

# Parameter Learning and Prediction with Bayesian Belief Networks

Georg Heimel and Viktor Atalla

February 20, 2018

## **Abstract**

Some thoughts on the matter.

# Chapter 1

## One Random Variable

### 1.1 Parameter or parameter distribution?

Suppose we have a single, discrete random variable  $X$  with probability mass function  $P(X)$ . The trivial belief network illustrating this variable is shown in Fig. 1.1. Strictly speaking, this graph illustrates the situation where  $P(X)$  is



Figure 1.1: Belief network of a single, discrete random variable  $X$  with probability mass function  $P(X)$ .

given in free, unparameterized form, one number for each discrete value  $x$  that  $X$  can take. If we want to visually represent a *parameterized* probability mass function with parameter  $\theta$ , then we have to introduce another node for it into the graph, as shown in Fig. 1.2. However, because the nodes in belief networks



Figure 1.2: A single random variable  $X$ , whose probability mass function is parameterized by  $\theta$ , which is itself a random variate.

represent probability mass functions (or densities), not numbers, this graph represents a joint probability distribution  $P(X, \theta)$  of *two* random variables, not just one. In fact, what we have drawn designates  $\theta$  as a (usually continuous) random variate as well. Furthermore, the graph encodes a particular factorization of the joint probability distribution.

$$P(X, \theta) = P(X|\theta)P(\theta) \quad (1.1)$$

To recover the (now *marginal*) probability mass function  $P(X)$ , we would

have to integrate out  $\theta$ .

$$P(X) = \int d\theta P(X, \theta) = \int d\theta P(X|\theta)P(\theta) \quad (1.2)$$

In order to do that, however, we would need to know and specify the probability density  $P(\theta)$ . If, in contrast, we just wanted to express the concept that  $P(X)$  is a probability mass function parameterized by a *known*, well-defined value  $\theta^*$ , then we have to add one more element to the graph, as shown in Fig. 1.3. Because  $\theta^*$  is a constant and no probability distribution can be assigned to it,

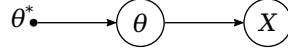


Figure 1.3: A single random variable  $X$ , whose probability mass function is parameterized by  $\theta$ , which is set to the number  $\theta^*$ .

above graph implies the following factorization of the overall joint probability.

$$P(X, \theta) = P(X|\theta)P(\theta|\theta^*) \quad (1.3)$$

In order to *set* the value of  $\theta$  to  $\theta^*$ , we chose a *delta distribution* centered on  $\theta^*$  for the last term on the right-hand side.

$$P(\theta|\theta^*) = \delta(\theta - \theta^*) \quad (1.4)$$

Substituting 1.4 into Eq. 1.3 and integrating out  $\theta$  yields:

$$P(X) = \int d\theta P(X|\theta) \delta(\theta - \theta^*) = P(X|\theta^*) \quad (1.5)$$

Trying to visually represent the result of this integration, *i.e.*, the very right-hand side term, leads to the graph in Fig. 1.4. It literally spells out the parameterized probability mass function  $P(X|\theta^*)$ .

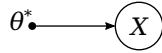


Figure 1.4: A single random variable  $X$ , whose probability mass function is parameterized by the *a priori* known and fixed number  $\theta^*$ .

Returning to Eq. 1.2, we could also keep  $\theta$  as a random variate but chose a parameterized distribution for it as well. Calling the parameters, again fixed and known numbers,  $\alpha$  and  $\beta$ , the graph in Fig. 1.5 then encodes the factorization:

$$P(X, \theta) = P(X|\theta)P(\theta|\alpha, \beta) \quad (1.6)$$

If we wanted, for example, to learn the distribution of  $\theta$  from the distribution of  $X$  with our *prior* belief encoded in  $P(\theta|\alpha, \beta)$ , we would factorize the left-hand side of Eq. 1.6 as  $P(\theta|X)P(X)$  and divide by  $P(X)$  to arrive at:

$$P(\theta|X) = \frac{P(X|\theta)P(\theta|\alpha, \beta)}{P(X)} \quad (1.7)$$

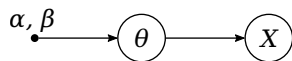


Figure 1.5: A single random variable  $X$ , whose probability mass function is parameterized by  $\theta$ , itself given by a distribution parameterized by  $\alpha$  and  $\beta$ .

## 1.2 Learning the parameter

Clearly, if we want to predict the outcome of a future *realization* of  $X$ , we need some idea of what the number  $\theta^*$  or the probability density  $P(\theta)$  should be. To learn either, we will first have to experiment and observe a number  $N$  of realizations  $x$  of  $X$  and the learned parameter (or its distribution) will depend on the specific outcome of our experiment. Trying to preserve a strict one-to-one correspondence between formulas and graphs, the repeated observations are modeled as a random process, where there are, in fact,  $N$  separate random variates  $X_{i=1\dots N}$ , all with the same probability mass function and all parameterized by the very same  $\theta$ , which we take to be non-parametrically distributed to keep things simple for now.

We start with a trivial factorization of the joint probability distribution of all these random variables,

$$P(\theta|X_1, X_2, \dots, X_N)P(X_1, X_2, \dots, X_N) = P(\theta, X_1, X_2, \dots, X_N) \quad (1.8)$$

and divide by the second term on the left-hand side.

$$P(\theta|X_1, X_2, \dots, X_N) = \frac{P(\theta, X_1, X_2, \dots, X_N)}{P(X_1, X_2, \dots, X_N)} \quad (1.9)$$

If each  $X_i$  is independent of all others *conditioned* on  $\theta$ , that is, if the observations are *i.i.d.*, we draw the graph for what we are about to see in our repeated experiment as shown in Fig. 1.6; the fact that circles representing the (conditional) probability mass functions of the  $X_i$  are shaded indicates that these are *observable* variables, while  $\theta$  is a *hidden* variable. Using the factorization implied by the graph for the joint probability distribution in the numerator of the

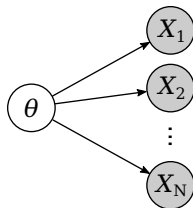


Figure 1.6: Multiple, independent instances of a single, observable random variable  $X$ , whose probability mass functions are all parameterized by the same  $\theta$ , which is itself a random variate, but with a non-parametric distribution.

right-hand side of Eq. 1.9, we obtain:

$$P(\theta|X_1, X_2, \dots, X_N) = \frac{P(\theta) \prod_{i=1}^N P(X_i|\theta)}{P(X_1, X_2, \dots, X_N)} \quad (1.10)$$

Now we start our experiment and observe actual values for  $X$ , that is, first we observe that  $X_1$  takes on the value  $x_1$ ,  $X_2 = x_2$  and so on until we see that  $X_N = x_N$ . Replacing all variables by their observed values in Eq. 1.10, leads us to:

$$P(\theta|x_1, x_2, \dots, x_N) = \frac{P(\theta) \prod_{i=1}^N P(x_i|\theta)}{P(x_1, x_2, \dots, x_N)} \quad (1.11)$$

The set of observed realizations  $\{x_i\}$  is commonly referred to as *the data*, denoted by  $\mathcal{D}$ . Using this notation, we finally obtain:

$$P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})} \quad (1.12)$$

Now, given a suitable expression for  $P(x_i|\theta)$ , the *likelihood* of  $\theta$  given the observation  $x_i$ , and a choice for the *prior* probability distribution  $P(\theta)$ , we can learn the *posterior* probability distribution of  $\theta$ . In general, however, computing the denominator in Eq. 1.12, the *evidence* for the data, is quite cumbersome. Therefore, one often resorts to taking the value  $\theta^*$ , where the product  $P(\theta|\mathcal{D})P(\theta)$  has a maximum, as a *point estimate* for  $\theta$ . Specifically, we take the *maximum likelihood estimate* (MLE) if we chose a flat prior  $P(\theta) = \text{const.}$

$$\theta_{\text{MLE}} = \text{argmax}_{\theta} P(\mathcal{D}|\theta) \quad (1.13)$$

and the *maximum posterior* (MAP) estimate if we don't.

$$\theta_{\text{MAP}} = \text{argmax}_{\theta} P(\mathcal{D}|\theta)P(\theta) \quad (1.14)$$

Mathematically, this again corresponds to *collapsing* the posterior to a delta distribution  $\delta(\theta - \theta^*)$  with its center  $\theta^*$  either  $\theta_{\text{MLE}}$  or  $\theta_{\text{MAP}}$ . Depending on whether we do or do not chose to replace the posterior distribution with a point estimate, we end up with the graphs shown in Figs. 1.4 and 1.2, respectively, with the difference that we now have *learned* the parameter  $\theta$  (or its probability density), and not just set it.

Note that, if we wanted to express our prior belief of  $P(\theta)$  with a parameterized probability density, we would extend all graphs in the present Section like we demonstrated in Fig. 1.5 and replace  $P(\theta)$  by  $P(\theta|\alpha, \beta)$  in all formulas accordingly (*cf.* Eq. 1.7).

## 1.3 Predicting on new data

### 1.3.1 The short answer

We can now proceed to predict the probability mass function for the outcome of a new realization of our univariate random process. To distinguish this event from

past observations, we introduce yet another random variate,  $X^*$  this time. Now, if we chose a (learned) point estimate  $\theta^*$  for the parameter of the probability mass function (*e.g.*, MLE or MAP), we have to look no further than Fig. 1.4 and Eq. 1.5.

$$P(X^*) = P(X^*|\theta^*) \quad (1.15)$$

If, however, we decided to keep the entire posterior distribution of  $\theta$ , then we are in the situation depicted in Fig. 1.2 and, consequently, have to employ Eq. 1.2.

$$P(X^*) = \int d\theta P(X^*|\theta)P(\theta) \quad (1.16)$$

### 1.3.2 The long answer

To explicitly express that  $P(\theta)$  is indeed our *posterior* belief about the distribution of  $\theta$ , and that we are indeed predicting the probability mass function for a new observation *after* having seen the data, we explicitly introduce  $X^*$  into the graph in Fig. 1.6 to produce the one shown in Fig. 1.7. Collecting all the  $X_i$

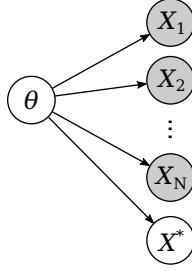


Figure 1.7: Same as Figure 1.6 with  $X^*$  representing a new, unobserved instance of  $X$ .

as  $\mathcal{D}$ , the joint probability of all three parameters  $X^*$ ,  $\theta$ , and  $\mathcal{D}$ , can then be factorized in two ways which are, of course, equal to each other. Specifically, we chose:

$$P(X^*, \theta|\mathcal{D})P(\mathcal{D}) = P(X^*, \mathcal{D}|\theta)P(\theta) \quad (1.17)$$

Because all observations of  $X$ , past or future, are independent of each other conditioned on  $\theta$ ,  $X^*$  is also independent from  $\mathcal{D}$ . We can, therefore, rewrite the right-hand side,

$$P(X^*, \theta|\mathcal{D})P(\mathcal{D}) = P(X^*|\theta)P(\mathcal{D}|\theta)P(\theta) \quad (1.18)$$

join its last two terms,

$$P(X^*, \theta|\mathcal{D})P(\mathcal{D}) = P(X^*|\theta)P(\mathcal{D}, \theta) \quad (1.19)$$

take them apart again the other way around,

$$P(X^*, \theta|\mathcal{D})P(\mathcal{D}) = P(X^*|\theta)P(\theta|\mathcal{D})P(\mathcal{D}) \quad (1.20)$$

and divide by  $P(\mathcal{D})$ .

$$P(X^*, \theta | \mathcal{D}) = P(X^* | \theta) P(\theta | \mathcal{D}) \quad (1.21)$$

Now integrating over  $\theta$  brings us to the formal expression for the *posterior predictive*.

$$P(X^* | \mathcal{D}) = \int d\theta P(X^* | \theta) P(\theta | \mathcal{D}) \quad (1.22)$$



## Chapter 2

# Two Random Variables

### 2.1 Repeat after me ...

Let  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{Z}$  be three non-intersecting sets of nodes in the directed graph  $G$ , that is  $\mathcal{X} \cap \mathcal{Y} = \mathcal{X} \cap \mathcal{Z} = \mathcal{Y} \cap \mathcal{Z} = \emptyset$ . We define  $\mathcal{X}$  and  $\mathcal{Y}$  to be *d-separated* by  $\mathcal{Z}$  in  $G$  iff all paths between every  $x \in \mathcal{X}$  and every  $y \in \mathcal{Y}$  are *blocked*. A path  $U$  is said to be blocked iff there is a node  $w$  on  $U$  such that either

- $w$  is a collider and neither  $w$  nor any of its descendants is in  $\mathcal{Z}$ , or
- $w$  is not a collider on  $U$  and  $w$  is in  $\mathcal{Z}$ .

If, in particular,  $G$  represents a *belief network* and, thus, the nodes represent random variates, then  $\mathcal{X}$  and  $\mathcal{Y}$  are independent conditional on  $\mathcal{Z}$  in all probability distributions  $G$  can represent iff  $\mathcal{X}$  and  $\mathcal{Y}$  are *d-separated* by  $\mathcal{Z}$ .

### 2.2 Parameter separation

Moving on, we now suppose we have two random variables,  $X$  and  $Y$ , parameterized by  $\theta_x$  and  $\theta_y$ , respectively. To make things more interesting, we will further suppose that they are not independent of each other but that, on the contrary, the value that  $X$  takes impacts the distribution of values observed for  $Y$ . Because our distributions are parameterized, affecting the probability mass function of  $Y$  means affecting the parameter  $\theta_y$  that defines it. The graph encoding the situation just described, shown in Fig. 2.1, suggest the following factorization of the joint distribution of all involved random variates.

$$P(Y, X, \theta_y, \theta_x) = P(Y|\theta_y, X)P(\theta_y|X)P(X|\theta_x)P(\theta_x) \quad (2.1)$$

In an attempt to find a second, complementary factorization that would allow us to find the analogue to Eq. 1.12 for two parameters, we start by conditioning on  $X$ .

$$P(Y, \theta_y, \theta_x|X)P(X) = rhs \quad (2.2)$$

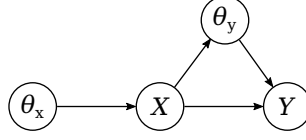


Figure 2.1: Two correlated random variables,  $X$  and  $Y$ , with their probability mass functions parameterized by  $\theta_x$  and  $\theta_y$ , which are themselves random variables, but with non-parametric distributions.

Using the recipe from Section 2.1, we then realize that both  $Y$  and  $\theta_y$  are independent of  $\theta_x$  conditioned on  $X$  (because  $X$  is on any path between them, is *not* a collider, and *is* in the conditioning set).

$$P(Y, \theta_y | X) P(\theta_x | X) P(X) = rhs \quad (2.3)$$

Next, we combine the lone  $P(X)$  with the first factor,

$$P(Y, X, \theta_y) P(\theta_x | X) = rhs \quad (2.4)$$

and tear that factor apart in a different way.

$$P(\theta_y | Y, X) P(Y, X) P(\theta_x | X) = rhs \quad (2.5)$$

Writing out again the right-hand side of Eq. 2.1 and dividing by  $P(Y, X) = P(Y|X)P(X)$  then yields:

$$P(\theta_y | Y, X) P(\theta_x | X) = \frac{P(Y | \theta_y, X) P(\theta_y | X) P(X | \theta_x) P(\theta_x)}{P(Y | X) P(X)} \quad (2.6)$$

Comparing this to Eq. 1.12, we see that we got what we asked for, sort of. But how are we going to learn the parameters  $\theta_x$  and  $\theta_y$  at the same time? The answer is, we aren't because Eq. 2.6 is actually the product of two equations, one for each parameter, and we need to separate it to learn them one by one.

First, we combine the last two terms in the numerator on the right-hand side of Eq. 2.6 to  $P(\theta_x, X)$  and take that apart the other way around.

$$P(\theta_y | Y, X) P(\theta_x | X) = \frac{P(Y | \theta_y, X) P(\theta_y | X) P(\theta_x | X) P(X)}{P(Y | X) P(X)} \quad (2.7)$$

Then we realize that the newly created factors cancel with existing ones only to yield,

$$P(\theta_y | Y, X) = \frac{P(Y | \theta_y, X) P(\theta_y | X)}{P(Y | X)} \quad (2.8)$$

which we recognize as the "first half" of Eq. 2.6. We will use this expression to learn the  $X$ -dependent parameter  $\theta_y$ .

To isolate an equivalent expression for  $\theta_x$ , we start again with Eq. 2.6 and multiply with  $P(X)$ ,

$$P(\theta_y | Y, X) P(X) P(\theta_x | X) = \frac{P(Y | \theta_y, X) P(\theta_y | X) P(X) P(X | \theta_x) P(\theta_x)}{P(Y | X) P(X)} \quad (2.9)$$

combine the first three terms in the numerator on the right-hand side,

$$P(\theta_y|Y, X)P(X)P(\theta_x|X) = \frac{P(Y, X, \theta_y)P(X|\theta_x)P(\theta_x)}{P(Y|X)P(X)} \quad (2.10)$$

multiply by  $P(Y|X)$ ,

$$P(\theta_y|Y, X)P(Y|X)P(X)P(\theta_x|X) = \frac{P(Y, X, \theta_y)P(X|\theta_x)P(\theta_x)}{P(X)} \quad (2.11)$$

and combine the first three terms on the left-hand side,

$$P(Y, X, \theta_y)P(\theta_x|X) = \frac{P(Y, X, \theta_y)P(X|\theta_x)P(\theta_x)}{P(X)} \quad (2.12)$$

which cancel with their right-hand side equivalent to give us,

$$P(\theta_x|X) = \frac{P(X|\theta_x)P(\theta_x)}{P(X)} \quad (2.13)$$

which we recognize as the "second half" of Eq. 2.6.

## 2.3 Learning the parameters

### 2.3.1 The short answer

Looking at Eq. 2.2, we realize that we can proceed in exactly the same manner as for the univariate case discussed in Section 1.2 to learn the parameter  $\theta_x$ . We denote the obtained posterior  $P(\theta_x|\mathcal{D}_x)$  because we do not, in fact, need the values observed for  $Y$  to learn  $\theta_x$ , but only the set  $\{\mathbf{x}_i\} = \mathcal{D}_x$  of realizations of  $X$ . To learn  $\theta_y$ , we take another look at Eq. 2.8 and realize that everything is conditioned on  $X$ , by definition a *discrete* random variable that can take on the values  $x_1, x_2, \dots$ . So, instead of just one equation, we are actually looking at as many as there are different values for  $X$ .

$$\begin{aligned} P(\theta_y|Y, X = x_1) &= \frac{P(Y|\theta_y, X = x_1)P(\theta_y|X = x_1)}{P(Y|X = x_1)} \\ P(\theta_y|Y, X = x_2) &= \frac{P(Y|\theta_y, X = x_2)P(\theta_y|X = x_2)}{P(Y|X = x_2)} \\ &\vdots \end{aligned} \quad (2.14)$$

This suggest that, having collected paired realizations of  $X$  and  $Y$ , we group by the possible values of  $X$ , specify a prior  $P(\theta_y|X)$  for each of them, and compute a separate posterior for  $\theta_y$  from the  $Y$  realizations in each group, again following Section 1.2.

### 2.3.2 The long answer

Let's start by staring at Eq. 2.8 yet again.

$$P(\theta_y|Y, X) = \frac{P(Y|\theta_y, X)P(\theta_y|X)}{P(Y|X)}$$

It was derived by rigorously applying the rules for manipulating probability distributions but, to actually use it, we now must understand just what each term really means. Starting with the first of the two factors in the numerator, we acknowledge that it gives us the possibility to specify a different likelihood for each value  $x$  that  $X$  can take. While we might not actually want to do that, we still note that the product

$$P(Y|\theta_y, X) = \prod_k p_k(Y|\theta_y)^{[X=x_k]} \quad (2.15)$$

does the trick of filtering out the likelihood  $p_k(Y|\theta_y)$  by setting  $X$  to  $x_k$ . This is thanks to the *Iverson bracket* in the exponent, which evaluates to 1 if its contents are true and to 0 otherwise.

Likewise, the second term in the numerator intuitively expresses that a different prior  $\pi(\theta_y)$  can be chosen for every value  $x$  that  $X$  can take. This is again expressed in the same way,

$$P(\theta_y|X) = \prod_k \pi_k(\theta_y)^{[X=x_k]} \quad (2.16)$$

which yields for the entire numerator:

$$P(Y|\theta_y, X)P(\theta_y|X) = \prod_k \left\{ p_k(Y|\theta_y) \pi_k(\theta_y) \right\}^{[X=x_k]} \quad (2.17)$$

Finally, the denominator is the integral of the numerator over  $\theta_y$ .

$$\begin{aligned} P(Y|X) &= \int d\theta_y \prod_k \left\{ p_k(Y|\theta_y) \pi_k(\theta_y) \right\}^{[X=x_k]} \\ &= \prod_k \left\{ \int d\theta_y p_k(Y|\theta_y) \pi_k(\theta_y) \right\}^{[X=x_k]} \\ &\stackrel{!}{=} \prod_k p_k(Y)^{[X=x_k]} \end{aligned} \quad (2.18)$$

Putting it all together, then leads us to realize that the conditional posterior

$$\begin{aligned} P(\theta_y|Y, X) &= \prod_k \left\{ \frac{p_k(Y|\theta_y) \pi_k(\theta_y)}{p_k(Y)} \right\}^{[X=x_k]} \\ &\stackrel{!}{=} \prod_k p_k(\theta_y|Y)^{[X=x_k]} \end{aligned} \quad (2.19)$$

is actually a product of as many posteriors as there are values that  $X$  can take. Our goal, of course, is to learn every single one of them.

### Extending the graph

Having grasped the multitude of posteriors we want to learn as well as their functional relationships, we now take a step back, and extend the graph in Fig. 2.1 as we did in Section 1.2 to include multiple instances of the two random variates that are to be observed. Note that the monstrosity in Fig. 2.2 seems

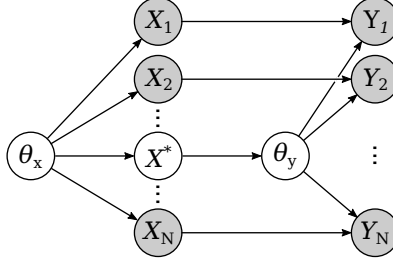


Figure 2.2: Multiple, independent instances of two correlated random variables,  $X$  and  $Y$ , with their probability mass functions parameterized by  $\theta_x$  and  $\theta_y$ , which are themselves random variables, but with non-parametric distributions.

to lie somewhere between Figs. 1.6 and 1.7. We will see soon enough that the extra, unobserved  $X$  is necessary to capture the conditional dependence of the posterior in  $\theta_y$ .

This graph then encodes the following factorization of the joint probability of all involved variables (*cf.* Eq. 2.1):

$$lhs = \left\{ \prod_{i=1}^N P(Y_i | \theta_y, X_i) \right\} P(\theta_y | X^*) \left\{ \prod_{i=1}^N P(X_i | \theta_x) \right\} P(X^* | \theta_x) P(\theta_x) \quad (2.20)$$

To find the matching left-hand side, we first redo the exact same derivation as in Section 2.2, with the entire set of all  $\{Y_i\}$  taking the place of the lone  $Y$  and, likewise, the entire set of all  $\{X_i\}$  plus the  $X^*$  substituted for the lone  $X$ . This first leads us to the equivalent of Eq. 2.5.

$$P(\theta_y | \{Y_i\}, \{X_i\}, X^*) P(\theta_x | \{X_i\}, X^*) P(\{Y_i\}, \{X_i\}, X^*) = rhs \quad (2.21)$$

Also the next step is analogous to the preceding Section 2.2 in that we factorize the last term into  $P(\{Y_i\} | \{X_i\}, X^*) P(\{X_i\}, X^*)$ , which allows us again to separate into two parts.

$$P(\theta_y | \{Y_i\}, \{X_i\}, X^*) = \frac{P(\theta_y | X^*) \prod_{i=1}^N P(Y_i | \theta_y, X_i)}{P(\{Y_i\} | \{X_i\}, X^*)} \quad (2.22)$$

$$P(\theta_x | \{X_i\}, X^*) = \frac{P(X^* | \theta_x) P(\theta_x) \prod_{i=1}^N P(X_i | \theta_x)}{P(\{X_i\}, X^*)} \quad (2.23)$$

### Posterior of $\theta_x$

The appearance of the auxiliary quantity  $X$  in Eq. 2.23 is sitting somewhat oddly in the above expression for the posterior of  $\theta_x$ . Why should one of the instances of  $X$  stick out from the rest? Fortunately, it is not hard to make it disappear. To do so, we multiply with the denominator of the right-hand side,

$$P(\theta_x, \{X_i\}, X^*) = P(X^*|\theta_x) P(\theta_x) \prod_{i=1}^N P(X_i|\theta_x) \quad (2.24)$$

and pluck apart the left-hand side by conditioning on  $\theta_x$ .

$$P(\{X_i\}, X^*|\theta_x) P(\theta_x) = rhs \quad (2.25)$$

Examining the graph in Fig. 2.2 in the light of Section 2.1, we realize that paths from  $X^*$  to any of the  $\{X_i\}$  either lead us through  $\theta_x$ , which is not a collider but is in the conditioning set, or through one of the  $\{Y_i\}$ , which are all colliders and are not in the conditioning set. Thus, all these paths are blocked and we can further write:

$$\begin{aligned} P(\{X_i\}|\theta_x) P(X^*|\theta_x) P(\theta_x) &= rhs \\ P(\theta_x, \{X_i\}) P(X^*|\theta_x) &= rhs \\ P(\theta_x|\{X_i\}) P(\{X_i\}) P(X^*|\theta_x) &= rhs \end{aligned} \quad (2.26)$$

Now that we see how the offensive term  $P(X^*|\theta_x)$  cancels with its equivalent on the right-hand side, we can divide again by  $P(\{X_i\})$  and replace the variables by their realizations  $\{x_i\}$  to arrive at:

$$P(\theta_x|\{x_i\}) = \frac{P(\theta_x) \prod_{i=1}^N P(x_i|\theta_x)}{P(\{x_i\})} \quad (2.27)$$

If we denote, as before,  $\{x_i\} = \mathcal{D}_x$ , we finally end up with the familiar expression for the posterior of  $\theta_x$ .

$$P(\theta_x|\mathcal{D}_x) = \frac{P(\mathcal{D}_x|\theta_x) P(\theta_x)}{P(\mathcal{D}_x)} \quad (2.28)$$

### Posterior(s) of $\theta_y$

To make sense of Eq. 2.22, we first plug Eqs. 2.15 and 2.16 into the numerator (*num*) of the right-hand side.

$$num = \prod_j \pi_j(\theta_y)^{[X^*=x_k]} \prod_{i=1}^N \left\{ \prod_k p_k(Y_i|\theta_y)^{[X_i=x_k]} \right\} \quad (2.29)$$

Rearranging the products, we get,

$$\begin{aligned} num &= \prod_k \left\{ \pi_k(\theta_y) \prod_{i=1}^N p_k(Y_i | \theta_y)^{[X_i=x_k]} \right\}^{[X^*=x_k]} \\ &\stackrel{!}{=} \prod_k \left\{ p_k(\{Y_i\} | \theta_y, \{X_i\}) \pi_k(\theta_y) \right\}^{[X^*=x_k]} \end{aligned} \quad (2.30)$$

which leaves us with

$$\begin{aligned} den &= \prod_k \left\{ \int d\theta_y \pi_k(\theta_y) \prod_{i=1}^N p_k(Y_i | \theta_y)^{[X_i=x_k]} \right\}^{[X^*=x_k]} \\ &\stackrel{!}{=} \prod_k \left\{ p_k(\{Y_i\} | \{X_i\}) \right\}^{[X^*=x_k]} \end{aligned} \quad (2.31)$$

for the denominator (*den*). Putting it all together and substituting the actual realizations,  $\{y_i\}$  and  $\{x_i\}$ , of  $Y$  and  $X$ , respectively, we then arrive at:

$$\begin{aligned} P(\theta_y | \{y_i\}, \{x_i\}, X^*) &= \prod_k \left\{ \frac{p_k(\{y_i\} | \theta_y, \{x_i\}) \pi_k(\theta_y)}{p_k(\{y_i\} | \{x_i\})} \right\}^{[X^*=x_k]} \\ &\stackrel{!}{=} \prod_k \left\{ p_k(\theta_y | \{y_i\}, \{x_i\}) \right\}^{[X^*=x_k]} \end{aligned} \quad (2.32)$$

Introducing again  $\mathcal{D}_y = \{y_i\}$ ,  $\mathcal{D}_x = \{x_i\}$ , and  $\mathcal{D} = \{y_i, x_i\}$ , we could also rewrite this as,

$$P(\theta_y | \mathcal{D}, X) = \prod_k \left\{ p_k(\theta_y | \mathcal{D}) \right\}^{[X^*=x_k]} \quad (2.33)$$

where the relevant factors are given by:

$$p_k(\theta_y | \mathcal{D}) = \frac{p_k(\mathcal{D}_y | \theta_y, \mathcal{D}_x) \pi_k(\theta_y)}{p_k(\mathcal{D}_y | \mathcal{D}_x)} \quad (2.34)$$

**Proceed at your own risk from here!**

## 2.4 Predicting on new data

To predict the probability mass function for a new observation  $X^*$ , we proceed in exactly the same way as for the single-variable case outlined in Section 1.3 and rewrite Eq. 1.22.

$$P(X^* | \mathcal{D}) = \int d\theta_x P(X^* | \theta_x) P(\theta_x | \mathcal{D}) \quad (2.35)$$

The same is true if we wanted to predict the probability mass function for a new observation  $Y^*$  given a certain value of  $X^*$ .

$$P(Y^*|X^*, \mathcal{D}) = \int d\theta_y P(Y^*|\theta_y, X^*)P(\theta_y|\mathcal{D}, X^*) \quad (2.36)$$

Again, Eq. 2.36 is actually many equations, one for every value  $x^*$  that  $X^*$  can take. Multiplying with Eq. 2.35,

$$P(Y^*, X^*|\mathcal{D}) = P(Y^*|X^*, \mathcal{D})P(X^*|\mathcal{D}) \quad (2.37)$$

we get the joint posterior predictive,

$$P(Y^*, X^*|\mathcal{D}) = \int d\theta_y d\theta_x P(Y^*|\theta_y, X^*)P(\theta_y|\mathcal{D}, X^*)P(X^*|\theta_x)P(\theta_x|\mathcal{D}) \quad (2.38)$$

and that for a new observation  $X^*$  given a certain value of  $Y^*$ .

$$P(X^*|Y^*, \mathcal{D}) = \frac{P(Y^*, X^*|\mathcal{D})}{P(Y^*|\mathcal{D})} \quad (2.39)$$

The denominator on the right-hand side of this expression is obtained by summing the joint posterior predictive in Eq. 2.38 over all values  $\{x_k\}$  that  $X^*$  can possibly take.

$$P(Y^*|\mathcal{D}) = \sum_k P(Y^*, x_k|\mathcal{D}) \quad (2.40)$$

We note in passing that, if the discrete random variable  $X$  is relabeled  $C$  and "... has the value  $x$ ." is substituted with "... belongs to class  $c$ .", then Eq. 2.39 describes the recipe for predicting the *class probabilities* of new observation  $Y^*$ , and *classification* would mean assigning it to the class  $c^*$  with the highest probability.

## 2.5 Hyperparameters

Now that we understand the case of two dependent, parameterized random variables, we might want to include the possibility of specifying the (prior) distributions of these parameters with *hyperparameters*. Like in Fig. 1.5, adding two of them (why not  $\alpha$  and  $\beta$ ) for  $\theta_x$  is straightforward. When we want to express our prior belief about  $\theta_y$ , however, we should reserve the option to specify different parameters (why not  $\mu$  and  $\nu$ ) for each value  $x$  that  $X$  can take. This means adding an arrow from  $X$  to  $\mu$  and  $\nu$ , as shown in Fig. 2.3.

Without going through every step of our derivation again, we list here only how introducing hyperparameters changes a few key formulas.

## 2.6 Bayesian hypothesis test

Suppose  $X$  is a binary random variable, representing if a patient has recieved a certain treatment or not, and  $Y$  is also binary random variable, indicating



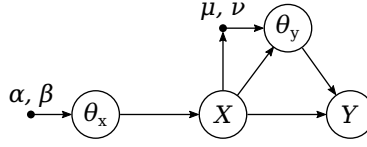


Figure 2.3: Two correlated random variables,  $X$  and  $Y$ , with their probability mass functions parameterized by  $\theta_x$  and  $\theta_y$ , which are themselves random variables defined through their respective hyperparameters  $(\alpha, \beta)$  and  $(\mu, \nu)$ .

whether patient's symptoms disappeared or not. Our goal is then to learn the distributions of the parameters,  $\theta_x$  and  $\theta_y$ , describing the probabilities of  $X$  and  $Y$  occurring, respectively.

As we are dealing with binary variables, we can choose between a *Beta-Bernoulli* and a *Beta-Binomial* model. As the names imply, the priors on  $\theta_x$  and  $\theta_y$  are *Beta* distributions in both models, each characterized by two parameters.

## Chapter 3

# Good Things Come in 3's

Example binary classification

## Chapter 4

# Continuous Variables

Density vs. mass function and the mystery of observation and hidden variables.