

# BE/Bi 103: Data Analysis in the Biological Sciences

Justin Bois

Caltech

Fall, 2016

## 8 Hierarchical models

In this lecture, we will investigate **hierarchical models**, in which some model parameters are dependent on others in specific ways. This is best learned by example.

### 8.1 A hierarchical model example

In [homework problem 4.2](#), we studied reversals under exposure to blue light in *C. elegans* with Channelrhodopsin in two different neurons. Let's consider one of the strains which contains a Channelrhodopsin in the ASH sensory neuron. The experiment was performed by the students of [Bi 1x](#) in 2015 and again in 2016. In 2015, we found that 9 out of 35 worms reversed under exposure to blue light. In 2016, 12 out of 35 reversed.

Considering for a moment only the 2015 experiment, we can use this measurement to estimate the probability  $p$  of reversal. Specifically, we found that the posterior probability of reversal given  $r$  out of  $n$  trials showed reversals was

$$P(p \mid r, n, I) = \begin{cases} \frac{(n+1)!}{(n-r)!r!} p^r (1-p)^{n-r} & 0 \leq p \leq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (8.1)$$

This posterior assumed a uniform prior  $P(p \mid I)$  on  $0 \leq p \leq 1$ , and a binomial likelihood,  $P(r \mid n, p, I)$ .

We did the experiment again in 2016, getting  $r = 12$  and  $n = 35$ . Actually, we could imagine doing the experiment over and over again, say  $k$  times, each time getting a value of  $r$  and  $n$ . Conditions may change from experiment to experiment. For example, we may have different lighting set-ups, slight differences in the strain of worms we're using, etc. We are left with some choices on how to model the data.

#### 8.1.1 Pooled data: identical parameters

We could pool all of the data together. In other words, let's say we measure  $r_1$  out of  $n_1$  reversals in the first set of experiments,  $r_2$  out of  $n_2$  reversals in the second set, etc., up to  $k$  total experiments. We could pool all of the data together to get

$$r = \sum_{i=1}^k r_i \text{ out of } n = \sum_{i=1}^k n_i \text{ reversals.} \quad (8.2)$$

We then compute our posterior as in equation (8.1). Here, the assumption is that the result in each experiment are governed by *identical parameters*. That is to say that we assume  $p_1 = p_2 = \dots = p_k = p$ .

This is similar to what we did in section 1.9, in which we looked at how a single hypothesis (or parameter value) is informed by more data. This is also exactly what you did in homework 4.

### 8.1.2 Independent parameters

As an alternative, we could instead say that the parameters in each experiment are totally independent of each other. In this case, we assume that  $p_1, p_2, \dots, p_k$  are all independent of each other. Thus, the posterior probability is

$$P(\mathbf{p} \mid \mathbf{r}, \mathbf{n}, I) = \prod_{i=1}^k \frac{(n_i + 1)!}{(n_i - r_i)! r_i!} p_i^{r_i} (1 - p_i)^{n_i - r_i}, \quad (8.3)$$

where  $\mathbf{p} = \{p_1, p_2, \dots, p_k\}$ , with  $\mathbf{n}$  and  $\mathbf{r}$  similarly defined, and the posterior is understood to be zero if any the  $p_i$ 's fall out of the interval  $[0, 1]$ .

When we make this assumption, we often report a value of  $p$  that is given by the mean of the  $p_i$ 's with some error bar.

### 8.1.3 Best of both worlds: a hierarchical model

Each of these extremes have their advantages. We are often trying to estimate a parameter that is more universal than our experiments, e.g., something that describes worms with Channelrhodopsin in the ASH neuron generally. So, pooling the experiments makes sense. On the other hand, we have reason to assume that there is going to be a different value of  $p$  in different experiments, as biological systems are highly variable, not to mention measurement variations. So, how can we capture both of these effects?

We can consider a model in which there is a “master” reversal probability, which we’ll call  $q$  to avoid too many  $p$ 's, and the values of  $p_i$  may vary from this  $q$  according to some probability distribution,  $P(p_i \mid q, I)$ . So now, we have parameters  $p_1, p_2, \dots, p_k$  and  $q$ . So, the posterior can be written using Bayes’s theorem,

$$P(q, \mathbf{p} \mid \mathbf{r}, \mathbf{n}, I) = \frac{P(\mathbf{r}, \mathbf{n} \mid q, \mathbf{p}, I) P(q, \mathbf{p} \mid I)}{P(\mathbf{n}, \mathbf{r} \mid I)}. \quad (8.4)$$

Note, though, that the observed values of  $r$  do not depend directly on  $q$ , only on  $\mathbf{p}$ . In other words, they only depend indirectly on  $q$ . So, we can write  $P(\mathbf{r}, \mathbf{n} \mid q, \mathbf{p}, I) = P(\mathbf{r}, \mathbf{n} \mid \mathbf{p}, I)$ . Thus, we have

$$P(q, \mathbf{p} \mid \mathbf{r}, \mathbf{n}, I) = \frac{P(\mathbf{r}, \mathbf{n} \mid \mathbf{p}, I) P(q, \mathbf{p} \mid I)}{P(\mathbf{n}, \mathbf{r} \mid I)}. \quad (8.5)$$

Next, we can rewrite the prior using the definition of conditional probability.

$$P(q, \mathbf{p} \mid I) = P(\mathbf{p} \mid q, I) P(q \mid I). \quad (8.6)$$

Substituting this back into our expression for the posterior, we have

$$P(q, \mathbf{p} \mid \mathbf{r}, \mathbf{n}, I) = \frac{P(\mathbf{r}, \mathbf{n} \mid \mathbf{p}, I) P(\mathbf{p} \mid q, I) P(q \mid I)}{P(\mathbf{n}, \mathbf{r} \mid I)}. \quad (8.7)$$

Now, if we read off the numerator of this equation, we see a chain of dependencies. The experimental results  $\mathbf{r}$  depend on parameters  $\mathbf{p}$ . Parameters  $\mathbf{p}$  depend on *hyperparameter*  $q$ . Hyperparameter  $q$  then has some prior distribution. Any model that can be written as a chain of dependencies like this is called a **hierarchical model**, and the parameters that do not *directly* influence the data are called **hyperparameters**.

So, the hierarchical model captures both the experiment-to-experiment variability, as well as the master regulator of outcomes. Note that the product  $P(\mathbf{p} \mid q, I) P(q \mid I)$  comprises the prior, as it is therefore independent of the observed data.

## 8.2 Exchangeability

The conditional probability,  $P(\mathbf{p} \mid q, I)$ , can take any reasonable form. In the case where we have no reason to believe that we can distinguish any one  $p_i$  from another prior to the experiment, then the label “ $i$ ” applied to the experiment may be exchanged with the label of any other experiment. I.e.,  $P(p_1, p_2, \dots, p_k \mid q, I)$  is invariant to permutations of the indices. Parameters behaving this way are said to be **exchangeable**. A common (simple) exchangeable distribution is

$$P(\mathbf{p} \mid q, I) = \prod_{i=1}^k P(p_i \mid q, I), \quad (8.8)$$

which means that each of the parameters is an independent sample out of a distribution  $P(p_i \mid q, I)$ , which we often take to be the same for all  $i$ . This is reasonable to do in the worm reversal example.

### 8.3 Choice of the conditional distribution/prior

We need to specify our prior, which for this hierarchical model means that we have to specify the conditional distribution,  $P(p_i | q, I)$ , as well as  $P(q | I)$ . For the latter, we will take it to be uniform on  $[0, 1]$ . For the conditional distribution, we will assume it is Beta-distributed, which is defined on the interval  $[0, 1]$  and can be peaked. The Beta distribution can be written as

$$P(p | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1}, \quad (8.9)$$

where it is parametrized by positive constants  $\alpha$  and  $\beta$ . The Beta distribution has mean and variance, respectively, of

$$q = \frac{\alpha}{\alpha + \beta}, \quad (8.10)$$

$$\sigma^2 = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}. \quad (8.11)$$

So, we could also parametrize the Beta distribution with these, noting that we can convert back to  $\alpha$  and  $\beta$  using

$$\alpha = \frac{q}{\sigma^2} (q(1-q) - \sigma^2), \quad (8.12)$$

$$\beta = \frac{1-q}{\sigma^2} (q(1-q) - \sigma^2). \quad (8.13)$$

We have  $0 < q < 1$  and  $0 < \sigma^2 < q(1-q)$ . So, we have an additional hyperparameter,  $\sigma^2$ , which describes experiment-to-experiment variability. In analogy to Gaussian distributions (which the Beta approximates for large  $\alpha$  and  $\beta$ ), we will take  $P(\sigma^2 | I) \propto 1/\sigma^2$ . Thus, our full posterior is

$$P(q, \sigma^2, \mathbf{p} | \mathbf{r}, \mathbf{n}, I) \propto P(\mathbf{r}, \mathbf{n} | \mathbf{p}, I) \sigma^{-2} \left( \prod_{i=1}^k P(p_i | q, \kappa) \right), \quad (8.14)$$

nonzero on  $0 \leq q, \mathbf{p} \leq 1$  and  $0 < \sigma^2 < q(1-q)$ , where

$$P(p_i | q, \sigma^2) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p_i^\alpha (1-p_i)^\beta, \quad (8.15)$$

where  $\alpha$  and  $\beta$  are understood to be functions of  $q$  and  $\sigma^2$  according to equations (8.12) and (8.13). As before, we have a binomial likelihood, where we assume the experiments are independent.

$$P(\mathbf{r}, \mathbf{n} | \mathbf{p}, I) = \prod_{i=1}^k \frac{n_i!}{(n_i - r_i)! r_i!} p_i^{r_i} (1-p_i)^{n_i - r_i}. \quad (8.16)$$

## 8.4 Implementation

In some cases, we can do some gnarly integration and work out analytical results for the posterior of a hierarchical model. This usually involves choosing conjugate priors. Most often, though, we need to resort to numerical methods, MCMC as usual being the most powerful. To see the worm reversal problem solved with a hierarchical model, see the implementation [here](#).

## 8.5 Generalization

The worm reversal problem is easily generalized. You can imagine having more levels of the hierarchy. This is just more steps in the chain of dependencies that are factored in the prior. For general parameters  $\boldsymbol{\theta}$  and hyperparameters  $\boldsymbol{\phi}$ , we have

$$P(\boldsymbol{\theta}, \boldsymbol{\phi} \mid D, I) = \frac{P(D \mid \boldsymbol{\theta}, I) P(\boldsymbol{\theta} \mid \boldsymbol{\phi}, I) P(\boldsymbol{\phi} \mid I)}{P(D \mid I)} \quad (8.17)$$

for a two-level hierarchical model. For a three-level hierarchical model, we can consider hyperparameters  $\boldsymbol{\xi}$  that depend on  $\boldsymbol{\phi}$ , giving

$$P(\boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\xi} \mid D, I) = \frac{P(D \mid \boldsymbol{\theta}, I) P(\boldsymbol{\theta} \mid \boldsymbol{\phi}, I) P(\boldsymbol{\phi} \mid \boldsymbol{\xi}, I) P(\boldsymbol{\xi} \mid I)}{P(D \mid I)}, \quad (8.18)$$

and so on for four, five, etc., level hierarchical models. As we have seen in the course, the work is all in coming up with the models for the likelihood  $P(D \mid \boldsymbol{\theta}, I)$ , and prior,  $P(\boldsymbol{\theta} \mid \boldsymbol{\phi}, I) P(\boldsymbol{\phi} \mid I)$ , in this case for a two-level hierarchical model. For coming up with the conditional portion of the prior,  $P(\boldsymbol{\theta} \mid \boldsymbol{\phi}, I)$ , we often assume a Gaussian distribution because this often describes experiment-to-experiment variability. (The Beta distribution we used in our example is approximately Gaussian and has the convenient feature that it is defined on the interval  $[0, 1]$ .) Bayes's theorem gives you the posterior, and it is then “just” a matter of computing it by sampling from it.