching processes

Assume you want to use a Branching process, as defined in Dobrow, as a model in some applied setting. The parameter defining such a Branching process is the vector a of probabilities in the offspring distribution of having  $0, 1, 2, \ldots$ , offspring. Various types of data that one might learn about this parameter are conceivable. Below, we will simply assume that the data consists of nindependent observations of the offspring distribution, i.e., counts  $y_1, y_2, \ldots, y_n$ of the number of actual offspring in n different cases.

Within the Bayesian paradigm, we start inference with defining a probabilistic model for the data. One possibility is to view a as the model parameter. The probability of observed data  $y_1, y_2, \ldots, y_n$  given this parameter is then simply  $\prod_{i=1}^{n} a_{y_i}$ . However, defining a prior on the infinite-dimensional space of possible vectors a is not trivial, and we will not pursue this possibility further here.

Instead, we will look at situations where the context indicates that that the offspring distribution has a certain parametric form. We will consider two different such forms below, and a third in the Exercise. In either case, we find a posterior for the model parameter, and briefly show how to use this posterior in predictions, for example of computation of coming generation sizes or the probability of extinction.

## 4.1 Using a Binomial likelihood

Assume from the context it is natural to assume that the number of offspring is between 0 and N, distributed according to a Binomial distribution with some parameter p. This means that the likelihood for the data is

$$\pi(y_1, \dots, y_n \mid p) = \prod_{i=1}^n \text{Binomial}(y_i; N, p)$$

and to complete the model, we need a prior for  $p \in [0, 1]$ . A possibility is to use a Beta $(\alpha, \beta)$  prior for p, as we did in Chapter 1 of this Compendium. Using the theory of that chapter, we find that the posterior for p after the observations  $y_1, \ldots, y_n$  becomes Beta  $(\alpha + S, \beta + nN - S)$  where  $S = \sum_{i=1}^n y_i$ . This distribution can then be used in predictions for the further growth of the branching process. Instead of pursuing further the use of a Beta prior for p, let us consider how to compute in the more general case when the prior for p is not conjugate. We may then use numerical calculations.

## Example with a general prior

Assume the prior is any non-negative function f(p) on [0, 1] that integrates to 1. We then get for the posterior

$$\pi(p \mid \text{data}) \propto_{p} \pi(\text{data} \mid p)\pi(p) \propto_{p} \left(\prod_{i=1}^{n} \text{Binomial}(y_{i}; N, p)\right) \text{Beta}(p; 1, 1)f(p)$$
$$\propto_{p} \text{Beta}(p; 1 + S, 1 + nN - S)f(p).$$
(4.1)

In the computations above, we used the fact that Beta (p; 1, 1) = 1 for  $0 \le p \le 1$ and that the conjugate computations in this section give us

$$\left(\prod_{i=1}^{n} \operatorname{Binomial}(y_i; N, p)\right) \operatorname{Beta}(p; 1, 1) \propto_p \operatorname{Beta}(p; 1 + S, 1 + nN - S).$$

We can use Equation 4.1 in numerical computations. As a concrete example, suppose that we use a prior  $\pi(p) = f(p)$  with

$$f(p) = \begin{cases} 100(p-0.1) & 0.1 \le p \le 0.2\\ 100(0.3-p) & 0.2 \le p \le 0.3\\ 0 & \text{otherwise} \end{cases}$$
(4.2)

This is a "triangle prior", and it indicates that we believe the most likely value for p is 0.2. Further, p could be as small as 0.1 or as large as 0.3, but it could not be outside this interval. Note how f is defined so that it integrates to 1 over the interval [0, 1].

Let us further assume there can be a maximum of N = 6 offspring and that the number of offspring is distributed according to a Binomial(6, p) distribution. Finally, assume the number of offspring in 342 observed cases are

Number of offspring	0	1	2	3	4	5	6
Number of cases	117	138	58	25	3	0	1

so that the total number of offspring is S = 347. Equation 4.2 shows that the posterior for p is proportional to  $\text{Beta}(p; 1 + 347, 1 + 342 \cdot 6 - 347)f(p) = \text{Beta}(p; 348, 1706)f(p)$ .

What is the probability that the branching process is supercritical? As the offspring process is Binomial its expectation is 6p. We would like to compute the posterior probability that 6p > 1, i.e., that p > 1/6. This can be done for example with

```
prior <- function(p) {
    if (p<0.1) 0 else if (p<0.2) 100*(p-0.1) else
        if (p<0.3) 100*(0.3-p) else 0
}
g <- function(p) dbeta(p, 348, 1706)*Vectorize(prior)(p)
integrate(g, 1/6, 1)$value/integrate(g, 0, 1)$value</pre>
```

producing 0.671135. Note the use of the R function Vectorize: It changes the function **prior** to a similar function that can take *vectorized* input, in other words, instead of just outputing the prior density for a single input, the changed function will also output a vector of prior densities given an input vector of values. This vectorization is necessary as R's integrate function expects as input functions that are vectorized.

Notice also how we compute the result as a quotient of integrals. This is because Equation 4.1 only gives the posterior density up to a proportionality constant. Thus we must compute

$$\int_{1/6}^{1} \pi(p \mid \text{data}) \, dp = \frac{\int_{1/6}^{1} \text{Beta}(p; 348, 1706) f(p) \, dp}{\int_{0}^{1} \text{Beta}(p; 348, 1706) f(p) \, dp}.$$

## 4.2 Using a Multinomial likelihood

We now look at the situation where we assume there is a fixed maximum number offspring N, but where we otherwise don't make assumptions about the offspring distribution. Let p be a probability vector listing the probabilities for the offspring outcomes  $0, 1, \ldots, N$ . Then the likelihood for data  $y_1, y_2, \ldots, y_n$ becomes

$$\pi(\text{data} \mid p) = \text{Multinomial}(c; p)$$

where  $c = (c_0, c_1, \ldots, c_N)$  lists the counts of each type of outcome in the data, in other words,

$$c_j = \sum_{i=1}^n I(y_i = j)$$

for  $j = 0 \dots, N$ , where I is the indicator function. To complete the model, we need a prior for p.

Based on our previous experience with the Multinomial Dirichlet conjugacy, it is natural to try a Dirichlet prior. So assume we use the prior

$$p \sim \text{Dirichlet}(\alpha)$$

where  $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_N)$  is a vector of pseudocounts. The results from Section 3.1 gives the posterior

$$p \mid \text{data} \sim \text{Dirichlet}(\alpha + c).$$

We may then use this posterior for predictions about the future growth or extinction of the branching process.

As a concrete example, let us reconsider the data from the previous section, but now assuming that there is a maximum of 6 offspring, and that the prior on the probability vector p for the outcomes  $0, 1, \ldots, 6$  is uniform. The uniform distribution on the set of probability vectors for 7 outcomes is Dirichlet(1, 1, 1, 1, 1, 1, 1): This may be seen by looking up the Dirichlet density function. Using this prior, i.e., setting  $\alpha = (1, 1, 1, 1, 1, 1, 1)$  and using the data yields the posterior

$$\pi(p \mid \text{data}) = \text{Dirchlet}(p; (118, 139, 59, 26, 4, 1, 2)).$$

So, for example, we get the posterior expectation

$$\mathbf{E}(p \mid \text{data}) = \left(\frac{118}{349}, \frac{139}{349}, \frac{59}{349}, \frac{26}{349}, \frac{4}{349}, \frac{1}{349}, \frac{2}{349}\right).$$

From this we may compute the posterior expectation of the offspring random variable X as

$$\mu = \sum_{i=0}^{6} i \Pr(X = i \mid \text{data}) = \sum_{i=0}^{6} i \operatorname{E}(p_i \mid \text{data})$$
  
=  $0 \cdot \frac{118}{349} + 1 \cdot \frac{139}{349} + 2 \cdot \frac{59}{349} + 3 \cdot \frac{26}{349} + 4 \cdot \frac{4}{349} + 5 \cdot \frac{1}{349} + 6 \cdot \frac{2}{349}$   
=  $1.054441$ 

With our current assumptions, what is the posterior probability that the branching process is supercritical? We have seen that the posterior expectation of the offspring distribution is above 1, but there is some uncertainty in this value. One way to answer the question is by simulation: We simulate posterior values for p, compute for each such p whether the offspring distribution is supercritical, and take the average. In R we may write

#### mean(rdirichlet(1000000, c(118,139,59,26,4,1,2))%\*%(0:6)>1)

which produces 0.835042. The function rdirichlet simulates from the Dirichlet distribution and can be found in the R package LearnBayes.

Note that under the assumptions used in these calculations, the probability that the process is supercritical is higher than the probability we found in the previous section, where the assumptions were different.

## 4.3 Exercise

- 1. In this exercise we use the same data as in the examples above. However, we now assume that the offspring distribution is Geometric(p), with a uniform prior  $\pi(p) \sim \text{Uniform}(0, 1)$ .
  - a) If  $y \sim \text{Geometric}(p)$  find a family of distributions for p that is a conjugate family. Do this by guessing and trying out if your guess is correct. (If this fails you may also look up a conjugate family for the Geometric distribution).
  - b) Compute the posterior distribution for the parameter p of the offspring distribution when the prior above is used.
  - c) Compute the probability that the process is supercritical under the assumptions if this exercise.

## CHAPTER 5

## Markov chain Monte Carlo (MCMC)

As described in Section 1.4, the overall idea of Bayesian inference is to describe a stochastid model with variables representing both y, the observed data, and  $y_{new}$ , whatever you would like to predict, and then use the resulting conditional distribution of  $y_{new} \mid y$  for predictions. Usually, one describes the model in terms of a vector of parameters  $\theta$ , describing models for likelihoods  $y \mid \theta$  and  $y_{new} \mid \theta$  and a prior distribution for  $\theta$  in such a way that y and  $y_{new}$  are conditionally independent given  $\theta$ , i.e.,

$$\pi(y_{new} \mid \theta, y) = \pi(y_{new} \mid \theta).$$

We can then write

$$\pi(y_{new} \mid y) = \int_{\theta} \pi(y_{new} \mid \theta) \pi(\theta \mid y) \, d\theta \tag{5.1}$$

and Bayesian inference resolves into three steps: First describing the stochastic model (i.e., the likelihoods and the prior), second deriving from this the posterior  $\pi(\theta \mid y)$ , and finally using the posterior to to make predictions according to Equation 5.1.

Computing the integral in Equation 5.1 can be a big challenge in many types of models. In Chapter 1, we saw how it can be done numerically when the number of dimensions of  $\theta$  is very low (in practice 1-3). We have also seen a number of cases where the integral can be computed analytically, i.e., where we can use conjugacy. However, in most realistic models, none of these options are available, and one needs to turn to approximate numerical approaches.

Assume we can generate a sample  $\theta_1, \theta_2, \ldots, \theta_m$  from the posterior distribution  $\pi(\theta \mid y)$ . Then we can approximate the integral above with the average of the numbers  $\pi(y_{new} \mid \theta_i)$ . More precisely, the Strong Law of Large Numbers gives that, with probability 1,

$$\int_{\theta} \pi(y_{new} \mid \theta) \pi(\theta \mid y) \, d\theta = \mathcal{E}_{\theta \mid y} \left[ \pi(y_{new} \mid \theta) \right] = \lim_{m \to \infty} \frac{1}{m} \sum_{i=1}^{m} \pi(y_{new} \mid \theta_i).$$
(5.2)

Importantly for us, Equation 5.2 holds even in the case where  $\theta_1, \theta_2, \ldots, \theta_i, \ldots$  is not a random sample but instead the values in a Markov chain with limiting distribution  $\pi(\theta \mid y)$ . Equation 5.2 is then called the Strong Law of Large Numbers for Markov chains. The second and third steps of Bayesian inference

can now be done as follows: Generate a sequence  $\theta_1, \ldots, \theta_m$  from a Markov chain with the posterior  $\pi(\theta \mid y)$  as limiting distribution, and approximate

$$\pi(y_{new} \mid y) \approx \frac{1}{m} \sum_{i=1}^{m} \pi(y_{new} \mid \theta_i).$$
(5.3)

This technique for Bayesian inference is called Markov chain Monte Carlo, or MCMC. To use it, you need to

- 1. Define and simulate from a Markov chain with your posterior  $\pi(\theta \mid y)$  as a limiting distribution.
- 2. Do this in a way so that you know, or at least can argue, that the approximation in Equation 5.3 is sufficiently accurate.

It turns out that the first point above is often surprisingly easy. We will briefly describe the Metropolis Hastings algorithm in the next section. However, the second point is surprisingly difficult. We will discuss it in the remaining parts of this chapter.

## 5.1 Bayesian inference with the Metropolis Hastings algorithm

Assume given a likelihood function  $\pi(y \mid \theta)$  and a prior  $\pi(\theta)$ . Define a *proposal* function which, for every  $\theta$  in the set  $\Omega$  of possible parameters describes a probability density  $q(\theta^* \mid \theta)$  on  $\theta^* \in \Omega$ . Assume you have an algorithm to simulate  $\theta^*$  from this proposal distribution. Then the Metropolis Hastings algorithm for Bayesian inference is:

- 1. Simulate  $\theta_0$  from some distribution on  $\Omega$ .
- 2. For i = 1, ..., m:
  - a) Generate  $\theta^*$  from  $q(\theta^* \mid \theta_{i-1})$ .
  - b) Generate  $U \sim \text{Uniform}(0, 1)$ .
  - c) If

$$U < a_{\theta,\theta^*} = \frac{\pi(y \mid \theta^*)\pi(\theta^*)q(\theta_{i-1} \mid \theta^*)}{\pi(y \mid \theta_{i-1})\pi(\theta_{i-1})q(\theta^* \mid \theta_{i-1})}$$

set  $\theta_i = \theta^*$ , otherwise set  $\theta_i = \theta_{i-1}$ .

This will generate a Markov chain  $\theta_0, \theta_1, \ldots$  In Dobrow it is proven that as long as this Markov chain is ergodic, it will have limiting distribution  $\pi(\theta \mid y)$ . Note how

$$\pi(y \mid \theta^*) = \frac{\pi(y \mid \theta^*)\pi(\theta^*)}{\pi(y)} \quad \text{and} \quad \pi(y \mid \theta_{i-1}) = \frac{\pi(y \mid \theta_{i-1})\pi(\theta_{i-1})}{\pi(y)}$$

imply

$$\frac{\pi(\theta^* \mid y)}{\pi(\theta_{i-1} \mid y)} = \frac{\pi(y \mid \theta^*)\pi(\theta^*)}{\pi(y \mid \theta_{i-1})\pi(\theta_{i-1})}.$$

This is why we can use the quotient on the right instead of the quotient on the left when computing  $a_{\theta,\theta^*}$ .

44

The data y will often consist of a random sample  $y_1, y_2, \ldots, y_n$  of observations that are independent given  $\theta$ . Then  $\pi(y \mid \theta) = \prod_{j=1}^n \pi(y_i \mid \theta)$  and numerically this number may come extremely close to zero, in particular when n is large, so that the computer may round it down to zero. Writing

$$a_{\theta,\theta^{*}} = \frac{\pi(y \mid \theta^{*})\pi(\theta^{*})q(\theta_{i-1} \mid \theta^{*})}{\pi(y \mid \theta_{i-1})\pi(\theta_{i-1})q(\theta^{*} \mid \theta_{i-1})} = \frac{\prod_{j=1}^{n}\pi(y_{j} \mid \theta^{*})\pi(\theta^{*})q(\theta_{i-1} \mid \theta^{*})}{\prod_{j=1}^{n}\pi(y_{j} \mid \theta_{i-1})\pi(\theta_{i-1})q(\theta^{*} \mid \theta_{i-1})} = \exp\left(\sum_{j=1}^{n}\left[\log \pi(y_{j} \mid \theta^{*}) - \log \pi(y_{j} \mid \theta_{i-1})\right] + \log \pi(\theta^{*}) - \log \pi(\theta_{i-1}) + \log q(\theta_{i-1} \mid \theta^{*}) - \log q(\theta^{*} \mid \theta_{i-1})\right)\right)$$
(5.4)

and computing  $a_{\theta,\theta^*}$  according to the last expression is a way to avoid such underflow problems on the computer.

## 5.2 Using a random walk proposal

A popular choice for the proposal density  $q(\theta^* \mid \theta_{i-1})$  is to write

$$\theta^* = \theta_{i-1} + \epsilon$$

where  $\epsilon$  is some random variable with symmetric density, i.e.,  $\pi(\epsilon) = \pi(-\epsilon)$ . With such a definition we get that  $q(\theta^* \mid \theta_{i-1}) = q(\theta_{i-1} \mid \theta^*)$  so Equation 5.4 simplifies to

$$a_{\theta,\theta^*} = \exp\left(\sum_{j=1}^n \left[\log \pi(y_j \mid \theta^*) - \log \pi(y_j \mid \theta_{i-1})\right] + \log \pi(\theta^*) - \log \pi(\theta_{i-1})\right).$$

## Example

We go back to the toy example used in Sections 1.5 and 1.6: Assume

$$p \sim \text{Beta}(2.3, 4.1)$$
  
 $y \mid p \sim \text{Binomial}(17, p)$   
 $y_{new} \mid p \sim \text{Binomial}(3, p).$ 

Assume we would like to compute the probability  $\pi(y_{new} = 1 \mid y = 4)$ . In the previous sections, we found, using either analytic computation, discretization, or numerical integration, the value 0.4033. Let us now use random walk Metropolis Hastings with  $\epsilon \sim \text{Normal}(0, 0.3^2)$ . To make sure the proposed value is still between 0 and 1, we subtract from it the largest integer less than it: Note that the proposal will still be symmetric, so that we can use the simplified computation of  $a_{\theta,\theta^*}$  above. We can now use the R code

```
N <- 1000
posterior <- rep(runif(1), N)
for (i in 2:N) {</pre>
```

```
Markov chain Monte Carlo (MCMC)
```

```
proposal <- posterior[i-1] + rnorm(1, 0, 0.3)
proposal <- proposal - floor(proposal)
loga <- dbinom(4, 17, proposal, log=T) -
    dbinom(4, 17, posterior[i-1], log=T) +
    dbeta(proposal, 2.3, 4.1, log=T) -
    dbeta(posterior[i-1], 2.3, 4.1, log=T)
if (runif(1)<exp(loga))
    posterior[i] <- proposal
else
    posterior[i] <- posterior[i-1]
}
print(mean(dbinom(1, 3, posterior)))</pre>
```

We get for example the result 0.4024178; the result will vary with every run of the code above.

## 5.3 Using independent proposals

Another choice for the proposals may be to simply generate them from some fixed density. Below, we recompute the result from the previous section using the independent proposal function Beta(2, 2). A possible code is then

```
N <- 1000
posterior <- rep(runif(1), N)</pre>
for (i in 2:N) {
  proposal <- rbeta(1, 2, 2)</pre>
  loga <- dbinom(4, 17, proposal, log=T) -</pre>
    dbinom(4, 17, posterior[i-1], log=T) +
    dbeta(proposal, 2.3, 4.1, log=T) -
    dbeta(posterior[i-1], 2.3, 4.1, log=T) +
    dbeta(posterior[i-1], 2, 2, log=T) -
    dbeta(proposal, 2, 2, log=T)
  if (runif(1)<exp(loga))</pre>
    posterior[i] <- proposal</pre>
  else
    posterior[i] <- posterior[i-1]</pre>
}
print(mean(dbinom(1, 3, posterior)))
```

A result is 0.4043354; the result will vary with every run of the code.

## 5.4 Gibbs sampling

Inference for the toy examples above could more easily be done for example with the methods of Sections 1.5 and 1.6. However when  $\theta$  is a vector, say  $\theta = (\theta_1, \theta_2, \ldots, \theta_k)$ , MCMC can be programmed in the same way, and the strengths of the algorithm is then much more apparent compared to the alternatives mentioned.

When  $\theta$  is a vector there is in addition to the choices above another popular choice for proposals which leads to the algorithm called Gibbs sampling. For each

dimension *i* in the vector, a proposal is made as follows: All  $\theta_j$  for  $j \neq i$  are kept fixed, while a new value for  $\theta_i$  is simulated from the conditional distribution of  $\theta_i$ given the remaining variables. It is fairly easy to show that using such a proposal in the Metropolis Hastings algorithm results in an acceptance probability of 1. If one somehow cycles through all the dimensions of  $\theta$ , one may obtain an ergodic Markov chain, and thus a valid MCMC algorithm.

## Example

Consider the model

$$y_1, y_2, \dots, y_n \sim \operatorname{Normal}(\mu, \tau^{-1})$$
  
 $\mu \sim \operatorname{Normal}(\mu_0, \tau_0^{-1})$   
 $\tau \sim \operatorname{Gamma}(\alpha, \beta).$ 

In other words, we have data  $y_1, \ldots, y_n$  that is normally distributed given parameters  $\theta = (\mu, \tau)$ , and there is a prior given by  $\pi(\theta) = \pi(\mu, \tau) =$ Normal $(\mu; \mu_0, \tau_0^{-1})$  Gamma $(\tau; \alpha, \beta)$ . We assume that  $\mu_0, \tau_0, \alpha, \beta$  are fixed and known numbers. If we plan a new observation  $y_{\text{new}} \sim \text{Normal}(\mu, \tau^{-1})$  what is the probability that, say,  $y_{\text{new}} < y_1$ ?

We can solve this using Gibbs sampling. For the two components of the parameter,  $\mu$  and  $\tau$ , we need to find their conditional distributions given that the data and the other parameter are fixed. Let us start with  $\tau$ . According to Appendix C we have a Normal Gamma conjugacy: If  $x \sim \text{Normal}(\mu, \theta^{-1})$  and  $\theta \sim \text{Gamma}(\alpha, \beta)$ , then  $\theta \mid x \sim \text{Gamma}(\alpha + 1/2, \beta + (x - \mu)^2/2)$ . Translating the names of the variables and applying this rule *n* times, we get that

$$\tau \mid y_1, \dots, y_n, \mu \sim \text{Gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2}\sum_{i=1}^n (y_i - \mu)^2\right).$$

For  $\mu$ , we can use the Normal-Normal conjugacy from Appendix C. Applying that rule *n* times we now get

$$\mu \mid y_1, \dots, y_n, \tau \sim \operatorname{Normal}\left(\frac{\tau_0 \mu_0 + \tau \sum_{i=1}^n y_i}{\tau_0 + n\tau}, \frac{1}{\tau_0 + n\tau}\right).$$

R code for Gibbs sampling for such an example is

## 5.5 Advantages and disadvantages with MCMC for Bayesian inference

MCMC is a very flexible technique with a huge range of applications. A big advantage is that, for many models, it is quite simple and fast to program an implementation of Metropolis Hastings that works with reasonable accuracy, at least if the Markov chain is run long enough. A big disadvantage is that some experience and skill, and some understanding of the posterior you are trying to simulate from, may be needed in order to select a proposal function that leads to reasonably fast and accurate results. An even bigger disadvantage may be that, except for some limited cases, there are no general mathematical proofs about the accuracy of results.

Although the Metropolis Hastings algorithm is very flexible, and has several important special cases such as Gibbs sampling, it is not the last word in how one can generate approximate samples from posteriors. Finding improved algorithms for such situations is an active research field.

Many software packages exist which tries to make it easier to run MCMC even for those without expert knowledge of the relevant algorithms. An important package is Stan (https://mc-stan.org).

Finally, it should be mentioned that MCMC algorithms tend to be less effective when the dimension of  $\theta$  increases above a few thousand. In modern machine learning applications, where parameter spaces can have millions of dimensions, other approaches, such as Variational Bayes, seem to be more effective.

## Appendices

## APPENDIX A

## Some solutions to some Exercises

## A.1 Exercises from Chapter 1

1. We want to compute the conditional probability that a person votes for B given that he or she lives in a rental flat. This is the quotient of the probability that the person lives in a rental flat and votes for B, divided by the probability that he or she lives in a rental flat:

$$\frac{0.03}{0.11 + 0.03 + 0.08 + 0.01} = 0.1304 = 13\%.$$

2. Let A denote that the person is affected, and P that the test is positive. Using Bayes formula, we get

$$\pi(A \mid P) = \frac{\pi(P \mid A)\pi(A)}{\pi(P \mid A)\pi(A) + \pi(P \mid A^c)\pi(A^c)} = \frac{0.95 \cdot 0.007}{0.95 \cdot 0.007 + 0.05 \cdot 0.993}$$
  
= 0.118 \approx 12%.

4	٦		
	,		
-		í.	
×			٠

a) A direct computation gives

$$\pi(\theta \mid \text{data}) \quad \propto_{\theta} \quad \pi(\text{data} \mid \theta)\pi(\theta) \propto_{\theta} \text{Binomial} (9; 12, \theta)$$
$$\propto_{\theta} \quad \theta^{9}(1-\theta)^{12-9} \propto_{\theta} \text{Beta} (\theta; 10, 4) \,.$$

A more direct argument uses that the uniform distribution on [0, 1] is identical to the Beta (1, 1) distribution; with this the formulas of the lecture notes can be used to derive the posterior Beta (1 + 9, 1 + 3)directly.

- b) The posterior becomes Beta(10 + 11, 4 + 19 11) = Beta(21, 12).
- c) The probability for success is the expectation of the Beta (21, 12) density, which is  $\frac{21}{21+12} = \frac{21}{33}$ .

4.

a)

$$\begin{pmatrix} y_H + y_T \\ y_H \end{pmatrix} (0.5 \cdot 0.3^{y_H} 0.7^{y_T} + 0.5 \cdot 0.7^{y_H} 0.3^{y_T})$$

b)  

$$\begin{pmatrix} y_H + y_T \\ y_H \end{pmatrix} \left( 0.5 \cdot 0.3^{y_H + 1} 0.7^{y_T} + 0.5 \cdot 0.7^{y_H + 1} 0.3^{y_T} \right)$$

c) We get, as before

$$\frac{0.3^{y_H+1}0.7^{y_T}+0.7^{y_H+1}0.3^{y_T}}{0.3^{y_H}0.7^{y_T}+0.7^{y_H}0.3^{y_T}}$$

5. For example

```
> fn <- function(alpha) {(pbeta(0.4, alpha, alpha)-0.05)^2}
> optimize(fn, c(1, 1000))
$'minimum'
[1] 33.38651
```

6. Assume  $p \sim \text{Beta}(\alpha_0, \beta_0)$ . Then

$$\begin{aligned} \pi(p \mid y) & \propto_p & \pi(y \mid p) \pi(p) \\ & \propto_p & (1-p)^y p^{\alpha} \cdot p^{\alpha_0 - 1} (1-p)^{\beta_0 - 1} \\ & \propto_p & p^{\alpha_0 + \alpha - 1} (1-p)^{\beta_0 + y - 1} \end{aligned}$$

 $\mathbf{SO}$ 

$$p \mid y \sim \text{Beta}(\alpha_0 + \alpha, \beta_0 + y)$$

and we have proved conjugacy.

7. Assume that  $\tau \sim \text{Gamma}(\alpha, \beta)$ . Then

$$\pi(\tau \mid y) \propto_{\tau} \pi(y \mid \tau)\pi(\tau)$$
  
$$\propto_{\tau} \tau^{1/2} \exp\left(-\frac{\tau}{2}(y-\mu)^{2}\right)\tau^{\alpha-1}\exp(-\beta\tau)$$
  
$$\propto_{\tau} \tau^{\alpha+1/2-1}\exp(-(\beta+\frac{1}{2}(y-\mu)^{2})\tau)$$

 $\mathbf{SO}$ 

$$\tau \mid y \sim \text{Gamma}\left(\alpha + \frac{1}{2}, \beta + \frac{1}{2}(y-\mu)^2\right)$$

and we have proved conjugacy.

8. Writing

$$k \mid p \sim \operatorname{Binomial}(n, p)$$
  
 $p \sim \operatorname{Beta}(\alpha, \beta)$ 

we get

$$\pi(k) = \frac{\pi(k \mid p)\pi(p)}{\pi(p \mid k)}$$

$$= \frac{\text{Binomial}(k; n, p) \cdot \text{Beta}(p; \alpha, \beta)}{\text{Beta}(\alpha + k, \beta + n - k)}$$

$$= \frac{\binom{n}{k}p^{k}(1-p)^{n-k}\frac{1}{B(\alpha,\beta)}p^{\alpha-1}(1-p)^{\beta-1}}{\frac{1}{B(\alpha+k,\beta+n-k)}p^{\alpha+k-1}(1-p)^{\beta+n-k-1}}$$

52

$$= \frac{B(\alpha+k,\beta+n-k)}{B(\alpha,\beta)} \binom{n}{k}$$

This is a probability mass function on the set of integers  $\{0, \ldots, n\}$  for all real numbers  $\alpha > 0$ ,  $\beta > 0$ , and integers n > 0. In fact, it is called the Beta-Binomial distribution.

9.

a) Writing (for i = 1, ..., 6)  $c_i$  for the count of throws with outcome i among the throws  $k_1, ..., k_n$ , we get

$$\begin{aligned} \pi(\theta \mid k_1, \dots, k_n) \\ \propto_{\theta} & \pi(k_1, \dots, k_n \mid \theta) \pi(\theta) \\ \propto_{\theta} & \left(\frac{1}{3}\theta_1\right)^{c_1} \left(\frac{1}{3}\theta_2\right)^{c_2} \left(\frac{1}{3}\theta_3\right)^{c_3} \left(\frac{1}{3}(1-\theta_3)\right)^{c_4} \left(\frac{1}{3}(1-\theta_2)\right)^{c_5} \left(\frac{1}{3}(1-\theta_1)\right)^{c_6} \\ & \theta_1^{20-1}(1-\theta_1)^{20-1}\theta_2^{20-1}(1-\theta_2)^{20-1}\theta_3^{20-1}(1-\theta_3)^{20-1} \\ \propto_{\theta} & \theta_1^{20+c_1-1}(1-\theta_1)^{20+c_6-1}\theta_2^{20+c_2-1}(1-\theta_2)^{20+c_5-1}\theta_3^{20+c_3-1}(1-\theta_3)^{20+c_4-1}. \end{aligned}$$

Thus the posterior is

$$Beta(\theta_1; 20+c_1, 20+c_6) \cdot Beta(\theta_2; 20+c_2, 20+c_5) \cdot Beta(\theta_3; 20+c_3, 20+c_4)$$

b) We have the counts

$$c_{1} = 1 \\ c_{2} = 3 \\ c_{3} = 2 \\ c_{4} = 2 \\ c_{5} = 0 \\ c_{6} = 3$$

It is possible to use the predictive distribution found in Exercise 8 to answer the question, but the easiest approach may be to use the same thinking as in the start of Chapter 1. Let  $k_{n+1}$  be the outcome of the (n + 1)'st throw. Then

$$\pi(k_{n+1} \mid k_1, \dots, k_n) = \int \pi(k_{n+1} \mid \theta) \pi(\theta \mid k_1, \dots, k_n) \, d\theta$$
$$= \begin{cases} \frac{1}{3} \int \theta_{k_{n+1}} \pi(\theta \mid k_1, \dots, k_n) \, d\theta & k_{n+1} = 1, 2, 3\\ \frac{1}{3} \left[ 1 - \int \theta_{7-k_{n+1}} \pi(\theta \mid k_1, \dots, k_n) \, d\theta \right] & k_{n+1} = 4, 5, 6 \end{cases}$$

As we have the Expectations, for i = 1, 2, 3,

$$\int \theta_i \pi(\theta \mid k_1, \dots, k_n) \, d\theta = \mathbf{E} \left[ \text{Beta}(20 + c_i, 20 + c_{7-i}) \right]$$
$$= \frac{20 + c_i}{20 + c_i + 20 + c_{7-i}},$$

we get

$$\begin{aligned} \pi(k_{n+1} = 1 \mid k_1, \dots, k_n) &= \frac{1}{3} \cdot \frac{20+1}{20+1+20+3} = 0.159 \\ \pi(k_{n+1} = 2 \mid k_1, \dots, k_n) &= \frac{1}{3} \cdot \frac{20+3}{20+3+20+0} = 0.178 \\ \pi(k_{n+1} = 3 \mid k_1, \dots, k_n) &= \frac{1}{3} \cdot \frac{20+2}{20+2+20+2} = 0.167 \\ \pi(k_{n+1} = 4 \mid k_1, \dots, k_n) &= \frac{1}{3} \left(1 - \frac{20+2}{20+2+20+2}\right) = 0.167 \\ \pi(k_{n+1} = 5 \mid k_1, \dots, k_n) &= \frac{1}{3} \left(1 - \frac{20+3}{20+3+20+0}\right) = 0.155 \\ \pi(k_{n+1} = 6 \mid k_1, \dots, k_n) &= \frac{1}{3} \left(1 - \frac{20+1}{20+1+20+3}\right) = 0.174 \end{aligned}$$

10.

a) Assume we use a prior  $\pi(\theta) = r(\theta)$  as defined in the exercise. For each  $i = 1, \ldots, k$ , define the prior predictive

$$r_i(y) = \int_{\theta} \pi(y \mid \theta) q_{\gamma_i}(\theta) \, d\theta = \frac{\pi(y \mid \theta) q_{\gamma_i}(\theta)}{q_{f(\gamma_i)}(\theta)}$$

and we can also write

$$r_i(y)q_{f(\gamma_i)}(\theta) = \pi(y \mid \theta)q_{\gamma_i}(\theta).$$

Then,

$$\pi(\theta \mid y) \propto_{\theta} \pi(y \mid \theta)\pi(\theta)$$

$$= \pi(y \mid \theta) \sum_{i=1}^{k} \lambda_{i} q_{\gamma_{i}}(\theta)$$

$$= \sum_{i=1}^{k} \lambda_{i} [\pi(y \mid \theta)q_{\gamma_{i}}(\theta)]$$

$$= \sum_{i=1}^{k} \lambda_{i} [r_{i}(y)q_{f(\gamma_{i})}(\theta)]$$

To get the actual posterior, we need to normalize this so that it integrates to 1 over  $\theta$ . But the densities  $q_{f(\gamma_i)}(\theta)$  all integrate to 1, so we get

$$\pi(\theta \mid y) = \frac{\sum_{i=1}^{k} \lambda_i r_i(y) q_{f(\gamma_i)}(\theta)}{\sum_{i=1}^{k} \lambda_i r_i(y)}.$$

We see that this is a new density in the family defined in the exercise, with  $\gamma_i'=f(\gamma_i)$  and

$$\lambda_i' = \frac{\lambda_i r_i(y)}{\sum_{j=1}^k \lambda_j r_j(y)}.$$

This proves conjugacy.

b) We have

$$\pi(y) = \int_{\theta} \pi(y \mid \theta) \sum_{i=1}^{k} \lambda_i q_{\gamma_i}(\theta) \, d\theta = \sum_{i=1}^{k} \lambda_i \int_{\theta} r_i(y) q_{f(\gamma_i)}(\theta) \, d\theta = \sum_{i=1}^{k} \lambda_i r_i(y).$$

- c) Conjugate priors are very simple and practical to use, but they are not very flexible. For example, in situations where the normal family is conjugate, it may not always be reasonable to use a normal prior. However, weighted sums of normals is represents a much more flexible class of densities, and thus can often be used when a single normal prior cannot.
- 11. Knowing that the marginal distribution for X is normal, we only have to compute the expection and variance of this random variable to find its distribution. We get

$$\mathbf{E}(X) = \mathbf{E}_{\mu} \left( \mathbf{E}_{X|\mu} \left( X \right) \right) = \mathbf{E}_{\mu}(\mu) = \mu_0$$

and

$$\operatorname{Var}(X) = \operatorname{E}_{\mu} \left( \operatorname{Var}_{X|\mu}(X) \right) + \operatorname{Var}_{\mu} \left( \operatorname{E}_{X|\mu}(X) \right) = \operatorname{E}_{\mu}(\tau^{-1}) + \operatorname{Var}_{\mu}(\mu) = \tau^{-1} + \tau_0^{-1}.$$

Thus we have the prior predictive

$$x \sim \text{Normal}(\mu_0, \tau^{-1} + \tau_0^{-1}).$$

## A.2 Exercises from Chapter 3

1	
T	

a) According to the formulas in this chapter, the poterior is

$$p \mid \text{data} \sim \text{Dirichlet}(1+3, 1+9, 1+1) = \text{Dirichlet}(4, 10, 2)$$

b) We get

$$\mathcal{E}(p \mid \text{data}) = \frac{c(4, 10, 2)}{4 + 10 + 2} = \left(\frac{4}{16}, \frac{10}{16}, \frac{2}{16}\right)$$

c) According to Equation 3.1 we get

$$\pi(x = (1, 2, 1)) = \frac{4!}{1!2!1!} \cdot \frac{\Gamma(4+1)}{\Gamma(4)} \cdot \frac{\Gamma(10+2)}{\Gamma(10)} \cdot \frac{\Gamma(2+1)}{\Gamma(2)} \frac{\Gamma(16)}{\Gamma(16+4)}$$
$$= 12 \cdot 4 \cdot 10 \cdot 11 \cdot 2 \cdot \frac{1}{16 \cdot 17 \cdot 18 \cdot 19} = 0.1135$$

2.

a)

## A. Some solutions to some Exercises

b) To compute the integral below, we use the change of variables  $\theta_i = K u_i$  for  $i = 1, \ldots, k - 1$ . Writing also  $\theta_k = K u_k$ , we get  $\sum_{i=1}^k u_i = 1$ . Note that  $\theta_k$  is completely determined by  $\theta_1, \ldots, \theta_{k-1}$ , so the density is (k-1)-dimensional. We get

$$\int \frac{1}{K^{\alpha_1 + \dots + \alpha_k - 1}} \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_k)} \theta_1^{\alpha_1 - 1} \cdots \theta_k^{\alpha_k - 1} d\theta_1 \dots d\theta_{k-1}$$

$$= \int \frac{1}{K^{\alpha_1 + \dots + \alpha_k - 1}} \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_k)} (Ku_1)^{\alpha_1 - 1} \cdots (Ku_k)^{\alpha_k - 1} K^{k-1} du_1 \dots du_{k-1}$$

$$= \int \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_k)} u_1^{\alpha_1 - 1} \cdots u_k^{\alpha_k - 1} du_1 \dots du_{k-1}$$

$$= 1$$

where in the last step we use that the standard Dirichlet density integrates to 1.

c) We get

$$\pi(\theta_1, \dots, \theta_{i-1} \mid \theta_i, \dots, \theta_k) \propto_{\theta_1, \dots, \theta_{i-1}} \pi(\theta_1, \dots, \theta_k) \propto_{\theta_1, \dots, \theta_{i-1}} \theta_1^{\alpha_1 - 1} \cdots \theta_{i-1}^{\alpha_{i-1} - 1}$$
  
Now,  $\theta_1 + \dots + \theta_{i-1} = 1 - \theta_i - \dots - \theta_k$ . Comparing with the densities defined in (b), we get that

$$\theta_1 \ldots, \theta_{i-1} \mid \theta_i, \ldots, \theta_k \sim \text{Dirichlet}_{1-\theta_i-\cdots-\theta_k}(\alpha_1, \ldots, \alpha_{i-1}).$$

d) We get

$$\pi(\theta_1, \dots, \theta_i) = \frac{\pi(\theta_1, \dots, \theta_k)}{\pi(\theta_{i+1}, \dots, \theta_k \mid \theta_1, \dots, \theta_i)}$$

$$= \frac{\text{Dirichlet}((\theta_1, \dots, \theta_k); (\alpha_1, \dots, \alpha_k))}{\text{Dirichlet}_{1-\theta_1 - \dots - \theta_i}((\theta_{i+1}, \dots, \theta_k); (\alpha_{i+1}, \dots, \alpha_k)))}$$

$$\propto_{\theta_1, \dots, \theta_i} \frac{\theta_1^{\alpha_1 - 1} \cdots \theta_i^{\alpha_i - 1}}{1/(1 - \theta_1 - \dots - \theta_i)^{\alpha_{i+1} + \dots + \alpha_k - 1}}$$

$$= \theta_1^{\alpha_1 - 1} \cdots \theta_i^{\alpha_i - 1}(1 - \theta_1 - \dots - \theta_i)^{\alpha_{i+1} + \dots + \alpha_k - 1}$$

e) When  $(\theta_1, \ldots, \theta_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k)$ , we have that  $\theta_{i+1} + \cdots + \theta_k$  is completely determined by  $(\theta_1, \ldots, \theta_i)$ : It is equal to 1 minus the sum of these numbers. Thus the density for the vector  $(\theta_1, \ldots, \theta_i, \theta_{i+1} + \cdots + \theta_k)$  is equal to the density found in (d). From this, and the fact that  $\theta_1 + \cdots + \theta_i + (\theta_{i+1} + \cdots + \theta_k) = 1$ , we can read off that

$$\theta_1, \ldots, \theta_i, \theta_{i+1} + \cdots + \theta_k \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_i, \alpha_{i+1} + \cdots + \alpha_k)$$

f) Applying the result from (e) twice, we get that

$$(\theta_1 + \dots + \theta_i, \theta_{i+1} + \dots + \theta_k) \sim \text{Dirichlet}(\alpha_1 + \dots + \alpha_i, \alpha_{i+1} + \dots + \alpha_k)$$

Using the result from (a) we see that

$$\theta_1 + \ldots \theta_i \sim \text{Beta}(\alpha_1 + \cdots + \alpha_i, \alpha_{i+1} + \cdots + \alpha_k)$$

- 3.
- a) A table with counts of transitions is

	1	2	3	4
1	2	3	0	0
2	3	1	4	0
3	0	4	1	3
4	0	0	3	1

Thus we get the estimate

$$\hat{P} = \begin{bmatrix} 2/5 & 3/5 & 0 & 0\\ 3/8 & 1/8 & 4/8 & 0\\ 0 & 4/8 & 1/8 & 3/8\\ 0 & 0 & 3/4 & 1/4 \end{bmatrix}.$$

b)

$$\mathbf{E}[P \mid \text{data}] = \begin{bmatrix} 3/9 & 4/9 & 1/9 & 1/9 \\ 4/12 & 2/12 & 5/12 & 1/12 \\ 1/12 & 5/12 & 2/12 & 4/12 \\ 1/8 & 1/8 & 4/8 & 2/8 \end{bmatrix}.$$

c) We have

 $P_1 \mid \text{data} \sim \text{Dirichlet}(3, 4, 1, 1).$ 

According to Exercise 2f, we thus get

 $P_{11} \mid \text{data} \sim \text{Beta}(3, 6).$ 

And pbeta(0.3, 3, 6, lower.tail=FALSE) produces 0.5517738.

d) Choosing to keep pseudocounts equal to 1 for the states that are possible, we get

$$\pi(P) = \text{Dirichlet}(P_1; (1, 1, 0, 0)) \text{Dirichlet}(P_2; (1, 1, 1, 0)) \text{Dirichlet}(P_3; (0, 1, 1, 1)) \text{Dirichlet}(0, 0, 1, 1)$$

With this prior, the expectation of the posterior for P becomes

$$\mathbf{E}[P \mid \text{data}] = \begin{bmatrix} 3/7 & 4/7 & 0 & 0\\ 4/11 & 2/11 & 5/11 & 0\\ 0 & 5/11 & 2/11 & 4/11\\ 0 & 0 & 4/6 & 2/6 \end{bmatrix}$$

and the command pbeta(0.3, 3, 4, lower.tail=FALSE) produces 0.74431.

4.

a) Assume  $p \sim \text{Dirichlet}(\alpha)$  where  $\alpha = (\alpha_1, \ldots, \alpha_k)$ . Then, exactly as for the Multinomial case,

 $p \mid z \sim \text{Dirichlet}(\alpha_1 + c_1, \dots, \alpha_k + c_k) = \text{Dirichlet}(\alpha + c)$ 

where  $c = (c_1, \ldots, c_k)$ , so the Dirichlet family is a conjugate family. Further,

$$\begin{aligned} \pi(z) &= \frac{\pi(z \mid p)\pi(p)}{\pi(p \mid z)} \\ &= \frac{p_1^{c_1} \cdots p_k^{c_k} \operatorname{Dirichlet}(p; \alpha)}{\operatorname{Dirichlet}(p; \alpha + c)} \\ &= \frac{p_1^{c_1} \cdots p_k^{c_k} \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} p_1^{\alpha_1 - 1} \cdots p_k^{\alpha_k - 1}}{\frac{\Gamma(\alpha_1 + \dots + \alpha_k + y_1 + \dots + y_k)}{\Gamma(\alpha_1 + c_1) \dots \Gamma(\alpha_k + c_k)} p_1^{\alpha_1 + c_1 - 1} \dots p_k^{\alpha_k + c_k - 1}} \\ &= \frac{\Gamma(\alpha_1 + c_1)}{\Gamma(\alpha_1)} \cdots \frac{\Gamma(\alpha_k + c_k)}{\Gamma(\alpha_k)} \cdot \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1 + \dots + \alpha_k + c_1 + \dots + c_k)} \end{aligned}$$

b) According to Equation 3.5 we have

$$\pi(x_{n+1}, \dots, x_{n+r} \mid x_0, \dots, x_n) = \int \left[\prod_{s=1}^r P_{x_{n+s}, x_{n+s-1}}\right] \pi(P \mid x_0, \dots, x_n) \, dP$$

The distribution  $\pi(P \mid x_0, \ldots, x_n)$  is a product over one independent distribution for each  $P_i$ ,  $i = 1, \ldots, k$ . Also we can write

$$\prod_{s=1}^{r} P_{x_{n+s}, x_{n+s-1}} = \prod_{i=1}^{k} \prod_{j=1}^{k} P_{ij}^{d_{ij}}$$

These two things together mean that the integral splits into a product  $Q_1 Q_2 \cdots Q_k$ , where

$$Q_i = \int \prod_{j=1}^k P_{ij}^{d_{ij}} \pi(P_i \mid x_0, \dots, x_n) \, dP_i = \int \prod_{j=1}^k P_{ij}^{d_{ij}} \operatorname{Dirichlet}(P_i; \alpha_i + c_i) \, dP_i.$$

- c) Comparing with (a), we see that  $Q_i$  is the same as the posterior predictive for a sequence of observations with counts  $c_i = (c_{i1}, \ldots, c_{ik})$ . Thus we can use the predictive distribution found in (a), but plugging in the posterior Dirichlet $(\alpha_i + c_i)$ . The result follows.
- d) We found in Question 3 the counts of transitions in the table below:

	1	2	3	4
1	2	3	0	0
2	3	1	4	0
3	0	4	1	3
4	0	0	3	1

Using the formula found in (c), we get that

$$Q_1 = \frac{\Gamma(1+3+1)}{\Gamma(1+3)} \cdot \frac{\Gamma(4+5)}{\Gamma(4+5+1)} = \frac{4!}{3!} \cdot \frac{8!}{9!} = \frac{4}{9}$$
$$Q_2 = \frac{\Gamma(1+1+1)}{\Gamma(1+1)} \cdot \frac{\Gamma(1+4+1)}{\Gamma(1+4)} \cdot \frac{\Gamma(4+8)}{\frac{4}{4}+8+2} = \frac{2!}{1!} \cdot \frac{5!}{4!} \cdot \frac{11!}{13!} = \frac{2 \cdot 5}{12 \cdot 13}$$

 $\begin{array}{rcl} Q_3 & = & 1 \\ Q_4 & = & 1 \end{array}$ 

Thus the answer becomes

$$Q_1 Q_2 Q_3 Q_4 = \frac{4 \cdot 2 \cdot 5}{9 \cdot 12 \cdot 13} = 0.02849$$

5.

6. For counts of transitions between states of X we get the table

	0	1	2
0	1	2	0
1	2	1	3
2	0	3	1

For counts of transitions from states of X to states of Y we get the table

	0	1	2
0	3	1	0
1	0	4	2
2	0	3	1

1. We get

$$\mathbf{E}(P) = \begin{bmatrix} 1.5/4.5 & 2.5/4.5 & 0.5/4.5\\ 2.5/7.5 & 1.5/7.5 & 3.5/7.5\\ 0.5/5.5 & 3.5/5.5 & 1.5/5.5 \end{bmatrix}$$

2. We get

$$\mathbf{E}(Q) = \begin{bmatrix} 5/10 & 3/10 & 2/10\\ 2/12 & 6/12 & 4/12\\ 2/10 & 5/10 & 3/10 \end{bmatrix}$$

- 3. As  $X_{13} = 0$  we can read off the expectation of P that the posterior is the probability vector (1.5/4.5, 2.5/4.5, 0.5/4.5).
- 4. We get the posterior probability distribution

 $\begin{bmatrix} \frac{1.5}{4.5}, \frac{2.5}{4.5}, \frac{0.5}{4.5} \end{bmatrix} \begin{bmatrix} 5/10 & 3/10 & 2/10\\ 2/12 & 6/12 & 4/12\\ 2/10 & 5/10 & 3/10 \end{bmatrix} = \begin{bmatrix} 0.28148, 0.43333, 0.28518 \end{bmatrix}$ 

## A.3 Exercises from Chapter 4

## APPENDIX B

# List of some probability distributions

## The Bernoulli distribution

If  $x \in \{0,1\}$  has a Bernoulli(p) distribution, with  $0 \le p \le 1$ , then the probability mass function is

$$\pi(x) = p^x (1-p)^{1-x}.$$

R: Use the Binomial with sample size 1.

## The Beta distribution

If  $x \ge 0$  has a Beta $(\alpha, \beta)$  distribution with  $\alpha > 0$  and  $\beta > 0$  then the density is

$$\pi(x \mid \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha - 1} (1 - x)^{\beta - 1}.$$

R:dbeta, p<br/>beta, q<br/>beta, r<br/>beta

## The Beta-binomial distribution

If  $x \in \{0, 1, 2..., n\}$  has a Beta-binomial $(n, \alpha, \beta)$  distribution with n a positive integer,  $\alpha > 0$ , and  $\beta > 0$ , then the probability mass function is

$$\pi(x \mid n, \alpha, \beta) = \binom{n}{x} \frac{B(x + \alpha, n - x + \beta)}{B(\alpha, \beta)}$$

where B is the Beta function.

## The Binomial distribution

If  $x \in \{0, 1, 2, ..., n\}$  has a Binomial(n, p) distribution, with n a positive integer and  $0 \le p \le 1$ , then the probability mass function is

$$\pi(x \mid n, p) = \binom{n}{x} p^x (1-p)^{n-x}.$$

R: dbinom, pbinom, qbinom, rbinom

## The Cauchy distribution

If  $x \ge 0$  has a Cauchy $(\mu, \gamma)$  distribution, with  $\gamma > 0$ , then the probability density is

$$\pi(x \mid \mu, \gamma) = \frac{1}{\pi \gamma \left(1 + \left(\frac{x-\mu}{\gamma}\right)^2\right)}.$$

The standard Cauchy distribution with  $\mu = 0$  is the t-distribution with  $\nu = 1$ . R: deauchy, peauchy, reauchy

## The Dirichlet distribution

A vector  $\theta = (\theta_1, \ldots, \theta_k)$  of non-negative real numbers satisfying  $\sum_{i=1}^k \theta_i = 1$  has a Dirichlet  $(\alpha_1, \ldots, \alpha_k)$  distribution with parameter vector  $\alpha = (\alpha_1, \ldots, \alpha_k)$ , with each  $\alpha_i > 0$ , if it has probability density function

$$\pi(\theta \mid \alpha) = \frac{\Gamma(\alpha_1 + \alpha_2 + \dots + \alpha_k)}{\Gamma(\alpha_1)\Gamma(\alpha_2) \cdot \Gamma(\alpha_k)} \theta_1^{\alpha_1 - 1} \theta_2^{\alpha_2 - 1} \cdots \theta_k^{\alpha_k - 1}$$

If  $\theta$  has the distribution above, the expectation is the vector  $\frac{\alpha}{\sum_{i=1}^{k} \alpha_i}$ .

## The Exponential distribution

If  $x \ge 0$  has an Exponential( $\lambda$ ) distribution with  $\lambda > 0$  as parameter, then the density is

$$\pi(x \mid \lambda) = \lambda \exp(-\lambda x)$$

and the cumulative distribution function is

$$F(x) = 1 - \exp(-\lambda x).$$

R: dexp, pexp, qexp, rexp

## The Gamma distribution

If x > 0 has a Gamma $(\alpha, \beta)$  distribution, with  $\alpha > 0$  and  $\beta > 0$ , then the density is

$$\pi(x \mid \alpha\beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x).$$

The expectation and variance are  $\alpha/\beta$  and  $\alpha/\beta^2$ , respectively, while the mode is  $(\alpha - 1)/\beta$  (when  $\alpha \ge 1$ ). R: dgamma, ggamma, ggamma, rgamma

## The Geometric distribution

If the non-negative integer x has a Geometric distribution with parameter  $p \in [0, 1]$ , its probability mass function is given by

$$\pi(x \mid p) = (1 - p)^x p.$$

R: dgeom, pgeom, qgeom, rgeom

62

## The Multinomial distribution

A vector  $x = (x_1, \ldots, x_k)$  of non-negative integers satisfying  $\sum_{i=1}^k x_i = n$  has a Multinomial  $(n, p_1, \ldots, p_k)$  distribution with parameters n and  $p = (p_1, \ldots, p_k)$ , where n > 0 is an integer and  $p_i \ge 0$  and  $\sum_{i=1}^k p_i = 1$ , if the probability mass function is given by

$$\pi(x \mid n, p) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}.$$

## The Negative Binomial distribution

A stochastic variable x taking on as possible values any positive integer has a Negative Binomial distribution if its probability mass function is given by

$$\pi(x \mid r, p) = \binom{x+r-1}{x} \cdot (1-p)^r p^x = \frac{\Gamma(x+r)}{\Gamma(x+1)\Gamma(r)} (1-p)^x p^r$$

where r > 0 and  $p \in (0, 1)$  are parameters. R: dnbinom, pnbinom, qnbinom, rnbinom

## The Normal distribution

If the real x has a Normal distribution with parameters  $\mu$  and  $\sigma^2$ , its density is given by

$$\pi(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

R: dnorm, pnorm, qnorm, rnorm

## The Pareto distribution

If the real number  $x \in [M, \infty)$  has a Pareto $(M, \alpha)$  distribution with parameters M > 0 and  $\alpha > 0$ , its density on this interval is given by

$$\pi(x \mid M, \alpha) = \alpha M^{\alpha} x^{-(\alpha+1)}$$

## The Poisson distribution

If the nonnegative integer x has a  $Poisson(\lambda)$  distribution with parameter  $\lambda > 0$ , its probability mass function is given by

$$\pi(x \mid \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

R: dpois, ppois, qpois, rpois

## The t-distribution

If the real number x has a  $t(\nu)$  distribution with parameter  $\nu > 0$ , its density is

$$\pi(x \mid \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}.$$

R: dt, pt, qt, rt

## The Uniform distribution

If  $x \in [a, b]$  has a Uniform(a, b) distribution with b > a, then the density is given by

$$\pi(x \mid a, b) = \frac{1}{b-a}.$$

R: dunif, punif, qunif, runif
## APPENDIX C

## List of some conjugacies

Note: More conjucacies can be found on the Wikipedia page "Conjugate priors".

#### The Beta Binomial conjugacy

Likelihood:  $x \sim \text{Binomial}(n, \theta)$  **P**rior:  $\theta \sim \text{Beta}(\alpha, \beta)$  **P**osterior:  $\theta \mid x \sim \text{Beta}(\alpha + x, \beta + n - x)$ **P**rior predictive:  $x \sim \text{Beta-binomial}(n, \alpha, \beta)$ 

#### The Exponential Gamma conjugacy

Likelihood:  $x \sim \text{Exponential}(\theta)$ Prior:  $\theta \sim \text{Gamma}(\alpha, \beta)$ Posterior:  $\theta \mid x \sim \text{Gamma}(\alpha + 1, \beta + x)$ 

#### The Multinomial Dirichlet conjugacy

Likelihood:  $x = (x_1, \ldots, x_k) \sim \text{Multinomial}(n, \theta_1, \ldots, \theta_k)$ Prior:  $\theta = (\theta_1, \ldots, \theta_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k)$ Posterior:  $\theta \mid x \sim \text{Dirichlet}(\alpha_1 + x_1, \ldots, \alpha_k + x_k)$ 

#### The Poisson Gamma conjucacy

 $\begin{array}{l} \text{Likelihood: } x \sim \text{Poisson}\left(\theta\right) \\ \textbf{Prior: } \theta \sim \text{Gamma}(\alpha,\beta) \\ \textbf{Posterior: } \theta \mid x \sim \text{Gamma}(\alpha+x,\beta+1) \\ \textbf{Prior predictive: } x \sim \text{Negative-Binomial}\left(\alpha,\beta/(1+\beta)\right) \end{array}$ 

#### The Normal-Gamma conjucacy

Likelihood:  $x \sim \text{Normal}(\mu, \theta^{-1})$  **P**rior:  $\theta \sim \text{Gamma}(\alpha, \beta)$ **P**osterior:  $\theta \mid x \sim \text{Gamma}(\alpha + \frac{1}{2}, \beta + \frac{1}{2}(x - \mu)^2)$ 

### The Normal-Normal conjugacy

Likelihood:  $x \sim \text{Normal}(\theta, \tau^{-1})$ Prior:  $\theta \sim \text{Normal}(\mu, \tau_0^{-1})$ Posterior:  $\theta \mid x \sim \text{Normal}\left(\frac{\tau x + \tau_0 \mu}{\tau + \tau_0}, \frac{1}{\tau + \tau_0}\right)$ Prior predictive:  $x \sim \text{Normal}(\mu, \tau^{-1} + \tau_0^{-1})$ Computations:

$$\pi(\theta \mid x) \propto_{\theta} \pi(x \mid \theta)\pi(\theta)$$

$$\propto_{\theta} \exp\left(-\frac{\tau}{2}(x-\theta)^{2}\right)\exp\left(-\frac{\tau_{0}}{2}(\theta-\mu)^{2}\right)$$

$$= \exp\left(-\frac{1}{2}\left[\tau x^{2} - 2\tau x\theta + \tau\theta^{2} + \tau_{0}\theta^{2} - 2\tau_{0}\theta\mu + \tau_{0}\mu^{2}\right]\right)$$

$$\propto_{\theta} \exp\left(-\frac{1}{2}\left[(\tau+\tau_{0})\theta^{2} - 2(\tau x + \tau_{0}\mu)\theta\right]\right)$$

$$\propto_{\theta} \exp\left(-\frac{1}{2}(\tau+\tau_{0})\left(\theta - \frac{\tau x + \tau_{0}\mu}{\tau+\tau_{0}}\right)^{2}\right)$$

$$\propto_{\theta} \operatorname{Normal}\left(\theta; \frac{\tau x + \tau_{0}\mu}{\tau+\tau_{0}}, \frac{1}{\tau+\tau_{0}}\right)$$

# Bibliography

- [Axe10] Axelson-Fisk, M. Comparative Gene Finding. Springer, 2010.
- [Dob16] Dobrow, R. P. Introduction to stochastic processes with R. John Wiley & Sons, 2016.
- [IRW12] Insua, D., Ruggeri, F. and Wiper, M. Bayesian analysis of stochastic process models. Vol. 978. John Wiley & Sons, 2012.