

Answers to Exercises

(P1) (a) If $P(A) = 0$ or $P(A) = 1$.

(b) Take A as “odd red die”, B as “odd blue die” and C as “odd sum”.

(P2) (a) $P(\text{homozygous}) = 1/3$; $P(\text{heterozygous}) = 2/3$.

(b) By Bayes’ Theorem

$$P(BB \mid 7 \text{ black}) = \frac{(1/3)(1)}{(1/3)(1) + (2/3)(1/2^7)} = \frac{64}{65}.$$

(P3) $P(k=0) = (1-\pi)^n = (1-\lambda/n)^n \rightarrow e^{-\lambda}$. More generally

$$\begin{aligned} p(k) &= \binom{n}{k} \pi^k (1-\pi)^{n-k} \\ &= \frac{\lambda^k}{k!} \left(1 - \frac{\lambda}{n}\right)^n \frac{n}{n} \frac{n-1}{n} \dots \frac{n-k+1}{n} \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &\rightarrow \frac{\lambda^k}{k!} \exp(-\lambda). \end{aligned}$$

(P4) **proof**

When a joint probability density function is well defined and the expectations are integrable, we write for the general case

$$\mathbf{E}(X) = \int x \Pr[X = x] \, dx$$

$$\mathbf{E}(X \mid Y = y) = \int x \Pr[X = x \mid Y = y] \, dx$$

$$\begin{aligned} \mathbf{E}(\mathbf{E}(X \mid Y)) &= \int \left(\int x \Pr[X = x \mid Y = y] \, dx \right) \Pr[Y = y] \, dy \\ &= \int \int x \Pr[X = x, Y = y] \, dx \, dy \\ &= \int x \left(\int \Pr[X = x, Y = y] \, dy \right) \, dx \\ &= \int x \Pr[X = x] \, dx \\ &= \mathbf{E}(X) . \end{aligned}$$

(S1) $\bar{x} = 16.35525$, so assuming uniform prior, posterior is $N(16.35525, 1/12)$. A 90% HDR is $16.35525 \pm 1.6449/\sqrt{12}$ or 16.35525 ± 0.47484 , that is, the interval $(15.88041, 16.83009)$.

(S2) $x - \theta \sim N(0, 1)$ and $\theta \sim N(16.35525, 1/12)$, so $x \sim N(16.35525, 13/12)$.

(S3)

Continuous distribution, continuous parameter space

For the normal distribution $\mathcal{N}(\mu, \sigma^2)$ which has probability density function

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

the corresponding probability density function for a sample of n independent identically distributed normal random variables (the likelihood) is

$$f(x_1, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n f(x_i | \mu, \sigma^2) = \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right).$$

This family of distributions has two parameters: $\theta = (\mu, \sigma)$; so we maximize the likelihood, $\mathcal{L}(\mu, \sigma^2) = f(x_1, \dots, x_n | \mu, \sigma^2)$, over both parameters simultaneously, or if possible, individually.

Since the logarithm function itself is a continuous strictly increasing function over the range of the likelihood, the values which maximize the likelihood will also maximize its logarithm (the log-likelihood itself is not necessarily strictly increasing). The log-likelihood can be written as follows:

$$\log(\mathcal{L}(\mu, \sigma^2)) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

(Note: the log-likelihood is closely related to information entropy and Fisher information.)

We now compute the derivatives of this log-likelihood as follows.

$$0 = \frac{\partial}{\partial \mu} \log(\mathcal{L}(\mu, \sigma^2)) = 0 - \frac{-2n(\bar{x} - \mu)}{2\sigma^2}.$$

where \bar{x} is the sample mean. This is solved by

$$\hat{\mu} = \bar{x} = \sum_{i=1}^n \frac{x_i}{n}.$$

This is indeed the maximum of the function, since it is the only turning point in μ and the second derivative is strictly less than zero. Its expected value is equal to the parameter μ of the given distribution,

$$\mathbb{E}[\hat{\mu}] = \mu,$$

which means that the maximum likelihood estimator $\hat{\mu}$ is unbiased.

Similarly we differentiate the log-likelihood with respect to σ and equate to zero:

$$0 = \frac{\partial}{\partial \sigma} \log(\mathcal{L}(\mu, \sigma^2)) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2.$$

which is solved by

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Inserting the estimate $\mu = \hat{\mu}$ we obtain

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n x_i x_j.$$

To calculate its expected value, it is convenient to rewrite the expression in terms of zero-mean random variables (statistical error) $\delta_i \equiv \mu - x_i$. Expressing the estimate in these variables yields

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mu - \delta_i)^2 - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\mu - \delta_i)(\mu - \delta_j).$$

Simplifying the expression above, utilizing the facts that $\mathbb{E} [\delta_i] = 0$ and $\mathbb{E} [\delta_i^2] = \sigma^2$, allows us to obtain

$$\mathbb{E} [\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2.$$

This means that the estimator $\hat{\sigma}^2$ is biased for σ^2 . It can also be shown that $\hat{\sigma}$ is biased for σ , but that both $\hat{\sigma}^2$ and $\hat{\sigma}$ are consistent.

Formally we say that the *maximum likelihood estimator* for $\theta = (\mu, \sigma^2)$ is

$$\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2).$$

Example

Suppose that we are given a sequence (x_1, \dots, x_n) of IID $N(\mu, \sigma_v^2)$ random variables and a prior distribution of μ is given by $N(\mu_0, \sigma_m^2)$. We wish to find the MAP estimate of μ . Note that the normal distribution is its own conjugate prior, so we will be able to find a closed-form solution analytically.

The function to be maximized is then given by^[3]

$$g(\mu)f(x | \mu) = \pi(\mu)L(\mu) = \frac{1}{\sqrt{2\pi}\sigma_m} \exp\left(-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_m}\right)^2\right) \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_v} \exp\left(-\frac{1}{2}\left(\frac{x_j - \mu}{\sigma_v}\right)^2\right),$$

which is equivalent to minimizing the following function of μ :

$$\sum_{j=1}^n \left(\frac{x_j - \mu}{\sigma_v}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_m}\right)^2.$$

Thus, we see that the **MAP estimator** for μ is given by^[3]

$$\hat{\mu}_{\text{MAP}} = \frac{\sigma_m^2 n}{\sigma_m^2 n + \sigma_v^2} \left(\frac{1}{n} \sum_{j=1}^n x_j\right) + \frac{\sigma_v^2}{\sigma_m^2 n + \sigma_v^2} \mu_0 = \frac{\sigma_m^2 \left(\sum_{j=1}^n x_j\right) + \sigma_v^2 \mu_0}{\sigma_m^2 n + \sigma_v^2}.$$

which turns out to be a linear interpolation between the prior mean and the sample mean weighted by their respective covariances.

The case of $\sigma_m \rightarrow \infty$ is called a non-informative prior and leads to an improper probability distribution; in this case $\hat{\mu}_{\text{MAP}} \rightarrow \hat{\mu}_{\text{MLE}}$.

Solution(54)

Denoting $x' = x_{n+1}$ for short, the posterior predictive is

$$\begin{aligned}
 p(x'|x_{1:n}) &= \int p(x'|\theta)p(\theta|x_{1:n})d\theta \\
 &= \int_0^\infty \theta e^{-\theta x'} \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} d\theta \\
 &= \frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^\infty \theta^{(\alpha+1)-1} e^{-(\beta+x')\theta} d\theta \\
 &= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{\Gamma(\alpha+1)}{(\beta+x')^{\alpha+1}} \int \text{Gamma}(\theta | \alpha+1, \beta+x') d\theta \\
 &= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{\Gamma(\alpha+1)}{(\beta+x')^{\alpha+1}}.
 \end{aligned}$$

The marginal likelihood is

$$\begin{aligned}
 p(x_{1:n}) &= \int p(x_{1:n}|\theta)p(\theta)d\theta \\
 &= \int_0^\infty \theta^n e^{-\theta \sum x_i} \frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta} d\theta \\
 &= \frac{b^a}{\Gamma(a)} \int_0^\infty \theta^{a+n-1} \exp(- (b + \sum x_i)\theta) d\theta \\
 &= \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+n)}{(b + \sum x_i)^{a+n}} \int \text{Gamma}(\theta | a+n, b + \sum x_i) d\theta \\
 &= \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+n)}{(b + \sum x_i)^{a+n}} = \frac{b^a}{\Gamma(a)} \frac{\Gamma(a)}{\beta^a}.
 \end{aligned}$$

The marginal likelihood can also be found by using Bayes' theorem: for any θ ,

$$p(x_{1:n}) = \frac{p(x_{1:n}|\theta)p(\theta)}{p(\theta|x_{1:n})} = \frac{\frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta}}{\text{Gamma}(\theta|\alpha, \beta)} = \frac{\frac{b^a}{\Gamma(a)}}{\frac{\beta^a}{\Gamma(\alpha)}}.$$