

Exam STK 4021/9021 December 2015:

Notes, by Nils Lid Hjort

The Project, Exercise 1

- (a) The log-density of the binomial is $y \log p + (m - y) \log(1 - p)$, with derivative

$$u(y, p) = \frac{y}{p} - \frac{m - y}{1 - p} = \frac{y - mp}{p(1 - p)},$$

which has variance $m/\{p(1-p)\}$. The Jeffreys prior is hence proportional to $p^{-1/2}(1-p)^{-1/2}$, which is the $\text{Beta}(\frac{1}{2}, \frac{1}{2})$.

- (b) With Jeffreys priors for p_0 and p_1 , we have

$$p_0 \mid \text{data} \sim \text{Beta}(\frac{1}{2} + y_0, \frac{1}{2} + m_0 - y_0) \quad \text{and} \quad p_1 \mid \text{data} \sim \text{Beta}(\frac{1}{2} + y_1, \frac{1}{2} + m_1 - y_1),$$

by classic results; here $(y_0, m_0) = (22, 79)$ and $(y_1, m_1) = (19, 78)$. Also, the two parameters are independent. I simulate 10^5 realisations of (p_0, p_1) given data through

$$\text{p0sim} = \text{rbeta}(\text{sim}, 0.5 + y_0, 0.5 + m_0 - y_0)$$

$$\text{p1sim} = \text{rbeta}(\text{sim}, 0.5 + y_1, 0.5 + m_1 - y_1)$$

and read off histograms or density estimates for $\gamma = \log(p_1/p_0)$ and for $\rho = \exp(\gamma) = p_1/p_0$. With 10^6 simulations, the 0.025, 0.50, 0.975 quantile points are 0.512, 0.877, 1.480. Hence $\rho = 1$, corresponding to $p_0 = p_1$, is smack in the middle of the soup (the 95% credibility interval is [0.512, 1.480]). It appears artificial to claim that there is any noticeable difference.

- (c) Since $\log p$ has derivative $1/p$, we have

$$\sqrt{m}(\log \hat{p} - \log p) \rightarrow_d N(0, \tau^2)$$

with $\tau^2 = (1/p)^2 p(1-p) = 1/p - 1$. This translates to

$$\text{Var} \log \hat{p} \approx_d N(0, (1/m)(1/p - 1)).$$

The variance here, say κ_m^2 , is naturally estimated using $\hat{\kappa}_m^2 = 1/y - 1/m$. One can easily show that $\hat{\kappa}_m^2/\kappa_m^2 \rightarrow_{\text{pr}} 1$. (Note, incidentally, that $\log \hat{p}$ is not defined if $y = 0$, but this happens with exceedingly low probability.)

- (d) From the above we have

$$\log(\hat{p}_1/\hat{p}_0) \approx_d N(0, \sigma^2),$$

say, with

$$\sigma^2 = (1/m_1)(1/p_1 - 1) + (1/m_0)(1/p_0 - 1),$$

which we estimate using

$$\hat{\sigma}^2 = 1/y_1 - 1/m_1 + 1/y_0 - 1/m_0 = 0.269^2.$$

- (e) The classic frequentist approximate 95% interval for γ becomes $\hat{\gamma} \pm 1.96 \hat{\sigma}$, which is $[-0.6620, 0.3943]$. Transforming to $\rho = \exp(\gamma)$ scale, the interval is $[0.516, 1.483]$. It is quite close to the Bayesian credibility interval based on Jeffreys, above, as we expect from the general normal approximation theory for posterior distributions.
- (f) We learn from putting up prior times likelihood that it takes the form $\exp(-\frac{1}{2}Q)$, with a quadratic function Q in $(\gamma, \hat{\gamma})$, so the joint distribution is binormal. It has the mean and variance matrix structure as given in the exercise. Here we use e.g. that

$$\text{Var } \hat{\gamma} = \text{E Var } (\hat{\gamma} | \gamma) + \text{Var E } (\hat{\gamma} | \gamma) = \text{E } \sigma^2 + \text{Var } \gamma = \sigma_0^2 + \sigma^2.$$

The conditional distribution

$$\gamma | \hat{\gamma} \sim \text{N}(\hat{\gamma}_B, \sigma_B^2),$$

with formulae as in the exercise, follows from the appropriate Nils Exercises.

- (g) Suppose $\gamma \sim \text{N}(0, \sigma_0^2)$ with $\sigma_0 = 0.35$ (which is Laptook et al.'s 'neutral prior'). The central 95% prior credibility interval is $\gamma_0 \pm 1.96 \sigma_0 = \pm 0.686$, which translates to the interval $[0.504, 1.986]$ on the ρ scale. The prior density is log-normal,

$$f_0(\rho) = \phi\left(\frac{\log \rho - \gamma_0}{\sigma_0}\right) \frac{1}{\sigma_0} \frac{1}{\rho} \quad \text{for } \rho > 0,$$

with $\phi(\cdot)$ the standard normal density.

- (h) Combining this with $\hat{\gamma} | \gamma \approx \text{N}(\gamma, \hat{\sigma}^2)$, with $\hat{\sigma}$ well estimated as 0.269, leads by the above to

$$\gamma | \hat{\gamma} \approx_d \text{N}(-0.084, 0.214^2).$$

So the posterior for γ is approximately normal, and by easy transformation the posterior for ρ is approximately

$$f(\rho) = \phi\left(\frac{\log \rho - \hat{\gamma}_B}{\hat{\sigma}_B}\right) \frac{1}{\hat{\sigma}_B} \frac{1}{\rho} \quad \text{for } \rho > 0.$$

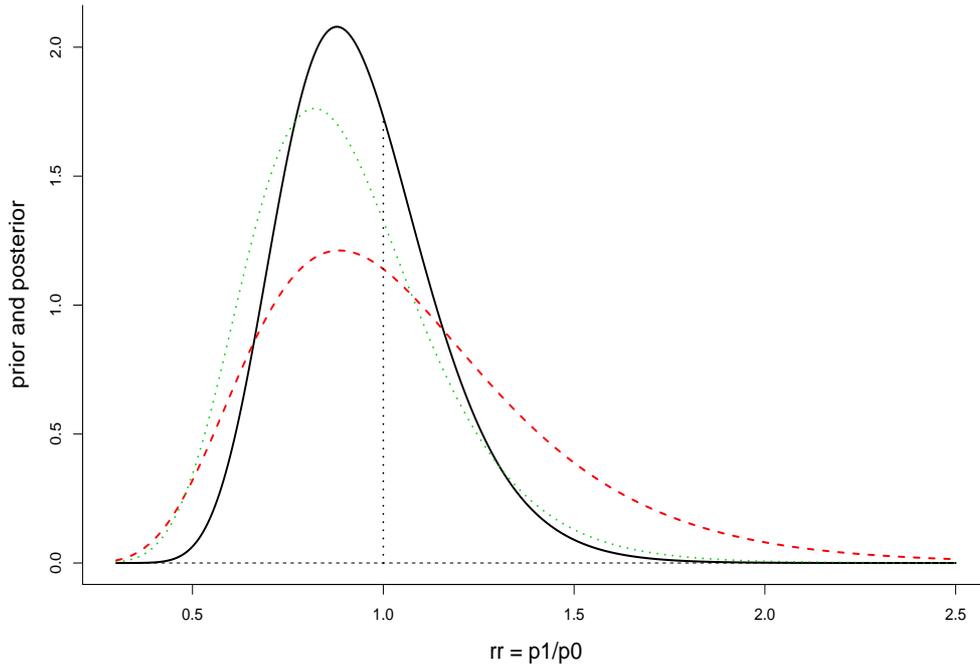
The posterior density for ρ with the Jeffreys priors, used above, can be well estimated via a histogram, or using a kernel density estimator, with say 10^6 simulations. A precise curve can also be computed numerically, however, bypassing approximations and simulations. Using g_0 and g_1 for the posterior Beta densities of p_0 and p_1 , with cumulatives G_0 and G_1 , the posterior cumulative for ρ is

$$H(\rho) = P\{p_1 \leq \rho p_0 \mid \text{data}\} = \int G_1(\rho p_0) g_0(p_0) dp_0$$

with derivative, i.e. the posterior density, equal to

$$h(\rho) = \int g_1(\rho p_0) p_0 f_0(p_0) dp_0.$$

This can be evaluated via numerical integration, using `dbeta` and indeed `integrate`. The figure displays the neutral prior for ρ (red, slanted); the posterior based on neutral prior and using the normal approximation above (black, full), which presumably what Laptok et al. have done; and the exact posterior density based on the Jeffreys priors (green, dotted). Again, the message is that $\rho = 1$ is smack in the middle, with no clear evidence that we can claim $\rho < 1$.



(i) I'm computing

$$Q(y_1) = P\{p_1 < p_0 \mid y_0 = 22, y_1\}$$

using the Jeffreys priors, and for each value $y_1 = 7, 8, 9, \dots, 18, 19$. This can be done via simulations, simulating for each value of y_1 a high number of (p_0, p_1) from the posterior distribution and then reading off the proportion of $p_1 < p_0$ cases. It may also be done via exact numerical integration, however:

$$Q(y_1) = \int G_1(p_0 \mid y_1) g_0(p_0 \mid y_0) dp_0,$$

with $g_0(p_0 \mid y_0)$ being the $\text{Beta}(0.5 + y_0, 0.5 + m_0 - y_0)$ density and $G_1(x \mid y_1)$ being the $\text{Beta}(0.5 + y_1, 0.5 + m_1 - y_1)$ cumulative. The result is a nice little plot $(y_1, Q(y_1))$, with $Q(19) = 0.690$ for the Laptok data (far too low, to really claim that $p_1 < p_0$), and $Q(13) = 0.954$, and even higher, of course, for y_1 even smaller than 13. So from the traditional point of view, one would claim significance only with $y_1 \leq 13$; $y_1 = 19$ is not at all statistically impressive.

(j) Laptok et al. work with their 'neutral prior' where $\gamma \sim N(0, 0.35^2)$, which we have seen corresponds to $\rho \in [0.504, 1.986]$ as the natural 95% prior credibility interval

(with median 1, rather than mean 1; the main figure 2.1 in Laptook et al. (2017) appears to be wrong. Finding Beta priors for (p_0, p_1) to match mean 0 and variance 0.35^2 for γ can be carried out in several ways. The easiest is to let both p_0 and p_1 come from the same $\text{Beta}(a_0, b_0)$, which ensures $\log(p_1/p_0)$ having mean zero, and fixing (a_0, b_0) to have variance

$$\text{Var } \log p = \int_0^1 (\log p)^2 g(p | a_0, b_0) dp - \left\{ \int_0^1 \log p g(p, a_0, b_0) dp \right\}^2$$

half of 0.35^2 . There are again many such Beta densities, but we may e.g. have it centred in 0.25, which means we should find k such that the the logarithm of $\text{Beta}(k \cdot 0.25, k \cdot 0.75)$ has variance precisely equal to $\frac{1}{2} 0.35^2 = 0.06125$. Computing this and fine-tuning for k leads to $k = 51.421$. In other words, if $p \sim \text{Beta}(a_0, b_0)$ with $(a_0, b_0) = (12.855, 38.566)$, then its logarithm has the right variance. – You may draw the corresponding prior density, in this setup common to both p_0 and p_1 , to see how ‘neutral’ the Laptook et al. neutral prior is.

(k) I do not embellish things here; have rather a look at my FocuStat blog post. But, briefly:

- [i] There’s really no serious statistical signpost here saying p_1 is smaller than p_0 ; both the frequentist and Bayesian methods lead to confidence and credibility intervals havint the value $\rho = p_1/p_0 = 1$ smack in the middle.
- [ii] This is a bigger discussion, not easy to sort out well. There is a role for Bayes and priors in clinical trials, but care needs to be exercised. And in case where the frequentist analysis is clear, as here, there’s a risk that Bayes and priors might lead to confusion. Again, see my blog post, *Cooling of Newborns and the Difference Between 0.244 and 0.278*.

The Project, Exercise 2

(a) Taking the logarithm of the likelihood function

$$\prod_{j=1}^{10} \binom{m}{y_j} p_j(a, b)^{y_j} \{1 - p_j(a, b)\}^{m-y_j}$$

leads to the log-likelihood function

$$\ell(a, b) = \sum_{j=1}^{10} \left[y_j \log p_j(a, b) + (m - y_j) \log \{1 - p_j(a, b)\} + \log \binom{m}{y_j} \right].$$

I’ve programmed `logL` in R, as a function of `ab`, complete with `a=ab[1]`, `b=ab[2]`, etc. Via `nlm` for `minuslogL` I find ML estimates $(-1.479, 0.949)$. Standard errors, via the inverse observed Fisher information matrix, are $(0.690, 0.546)$.

With the right plotting tools I then construct the figure given in the project.

- (b) I set up a MCMC scheme in the same fashion I've used for some exercises in the course, giving me a $\text{sim} \times 2$ matrix for a chain of (a, b) outcomes. I use an easy normal proposal, of the symmetric type

$$(a', b') = (a + \varepsilon_1, b + \varepsilon_2),$$

with $\varepsilon_1 \sim N(0, \sigma_1^2)$ and $\varepsilon_2 \sim N(0, \sigma_2^2)$. In principle there's a broad range of σ_1 and σ_2 which will work. I've used $\sigma_1 = \frac{1}{2} 0.690$ and $\sigma_2 = \frac{1}{2} 0.546$, half of the standard errors for the ML estimators (but, again, lots of other schemes will work).

- The posterior distribution for (a, b) obtained in this way is fairly binormal, with no drama. For a (first row) and b (second row) I record ML and standard error, then mean and standard deviation for the posterior distribution:

$$\begin{array}{cc} -1.479, 0.690 & -1.567, 0.713 \\ 0.949, 0.546 & 1.008, 0.554 \end{array}$$

There is clear and strong agreement, so we are in the Good Terrain of Bernshtein–von Mises, as in the book's Chapter 4.

- (c) I then used simulations to read off the 0.05 quantile $\text{low}(x)$ and 0.95 quantile $\text{up}(x)$, from the posterior distribution of $H(a + bx)$. In other words, first I got hold of 10^5 values of (a_j, b_j) , via the MCMC. Then I read off, for each given x , the consequent 10^5 numbers $H(a_j + b_j x)$. From these I read off the 0.05 and 0.95 quantiles, via `quantile`.
- (d) I can also generate 10^5 values of $p(x_{\text{new}}) = H(a_j + b_j x_{\text{new}})$, look at a histogram, etc. The histogram is skewed to the left. Quantiles for 0.05, 0.50, 0.95 are 0.415, 0.719, 0.907. This gives good assessment for $p(x)$ at this new position, provided the model assumption still holds, which might not be the case; cf. the Aukrust brothers.
- (e) I generate 10^5 values of (a_j, b_j) via the MCMC; then, for each of these pairs, I generate y_j from the binomial (m, p_j) , where $p_j = H(a_j + b_j x_{\text{new}})$. This gives the predictive distribution for y_{new} . From 10^5 simulations, I read off probabilities

$$0.0139, 0.0569, 0.1369, 0.2403, 0.3100, 0.2420$$

for values 0, 1, 2, 3, 4, 5.

- (f) Now the restriction is $b \geq 0$ for the (a, b) pair. I've set up a Metropolis MCMC where I make sure that if an (a, b) is proposed with negative b , then that proposal has zero chance of being accepted. This leads to a new MCMC, again of size $\text{sim} \times 2$, and with every (a_j, b_j) respecting $b_j \geq 0$. The posterior distribution is not binormal anymore, with a sharp boundary $b = 0$. I can read off $x_{0,j} = -a_j/b_j$, the LD50 parameter. I can display a histogram and compute quantiles. The 0.05, 0.50, 0.95 quantiles of the posterior distribution for x_0 are 1.051, 1.545, 3.114.

This is a parameter where the usual normal theory does not work well, but where the Bayesian scheme of things work out nicely (as long as the MCMC is set up correctly).

The four-hour exam

Exercise 1

- (a) First, $y_i | \theta$ is a binomial $(1, \theta)$, with mean θ and variance $\theta(1 - \theta)$. Secondly, $z = \sum_{i=1}^n y_i$ given θ is binomial (n, θ) , hence with mean and variance $n\theta$ and $n\theta(1 - \theta)$.
- (b) Now $\theta \sim \text{Beta}(2, 2)$.
- (i) The mean is $\frac{1}{2}$, the variance is $\frac{1}{2} \frac{1}{2} / 5 = 1/20$.
- (ii) $E y_i = E E(y_i | \theta) = E \theta = \frac{1}{2}$, and

$$\begin{aligned} \text{Var } y_i &= E \text{Var}(y_i | \theta) + \text{Var } E(y_i | \theta) \\ &= E \theta(1 - \theta) + \text{Var } \theta \\ &= \frac{1}{2} - E \theta^2 + E \theta^2 - \frac{1}{4} = \frac{1}{4}. \end{aligned}$$

- Also, y_i is a 0-1 variable, with $P(y_i = 1) = \frac{1}{2}$; hence its variance must be $\frac{1}{2} \frac{1}{2} = \frac{1}{4}$.
- (iii) We have $E(y_i y_j | \theta) = \theta^2$ and hence $\text{cov}(y_i, y_j) = E \theta^2 - (\frac{1}{2})^2 = \text{Var } \theta$, i.e. $1/20$. Hence the correlation is

$$\rho = \frac{1/20}{1/4} = \frac{4}{20} = 1/5 = 0.20.$$

- (iv) We have

$$E z = E n\theta = \frac{1}{2}n$$

and

$$\text{Var } z = E n\theta(1 - \theta) + \text{Var } n\theta = n(\frac{1}{2} - \frac{1}{4} - 1/20) + n^2/20 = 0.20 n + 0.05 n^2.$$

- (c) The distribution of z is

$$\begin{aligned} f(z) &= \int_0^1 \binom{n}{z} \theta^z (1 - \theta)^{n-z} 6\theta(1 - \theta) d\theta \\ &= \frac{n!}{z!(n-z)!} \frac{(z+1)!(n-z+1)!}{(n+2)!} \\ &= 6 \frac{(z+1)(n-z+1)}{n(n+1)} \end{aligned}$$

for $z = 0, 1, \dots, n$.

- (d) The posterior is proportional to

$$\theta^z (1 - \theta)^{n-z} \theta(1 - \theta) = \theta^{z+1} (1 - \theta)^{n-z+1},$$

which is a Beta($z + 2, n - z + 2$). The conditional mean is

$$\hat{\theta}_B = E(\theta | z) = \frac{z + 2}{n + 4}.$$

- (e) The standard estimator is $\tilde{\theta} = z/n$, which is unbiased. Hence the risk, the mean of $(\tilde{\theta} - \theta)^2$, is $\theta(1 - \theta)/n$. It starts and ends at zero, and is biggest for $\theta = \frac{1}{2}$.
- (f) The Bayes estimator has risk function

$$\begin{aligned} r(\theta) &= \text{Var} \hat{\theta}_B + (\text{E} \hat{\theta}_B - \theta)^2 \\ &= \frac{n\theta(1 - \theta)}{(n + 4)^2} + \left(\frac{n\theta + 2}{n + 4} - \frac{n\theta + 4\theta}{n + 4} \right)^2 \\ &= \frac{n}{(n + 4)^2} \theta(1 - \theta) + \frac{(4\theta - 2)^2}{(n + 4)^2}. \end{aligned}$$

It is better than the usual ML estimator for those θ where

$$\frac{16(\theta - \frac{1}{2})^2}{(n + 4)^2} \leq \left\{ \frac{1}{n} - \frac{n}{(n + 4)^2} \right\} \theta(1 - \theta),$$

or

$$16(\theta - \frac{1}{2})^2 \leq \frac{(n + 4)^2 - n^2}{n} \theta(1 - \theta) = \frac{16 + 8n}{n} \theta(1 - \theta),$$

or $(\theta - \frac{1}{2})^2 \leq (\frac{1}{2} + 1/n)\theta(1 - \theta)$. This means a certain interval around $\frac{1}{2}$, where Bayes is better. For large n , the inequality is close to $(\theta - \frac{1}{2})^2 \leq \frac{1}{2}\theta(1 - \theta)$, and this holds for θ inside $\frac{1}{2} \pm \sqrt{3}/6$, which is $[0.211, 0.789]$, i.e. a pretty wide interval.

Exercise 2

- (a) The likelihood function becomes

$$\begin{aligned} L &= (p(1 - q))^{212} ((1 - p)q)^{103} (pq)^{39} ((1 - p)(1 - q))^{148} \\ &= p^{212+39} (1 - p)^{103+148} q^{103+39} (1 - q)^{212+148} \\ &= p^{251} (1 - p)^{251} q^{142} (1 - q)^{360}. \end{aligned}$$

- (b) With independent uniforms for p and q , the posteriors are also independent, with

$$p \mid \text{data} \sim \text{Beta}(252, 252), \quad q \mid \text{data} \sim \text{Beta}(143, 361).$$

- (c) The expected number of AB cases, if this theory is correct, is

$$e_{AB} = \text{E}(npq \mid \text{data}) = n \frac{252}{502} \frac{143}{502} = 71.5.$$

But this is far off from the observed 39. So the theory looks very suspicious, indeed. As Landsteiner and others found out, about a hundred years ago, the two-loci theory stinks and sucks; the one-locus theory, however, is splendid. – One may put in more detail here, including computing the probability that one should get a number as far off as 39 (or more), as measured through the lens of the posterior distribution for npq . This will be a microscopic probability. The essence is simply to compare the observed

far too small 39 with the mean of 71.5 – and, of course, similar calculations for the other three cells.

Exercise 3

- (a) The posterior density $p(\theta)$ is the derivative of the cumulative $P(\theta)$, and one finds $\theta \exp(-\theta)$. This is also a gamma $(2, 1)$. Its mean is 2. The density is zero at zero, climbs to $\exp(-1) = 0.368$ at the value 1, and then decreases slowly to zero.
- (b) The Bayes decision is to pick among A, B, C the action that has the smallest expected posterior loss. These three expected posterior losses are

$$\begin{aligned} 2 \{1 - P(1.1)\} &= 2 \cdot 0.6990 = 1.3981, \\ 3 \{P(1.1) + 1 - P(3.3)\} &= 3 \cdot 0.4596 = 1.3787, \\ 4 P(3.3) &= 4 \cdot 9.8414 = 3.3656, \end{aligned}$$

for respectively A, B, C. So we take action B.

- (c) The risk function, for a credibility interval $[\hat{a}, \hat{b}]$ constructed from the data, becomes

$$r(\theta) = E_{\theta} L(\theta, [\hat{a}, \hat{b}]) = 0.10 E_{\theta} (\hat{b} - \hat{a}) + 1 - P_{\theta}\{\theta \in [\hat{a}, \hat{b}]\}.$$

A good method is one with short expected length and with high probability of containing the right parameter.

- (d) Again we ought to minimise posterior expected loss. This means minimising

$$E \{L(\theta, [a, b]) \mid \text{data}\} = 0.10 (b - a) + P(a) + 1 - P(b)$$

over (a, b) . Taking derivatives leads to the equations $p(a) = 0.10$, $p(b) = 0.10$. This again means finding a to the left of 1 and b to the right of 1, as solutions to $p(x) = 0.10$. I find $[a, b] = [0.112, 3.576]$.

Exercise 4

- (a) The likelihood for the data becomes

$$\prod_{i=1}^n \{(3\theta y_i^2) \exp(-\theta y_i^3)\} \propto \theta^n \exp(-n\theta W_n),$$

where $W_n = (1/n) \sum_{i=1}^n y_i^3$. Its logarithm is

$$\ell_n(\theta) = n \log \theta - n\theta W_n,$$

with first derivative $n/\theta - nW_n$ and second derivative $-n/\theta^2$. So the maximum likelihood (ML) estimator is

$$\hat{\theta} = 1/W_n.$$

(b) The prior times the likelihood is proportional to

$$\theta^{a-1} \exp(-b\theta) \theta^n \exp(-n\theta W_n) = \theta^{a+n-1} \exp\{-(b+nW_n)\theta\},$$

which means the posterior is a Gamma with parameters $(a+n, b+nW_n)$. Its mean, by the way, is the Bayes estimator

$$\hat{\theta}_B = \frac{a+n}{b+nW_n},$$

which is close to the ML estimator.

(c) First, by frequentist ML theory, and assuming the model is actually correct,

$$\hat{\theta} \approx_d N(\theta_0, \theta_0^2/n),$$

with θ_0 signalling the true parameter value. Secondly, by Bayes theory for larger sample sizes,

$$\theta | \text{data} \approx_d N(\hat{\theta}, \hat{\theta}^2/n).$$

So there's a mirror situation, and the Bayesian and the frequentist will have the same inferences, for large n .