

# Coursework 4: STAT 570

Philip Pham

November 5, 2018

1. Consider the so-called Neyman-Scott problem in which  $Y_{ij} \mid \mu_i, \sigma^2 \sim_{\text{ind}} \mathcal{N}(\mu_i, \sigma^2)$ ,  $i = 1, \dots, n, j = 1, 2$ .
  - (a) Obtain the MLE of  $\sigma^2$  and show that it is inconsistent. Why does this inconsistency arise in this example?

**Solution:** The likelihood is

$$\begin{aligned} L(\mu, \sigma) &= \prod_{i=1}^n \prod_{j=1}^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (Y_{ij} - \mu_i)^2\right) \\ &= \prod_{i=1}^n \frac{1}{2\pi\sigma^2} \exp\left(-\frac{1}{2\sigma^2} [(Y_{i1} - \mu_i)^2 + (Y_{i2} - \mu_i)^2]\right), \end{aligned} \quad (1)$$

so the log-likelihood is

$$l(\mu, \sigma) = -n \log(2\pi) - n \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n [(Y_{i1} - \mu_i)^2 + (Y_{i2} - \mu_i)^2]. \quad (2)$$

Taking the derivative with respect to  $\sigma^2$ , we have

$$\frac{\partial}{\partial \sigma^2} l(\mu, \sigma^2) = -\frac{n}{\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n [(Y_{i1} - \mu_i)^2 + (Y_{i2} - \mu_i)^2]. \quad (3)$$

Solving Equation 3, where  $\frac{\partial}{\partial \sigma^2} l(\hat{\mu}, \hat{\sigma}^2) = 0$ , we have

$$\hat{\sigma}^2 = \frac{1}{2n} \sum_{i=1}^n [(Y_{i1} - \hat{\mu}_i)^2 + (Y_{i2} - \hat{\mu}_i)^2]. \quad (4)$$

Taking the derivative of Equation 2 with respect to  $\mu_i$ , we have

$$\frac{\partial}{\partial \mu_i} l(\mu, \sigma^2) = \frac{1}{\sigma^2} (Y_{i1} + Y_{i2} - 2\mu_i). \quad (5)$$

Solving Equation 5, where  $\frac{\partial}{\partial \mu_i} l(\hat{\mu}, \hat{\sigma}^2) = 0$ , we have

$$\hat{\mu}_i = \frac{Y_{i1} + Y_{i2}}{2}. \quad (6)$$

Substituting Equation 6 into Equation 4, we have

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_{i1} - Y_{i2}}{2} \right)^2. \quad (7)$$

Taking the expected value of Equation 7, we have

$$\begin{aligned}\mathbb{E}[\hat{\sigma}^2] &= \frac{1}{4n} \sum_{i=1}^n \left( \mathbb{E}[Y_{i1}^2] + \mathbb{E}[Y_{i2}^2] - 2\mathbb{E}[Y_{i1}Y_{i2}] \right) \\ &= \frac{1}{4n} \sum_{i=1}^n \left( (\sigma^2 + \mu_i^2) + (\sigma^2 + \mu_i^2) - 2\mu_i^2 \right) \\ &= \frac{\sigma^2}{2}.\end{aligned}\tag{8}$$

Clearly,  $\mathbb{E}[\hat{\sigma}^2] = \sigma^2/2 \not\rightarrow \sigma^2$ , so the estimator is not consistent.

This is because the MLE estimate of  $\sigma^2$  depends on  $\mu_1, \dots, \mu_n$ , so the number of parameters being estimated increases with  $n$ . Thus, the model is not well-defined.

(b) Derive the posterior distribution corresponding to the prior

$$\pi(\mu_1, \dots, \mu_n, \sigma^2) \propto \sigma^{-n-2}\tag{9}$$

and show that

$$\mathbb{E}[\sigma^2 | Y] = \frac{1}{2(n-1)} \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2}.\tag{10}$$

**Solution:** Using the likelihood in Equation 1 and the prior in Equation 9. We have that

$$p(\mu, \sigma^2 | Y) \propto L(\mu, \sigma^2) \pi(\mu_1, \dots, \mu_n, \sigma^2).\tag{11}$$

We have that

$$\begin{aligned}p(Y) &= \int_0^\infty \int_{-\infty}^\infty \cdots \int_{-\infty}^\infty L(\mu, \sigma^2) \pi(\mu_1, \dots, \mu_n, \sigma^2) d\mu_1 \cdots d\mu_n d\sigma^2 \\ &= \int_0^\infty \frac{1}{2^n \pi^n (\sigma^2)^{(3n+2)/2}} (\pi \sigma^2)^{n/2} \prod_{i=1}^n \exp\left(-\frac{1}{4\sigma^2} (Y_{i1} - Y_{i2})^2\right) d\sigma^2 \\ &= \int_0^\infty \frac{1}{2^n \pi^{n/2} (\sigma^2)^{n+1}} \exp\left(-\frac{1}{4\sigma^2} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2\right) d\sigma^2 \\ &= -\frac{2^n}{\pi^{n/2}} \left( \sum_{i=1}^n (Y_{i1} - Y_{i2})^2 \right)^{-n} \int_0^\infty u^{n-1} \exp(-u) du \\ &= \frac{1}{\pi^{n/2}} \left( \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \right)^{-n} \Gamma(n).\end{aligned}\tag{12}$$

Normalizing Equation 11 with the evidence Equation 12, we have the posterior

$$p(\mu, \sigma^2 | Y) = \frac{(\sigma^2)^{-(3n+2)/2}}{2^n \pi^{n/2} \Gamma(n)} \left( \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \right)^{-n} \prod_{i=1}^n \prod_{j=1}^2 \exp\left(-\frac{1}{2\sigma^2} (Y_{ij} - \mu_i)^2\right).\tag{13}$$

Marginalizing  $\mu$  in Equation 13, we get

$$\begin{aligned}p(\sigma^2 | Y) &= \int_{-\infty}^\infty \cdots \int_{-\infty}^\infty p(\mu, \sigma^2 | Y) d\mu_1 \cdots d\mu_n \\ &= \frac{(\sigma^2)^{-n-1}}{2^n \Gamma(n)} \left( \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \right)^{-n} \exp\left(-\frac{1}{4\sigma^2} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2\right).\end{aligned}\tag{14}$$

Taking the expectation over the distribution in Equation 14, we have that

$$\begin{aligned}
\mathbb{E}[\sigma^2 | Y] &= \int_0^\infty \sigma^2 p(\sigma^2 | Y) d\sigma^2 \\
&= \frac{1}{\Gamma(n)} \int_0^\infty \left( \frac{1}{4\sigma^2} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2 \right)^n \exp\left(-\frac{1}{4\sigma^2} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2\right) d\sigma^2 \\
&= \frac{\sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{4\Gamma(n)} \int_0^\infty u^{n-1-1} \exp(u) du \\
&= \frac{\Gamma(n-1)}{2\Gamma(n)} \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \\
&= \frac{1}{2(n-1)} \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2}, \tag{15}
\end{aligned}$$

which is the desired result.

- (c) Hence, using Equation 15, show that  $\mathbb{E}[\sigma^2 | Y] \rightarrow \sigma^2/2$  as  $n \rightarrow \infty$ , so that the posterior mean is inconsistent.

**Solution:** From Equation 15, we have that

$$\begin{aligned}
\lim_{n \rightarrow \infty} \mathbb{E}[\sigma^2 | Y] &= \lim_{n \rightarrow \infty} \frac{1}{2(n-1)} \sum_{i=1}^n \frac{\mathbb{E}[(Y_{i1} - Y_{i2})^2]}{2} \\
&= \lim_{n \rightarrow \infty} \frac{1}{2(n-1)} \sum_{i=1}^n \frac{\text{Var}(Y_{i1} - Y_{i2})}{2} \\
&= \lim_{n \rightarrow \infty} \frac{n\sigma^2}{2(n-1)} \\
&= \frac{\sigma^2}{2} \neq \sigma^2, \tag{16}
\end{aligned}$$

so the posterior mean is inconsistent.

- (d) Examine the posterior distribution corresponding to the prior

$$\pi(\mu_1, \dots, \mu_n \sigma^2) \propto \sigma^{-2}. \tag{17}$$

**Solution:** If we use Equation 17, Equation 12 becomes

$$\begin{aligned}
p(Y) &= \int_0^\infty \frac{1}{2^n \pi^{n/2} (\sigma^2)^{n/2+1}} \exp\left(-\frac{1}{4\sigma^2} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2\right) d\sigma^2 \\
&= \frac{\Gamma(\frac{n}{2})}{\pi^{n/2}} \left( \sum_{i=1}^n (Y_{i1} - Y_{i2})^2 \right)^{-n/2}. \tag{18}
\end{aligned}$$

With Equation 18, the posterior becomes

$$p(\mu, \sigma^2 | Y) = \frac{(\sigma^2)^{-n-1}}{2^n \pi^{n/2} \Gamma(n/2)} \left( \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \right)^n \prod_{i=1}^n \prod_{j=1}^2 \exp\left(-\frac{1}{2\sigma^2} (Y_{ij} - \mu_i)^2\right). \tag{19}$$

Marginalizing Equation 19 over  $\mu$ , we have

$$p(\sigma^2 | Y) = \frac{1}{\sigma^2 \Gamma(n/2)} \left( \frac{\sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{4\sigma^2} \right)^{n/2} \exp\left(-\frac{\sum_{i=1}^n (Y_{i1} - Y_{i2})^2}{4\sigma^2}\right). \tag{20}$$

Length (mm)	0	1	2	3	4	5	6	7	8	9	10	11	12
1	2.247	2.640	2.842	2.908	3.099	3.126	3.245	3.328	3.355	3.383	3.572	3.581	3.681
10	1.901	2.132	2.203	2.228	2.257	2.350	2.361	2.396	2.397	2.445	2.454	2.454	2.474
20	1.312	1.314	1.479	1.552	1.700	1.803	1.861	1.865	1.944	1.958	1.966	1.997	2.006
50	1.339	1.434	1.549	1.574	1.589	1.613	1.746	1.753	1.764	1.807	1.812	1.840	1.852

Table 1: Failure stress data for four groups of fibers.

Equation 20 is quite similar to Equation 14, but with  $n$  replaced by  $n/2$  in the gamma function and the exponent of the sum of squares.

(e) Is the posterior mean for  $\sigma^2$  consistent in this case?

**Solution:** Yes. Taking the expectation with Equation 20, we have

$$\begin{aligned}
\mathbb{E}[\sigma^2 | Y] &= \int_0^\infty p(\sigma^2 | Y) d\sigma^2 \\
&= \frac{1}{4\Gamma(n/2)} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2 \int_0^\infty u^{n/2-1-1} \exp(-u) du \\
&= \frac{\Gamma(n/2-1)}{4\Gamma(n/2)} \sum_{i=1}^n (Y_{i1} - Y_{i2})^2 \\
&= \frac{1}{2(n/2-1)} \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2} \\
&= \frac{1}{n-2} \sum_{i=1}^n \frac{(Y_{i1} - Y_{i2})^2}{2}.
\end{aligned} \tag{21}$$

Taking the limit of Equation 21, we have

$$\begin{aligned}
\lim_{n \rightarrow \infty} \mathbb{E}[\sigma^2 | Y] &= \lim_{n \rightarrow \infty} \frac{1}{n-2} \sum_{i=1}^n \frac{\mathbb{E}[(Y_{i1} - Y_{i2})^2]}{2} \\
&= \lim_{n \rightarrow \infty} \frac{n}{(n-2)} \frac{1}{n} \sum_{i=1}^n \frac{\text{Var}(Y_{i1} - Y_{i2})}{2} \\
&= \lim_{n \rightarrow \infty} \frac{n}{(n-2)} \frac{1}{n} \sum_{i=1}^n \frac{2\sigma^2}{2} \\
&= \lim_{n \rightarrow \infty} \frac{n}{(n-2)} \sigma^2 \\
&= \sigma^2,
\end{aligned} \tag{22}$$

so the posterior mean is consistent when the prior doesn't depend on  $n$ .

2. The data in Table 1 contain data on a typical reliability experiment and give the failure stresses (in GPa) of four samples of carbon fibers of lengths 1, 10, 20 and 50mm.

(a) Consider a Bayesian analysis with an exponential likelihood and a gamma prior,  $\lambda \sim \text{Gamma}(a, b)$ . Derive the form of the posterior distribution for  $\lambda$ .

**Solution:** Suppose that we observe independent and identically distributed  $Y_i \sim \text{Exponential}(\lambda)$  for  $i = 1, \dots, n$ . Let  $Y = (Y_1, \dots, Y_n)$ . Then, we have

Length (mm)	$a'$	$b'$	Mean	Standard error
1	16.634063	48.275125	0.344568	0.084484
10	16.634063	37.320125	0.445713	0.109284
20	16.634063	30.025125	0.554005	0.135836
50	16.634063	28.940125	0.574775	0.140928

Table 2: The results of updating the prior belief in Part 2b with the data.

the likelihood function

$$L(\lambda) = p(Y | \lambda) = \lambda^n \exp \left( -\lambda \sum_{i=1}^n Y_i \right). \quad (23)$$

From Equation 23, we have the posterior

$$\begin{aligned} p(\lambda | Y) &\propto p(Y | \lambda) p(\lambda) \\ &\propto \left( \lambda^n \exp \left( -\lambda \sum_{i=1}^n Y_i \right) \right) \left( \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda) \right) \\ &\propto \lambda^{a+n-1} \exp \left( -\left( b + \sum_{i=1}^n Y_i \right) \lambda \right), \end{aligned} \quad (24)$$

which equal to the Gamma probability density function up to a constant factor independent of  $\lambda$ . Thus, we have that  $\lambda | Y \sim \Gamma(a + n, b + \sum_{i=1}^n Y_i)$ , and

$$p(\lambda | Y) = \frac{(b + \sum_{i=1}^n Y_i)^{a+n}}{\Gamma(a + n)} \lambda^{a+n-1} \exp \left( -\left( b + \sum_{i=1}^n Y_i \right) \lambda \right). \quad (25)$$

- (b) Choose  $a$  and  $b$  so that the prior probability that  $\lambda$  lies between 0.05 and 1 is 0.95.

**Solution:** For a specified mean  $\mu$ , we specify our prior as the gamma distribution  $\text{Gamma}(k\mu, k)$  for some  $k$ . Let  $F$  be the cumulative distribution function. We choose  $k$  such that

$$F(1) - F(0.05) = 0.95. \quad (26)$$

Choosing  $a \approx 3.634$  and  $b \approx 7.268$  satisfies Equation 26. These values were obtained numerically in `failure_stresses.ipynb`.

- (c) Obtain the posterior means and posterior standard deviations for  $\lambda$  and give histogram representations of  $p(\lambda | y)$  for each of the groups in Table 1. Compare inference with the frequentist analyses. Also give histogram representations of the posterior for  $\lambda^{-1}$ , again for each of the four groups.

**Solution:** The posterior parameters, means, and standard deviations can be seen in Table 2. The first two columns are the shape and rate parameters of the gamma posterior, respectively.

The means are similar to the estimates from frequentist analyses in the previous homework but are drawn towards the prior mean of 0.5. The

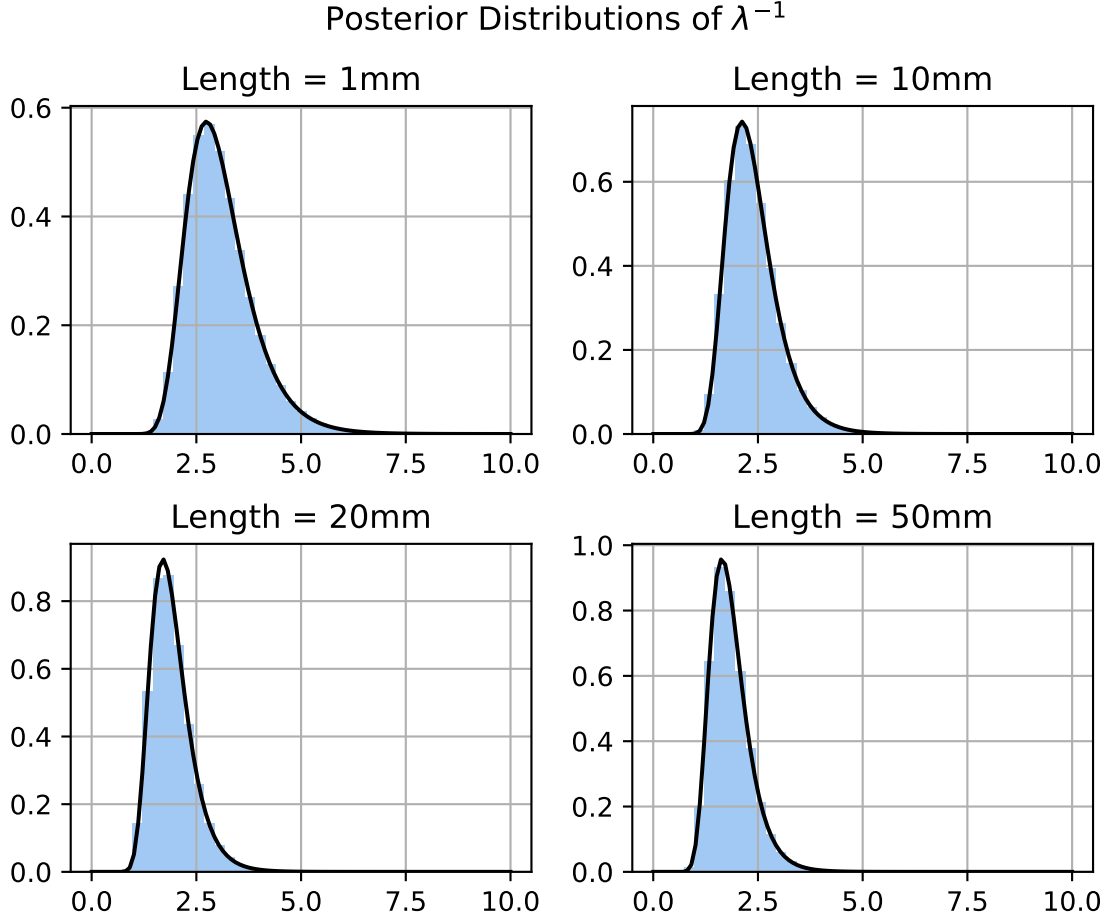


Figure 1: Histograms and probability density for samples drawn from the posteriors for  $\lambda^{-1}$ .

standard errors are slightly smaller than those obtained from fitting the exponential model but are much larger than either the quasi-likelihood model or sandwich estimates.

For the posterior distribution of  $\lambda^{-1}$ , we have that

$$\begin{aligned}
 p(\lambda^{-1} | Y) &= p(\lambda | Y) \left| \frac{d}{d\lambda^{-1}} (\lambda^{-1})^{-1} \right| = p(\lambda | Y) \lambda^2 \\
 &= \frac{(b + \sum_{i=1}^n Y_i)^{a+n}}{\Gamma(a+n)} \lambda^{a+n+1} \exp\left(-\left(b + \sum_{i=1}^n Y_i\right) \lambda\right), \quad (27)
 \end{aligned}$$

where we used Equation 25 and transformation of random variables, so  $\lambda^{-1} | Y \sim \text{InverseGamma}(a+n, b + \sum_{i=1}^n Y_i)$ .

Theoretical histograms along with their density are plotted in Figure 1.

Unlike the MLE, the posterior mean is not invariant under reparameterization:

$$\mathbb{E}[\lambda^{-1} | Y] = \frac{b'}{a' - 1} \geq \frac{b'}{a'} = \frac{1}{\mathbb{E}[\lambda | Y]}.$$

See `failure_stresses.ipynb` for code to reproduce plots.

3. In this question we will consider inference when the sampling model is multivariate hypergeometric. Suppose a population contains objects of  $K$  different types, with  $X_1, \dots, X_K$  being the number of each type,  $\sum_{k=1}^K X_k = N$ . A simple random sample of size  $n$  is taken and the number of each type,  $Y_1, \dots, Y_K$ , is recorded (so that  $\sum_{k=1}^K y_k = n$ ).

An obvious model for  $Y_1, \dots, Y_K$ , is the multivariate hypergeometric distribution:

$$\mathbb{P}(Y_1 = y_1, \dots, Y_K = y_K \mid x_1, \dots, x_K) = \frac{\prod_{k=1}^K \binom{x_k}{y_k}}{\binom{N}{n}}, \quad (28)$$

with means and variances:

$$\mathbb{E}[Y_k \mid x_k] = n \frac{x_k}{N} \quad (29)$$

$$\text{Var}(Y_k \mid x_k) = n \frac{x_k}{N} \left(1 - \frac{x_k}{N}\right) \frac{N-n}{N-1}. \quad (30)$$

Suppose we take a sample from a population of  $K$  distinct objects, and record  $y_1, \dots, y_K$ , but the numbers  $X_1, \dots, X_K$  are unknown (but  $N$  is known).

- (a) Using Equation 29, write down an estimator for  $X_k$ ,  $k = 1, \dots, K$ . We will refer to this a method of moments estimator. Using Equation 30, give a form for the variance of this estimator, along with the estimator of this variance.

**Solution:** A simple estimator for  $X_k$  can be obtained from Equation 29 by substituting  $\mathbb{E}[Y_k \mid x_k]$  with the observed  $y_k$ :

$$y_k = n \frac{\hat{x}_k}{N} \implies \hat{x}_k = \frac{N}{n} y_k. \quad (31)$$

If we treat  $\hat{x}_k$  like a random variable, then the variance of our estimator is

$$\begin{aligned} \text{Var}(\hat{x}_k) &= \text{Var}\left(\frac{N}{n} Y_k\right) = \left(\frac{N}{n}\right)^2 \text{Var}(Y_k) \\ &= \left(\frac{N}{n}\right)^2 \left(n \frac{x_k}{N} \left(1 - \frac{x_k}{N}\right) \frac{N-n}{N-1}\right) \\ &= x_k \left(\frac{N}{n}\right) \left(1 - \frac{x_k}{N}\right) \frac{N-n}{N-1} \end{aligned} \quad (32)$$

Since we don't know  $x_k$ , we substitute  $x_k$  with  $\hat{x}_k$  to get the estimated variance:

$$\begin{aligned} \hat{\text{Var}}(\hat{x}_k) &= \hat{x}_k \left(\frac{N}{n}\right) \left(1 - \frac{\hat{x}_k}{N}\right) \frac{N-n}{N-1} \\ &= y_k \left(\frac{N}{n}\right)^2 \left(1 - \frac{y_k}{n}\right) \frac{N-n}{N-1} \\ &= (N-n) \frac{y_k}{n} \left(1 - \frac{y_k}{n}\right) \left(\frac{N^2}{n(N-1)}\right) \end{aligned} \quad (33)$$

- (b) We now consider a Bayesian approach to inference. Consider a multinomial distribution for counts  $X_1, \dots, X_K$ ,

$$\mathbb{P}(X_1 = x_1, \dots, X_K = x_K \mid p_1, \dots, p_K) = \frac{N!}{\prod_{k=1}^K x_k!} \prod_{k=1}^K p_k^{x_k}, \quad (34)$$

with  $p_k > 0$  and  $\sum_{k=1}^K p_k = 1$ . Show that the Dirichlet:

$$\pi(p_1, \dots, p_K) = \frac{\Gamma(\alpha_+)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k-1}, \quad (35)$$

where  $\alpha_k > 0$ ,  $k = 1, \dots, K$  and  $\alpha_+ = \sum_{k=1}^K \alpha_k$  is the conjugate distribution to the multinomial sampling model.

**Solution:** We use the definition of conjugate distributions and show that  $(p_1, \dots, p_K) \mid (X_1 = x_1, \dots, X_K = x_K)$  is of the same distribution family as  $(p_1, \dots, p_K)$ . We have that

$$\begin{aligned} p(p_1, \dots, p_K \mid x_1, \dots, x_K) &\propto \mathbb{P}(x_1, \dots, x_K \mid p_1, \dots, p_K) \pi(p_1, \dots, p_K) \\ &\propto \prod_{k=1}^K p_k^{x_k + \alpha_k - 1}, \end{aligned} \quad (36)$$

which we recognize as the form of the Dirichlet family. Thus, we have that

$$p(p_1, \dots, p_K \mid x_1, \dots, x_K) = \frac{\Gamma(\alpha_+ + N)}{\prod_{k=1}^K \Gamma(\alpha_k + x_k)} \prod_{k=1}^K p_k^{\alpha_k + x_k - 1}, \quad (37)$$

so the Dirichlet distribution is the conjugate distribution to the multinomial sampling model.

(c) The compound multinomial distribution,  $\text{CMult}(N, \alpha)$ , is defined as

$$\mathbb{P}(X_1 = x_1, \dots, X_K = x_K) = \frac{N! \Gamma(\alpha_+)}{\Gamma(N + \alpha_+)} \prod_{k=1}^K \frac{\Gamma(x_k + \alpha_k)}{x_k! \Gamma(\alpha_k)} \quad (38)$$

where  $\alpha = (\alpha_1, \dots, \alpha_K)$ . Show that the prior predictive distribution, obtained as the marginal distribution when the likelihood is Equation 34 and the prior is Equation 35, is of the compound multinomial form with parameters that you should identify.

**Solution:** Let  $\mathbf{p} = (p_1, \dots, p_K)$ . We simply integrate:

$$\begin{aligned} \mathbb{P}(X_1 = x_1, \dots, X_K = x_K) &= \int \mathbb{P}(x_1, \dots, x_K \mid p_1, \dots, p_K) \pi(p_1, \dots, p_K) \, d\mathbf{p} \\ &= \frac{N!}{\prod_{k=1}^K x_k!} \frac{\Gamma(\alpha_+)}{\prod_{k=1}^K \Gamma(\alpha_k)} \int \prod_{k=1}^K p_k^{x_k + \alpha_k - 1} \, d\mathbf{p} \\ &= \frac{N! \Gamma(\alpha_+)}{\prod_{k=1}^K x_k! \Gamma(\alpha_k)} \left( \frac{\prod_{k=1}^K \Gamma(x_k + \alpha_k)}{\Gamma(N + \alpha_+)} \right) \\ &= \frac{N! \Gamma(\alpha_+)}{\Gamma(N + \alpha_+)} \prod_{k=1}^K \frac{\Gamma(x_k + \alpha_k)}{x_k! \Gamma(\alpha_k)}, \end{aligned}$$

where we have used that the Dirichlet probability density function must integrate to 1 to compute the integral.



- (d) Find the mean  $\mathbb{E}[X_k]$  and variance  $\text{Var}(X_k)$ ,  $k = 1, \dots, K$ , of a compound multinomial distribution.

**Solution:** We can use the law of total expectation for the mean:

$$\mathbb{E}[X_k] = \mathbb{E}_\pi[\mathbb{E}[X_k | p_k]] = \mathbb{E}_\pi[Np_k] = N \frac{\alpha_k}{\alpha_+}, \quad (39)$$

where we have taken advantage of the known means of the multinomial and Dirichlet distributions.

For the variance we can use the law of total variance:

$$\begin{aligned} \text{Var}(X_k) &= \mathbb{E}_\pi[\text{Var}(X_k | p_k)] + \text{Var}(\mathbb{E}[X_k | p_k]) \\ &= \mathbb{E}_\pi[Np_k(1 - p_k)] + \text{Var}(Np_k) \\ &= N \left( \mathbb{E}_\pi[p_k] - \mathbb{E}_\pi[p_k^2] \right) + N^2 \text{Var}(p_k) \\ &= N \left( \frac{\alpha_k}{\alpha_+} - \frac{\alpha_k(\alpha_+ - \alpha_k)}{\alpha_+^2(\alpha_+ + 1)} - \left( \frac{\alpha_k}{\alpha_+} \right)^2 \right) + N^2 \frac{\alpha_k(\alpha_+ - \alpha_k)}{\alpha_+^2(\alpha_+ + 1)} \\ &= N \frac{\alpha_k}{\alpha_+} \left( 1 - \frac{1 + \alpha_k}{1 + \alpha_+} \right) + N \frac{\alpha_k}{\alpha_+} \left( N \frac{\alpha_+ - \alpha_k}{\alpha_+(\alpha_+ + 1)} \right) \\ &= N \frac{\alpha_k}{\alpha_+} \left( \frac{\alpha_+(\alpha_+ - \alpha_k)}{\alpha_+(1 + \alpha_+)} + \frac{N(\alpha_+ - \alpha_k)}{\alpha_+(1 + \alpha_+)} \right) \\ &= N \frac{\alpha_k}{\alpha_+} \left( \frac{(\alpha_+ - \alpha_k)(N + \alpha_+)}{\alpha_+(1 + \alpha_+)} \right) \\ &= N \frac{\alpha_k}{\alpha_+} \left( 1 - \frac{\alpha_k}{\alpha_+} \right) \left( \frac{N + \alpha_+}{1 + \alpha_+} \right). \end{aligned} \quad (40)$$

- (e) Let  $W_k = X_k - y_k$  represent the unobserved counts,  $k = 1, \dots, K$ . Show that the posterior distribution  $\mathbb{P}(W_1, \dots, W_K | y_1, \dots, y_K)$  is compound multinomial  $\text{CMult}(N - n, \alpha + y)$ , where  $y = (y_1, \dots, y_K)$ .

**Solution:** Given  $y$ ,  $W = (W_1, \dots, W_K)$  is just a reparameterization of  $X = (X_1, \dots, X_K)$ . However, the domain is different since  $W_k \geq 0$ , so we need a normalization constant.

Using Equation 38 as our prior and Equation 28 as our likelihood, we can compute the posterior

$$\begin{aligned} &\mathbb{P}(W_1 = w_1, \dots, W_K = w_K | y_1, \dots, y_K) \\ &\propto \mathbb{P}(y_1, \dots, y_K | X) \mathbb{P}(X_1 = w_1 + y_1, \dots, X_K = w_K + y_K) \\ &\propto \left( \frac{\prod_{k=1}^K \binom{w_k + y_k}{y_k}}{\binom{N}{n}} \right) \left( \frac{N! \Gamma(\alpha_+)}{\Gamma(N + \alpha_+)} \prod_{k=1}^K \frac{\Gamma(w_k + y_k + \alpha_k)}{(w_k + y_k)! \Gamma(\alpha_k)} \right) \\ &\propto \frac{(N - n)! n! \Gamma(\alpha_+)}{\Gamma(N + \alpha_+)} \prod_{k=1}^K \frac{\Gamma(w_k + y_k + \alpha_k)}{w_k! y_k! \Gamma(\alpha_k)} \\ &\propto \prod_{k=1}^K \frac{\Gamma(w_k + y_k + \alpha_k)}{w_k!}, \end{aligned} \quad (41)$$

nn which we recognize as the from the compound multionomial distribution. Since Equation 41 must sum to 1 to be a proper probability distribution

over the support of the prior  $\sum_{k=1}^K W_k = N - n$ , where  $W_k$  are nonnegative integers, we get that

$$\begin{aligned} \mathbb{P}(W_1 = w_1, \dots, W_K = w_K \mid y_1, \dots, y_K) \\ = \frac{(N - n)! \Gamma(n + \alpha_+)}{\Gamma((N - n) + (n + \alpha_+))} \prod_{k=1}^K \frac{\Gamma(w_k + (y_k + \alpha_k))}{w_k! \Gamma(y_k + \alpha_k)}, \end{aligned} \quad (42)$$

which is the CMult( $N - n, y + \alpha$ ) distribution.

- (f) Write down the posterior mean and posterior variance of  $X_k$ ,  $k = 1, \dots, K$ . Comment on the case when  $\alpha_k = 0$ ,  $k = 1, \dots, K$ .

**Solution:** Using the results from Equation 42, 39, and 40, we have

$$\begin{aligned} \mathbb{E}[X_k \mid y_1, \dots, y_K] &= \mathbb{E}[W_k \mid y_1, \dots, y_K] + y_k \\ &= y_k + (N - n) \frac{y_k + \alpha_k}{n + \alpha_+}. \end{aligned} \quad (43)$$

$$\begin{aligned} \text{Var}(X_k \mid y_1, \dots, y_K) &= \text{Var}(W_k \mid y_1, \dots, y_K) \\ &= (N - n) \frac{y_k + \alpha_k}{n + \alpha_+} \left(1 - \frac{y_k + \alpha_k}{n + \alpha_+}\right) \left(\frac{N + \alpha_+}{1 + n + \alpha_+}\right). \end{aligned} \quad (44)$$

When  $\alpha = \mathbf{0}$ , the expectation and variance are similar to those of the multinomial distribution with  $p_k = y_k/n$ , but the variance has the additional  $\frac{N}{1+n}$  factor. These are similar to the method of the moments estimators in Equations 31 and 33. The mean estimate is same as posterior mean, and the variance has a  $\frac{N^2}{n(N-1)}$  has is quite close to  $\frac{N}{1+n}$ , especially for large values of  $N$  and  $n$ .

- (g) A certain infectious disease can be caused by one of three different pathogens, A, B, or C. Over a 1 year period population surveillance is carried out, and 750 individuals are observed to be infected. A random sample of 60 cases is selected for labtesting, i.e., to determine the pathogen responsible. Of these 60 selected cases, the numbers who were infected by pathogens A, B, C, were 42, 18, 0, respectively. We wish to estimate the numbers of the total population of cases that were infected by each of the pathogens.

- i. Calculate the method of moments estimators and the associated standard errors.

**Solution:** Let  $X_1$ ,  $X_2$ , and  $X_3$  correspond to the total population with pathogens A, B, and C, respectively. Let  $Y_1$ ,  $Y_2$ , and  $Y_3$  be the corresponding samples.

We substitute  $N = 750$ ,  $n = 60$ ,  $y_1 = 42$ ,  $y_2 = 18$ , and  $y_3 = 0$  into Equations 31 and 32.

For the estimated means, we have

$$\begin{aligned} \hat{x}_1 &= 750 \frac{42}{60} = 525 \\ \hat{x}_2 &= 750 \frac{18}{60} = 225 \\ \hat{x}_3 &= 750 \frac{0}{60} = 0. \end{aligned} \quad (45)$$

For estimated standard errors, we have

$$\begin{aligned}\sqrt{\hat{\text{Var}}(\hat{x}_1)} &\approx 42.587 \\ \sqrt{\hat{\text{Var}}(\hat{x}_2)} &\approx 42.587 \\ \sqrt{\hat{\text{Var}}(\hat{x}_3)} &= 0.\end{aligned}\tag{46}$$

- ii. Calculate the Bayesian posterior mean and posterior standard deviation, with prior specification,  $\alpha_k = 1$ ,  $k = 1, \dots, K$ . Which estimates are the most reasonable?

**Solution:** Here, we substitute into Equations 43 and 44.

We have that

$$\begin{aligned}\mathbb{E}[X_1 | y] &= 42 + (690) \frac{42 + 1}{60 + 3} = 42 + \frac{9890}{21} \approx 512.952 \\ \mathbb{E}[X_2 | y] &= 18 + (690) \frac{18 + 1}{60 + 3} = 18 + \frac{4370}{21} \approx 226.095 \\ \mathbb{E}[X_3 | y] &= 0 + (690) \frac{1}{63} = \frac{230}{21} \approx 10.952,\end{aligned}\tag{47}$$

and

$$\begin{aligned}\sqrt{\text{Var}[X_1 | y]} &\approx 41.941 \\ \sqrt{\text{Var}[X_2 | y]} &\approx 41.352 \\ \sqrt{\text{Var}[X_3 | y]} &\approx 11.261.\end{aligned}\tag{48}$$

The estimates in Equations 47 and 48 are similar to those in Equations 45 and 46 for  $X_1$  and  $X_2$  but are pulled slightly closer together since the prior specifies that the pathogens occur in about equal proportions. The estimate for  $X_3$  is much different, however. In Equations 45 and 46, the estimate and standard error is 0 because there were no observations. If we know that at least one person carries pathogen C, this estimate is deeply disturbing. It could simply be that no one in our sample carried the pathogen, so the Bayesian estimates are most reasonable.