

---

**MATH 80601A *Bayesian Modelling***

**Practice final examination**

Exam booklet

Examiner: Léo Belzile

---

**Instructions:** The time allotted for the examination is 180 minutes. You may answer in either English or French. No written material may be brought into the examination, but a simple (non-programmable) calculator may be used.

There are a total of 45 marks available in the exam paper, the distribution of which can be found in the right margin.

You must hand back the **exam booklet** at the end of the examination.

---

Question:	1	2	3	4	Total
Points:	10	9	10	16	45
Score:					

## Crib sheet

1.  $\text{binom}(n, p)$  with  $n$  fixed and  $p \in [0, 1]$ , has mass function  $\binom{n}{p} p^y (1-p)^{n-y}$  for  $y \in \{0, \dots, n\}$ . The expectation is  $np$  and the variance  $np(1-p)$ .
2.  $\text{beta}(a, b)$  random variable on  $[0, 1]$  has density  $f(x) = \Gamma(a+b) / \{\Gamma(a)\Gamma(b)\} x^{a-1} (1-x)^{b-1}$  and expectation  $a/(a+b)$ .
3.  $\text{Poisson}(\lambda)$  with expectation  $\lambda > 0$  has mass function  $f(x) = \lambda^x / x! \exp(-\lambda)$ ,  $x \in 0, 1, \dots$
4.  $\text{Gauss}(\mu, \sigma^2)$  with expectation  $\mu \in \mathbb{R}$  and variance  $\sigma^2 > 0$ , has density

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}, \quad x \in \mathbb{R}.$$

5.  $\text{Gauss}_p(\boldsymbol{\mu}, \mathbf{Q}^{-1})$  with expectation  $\boldsymbol{\mu} \in \mathbb{R}^p$  and precision (reciprocal variance)  $\mathbf{Q} \succeq 0$ , has density

$$f(\mathbf{x}) = (2\pi)^{-p/2} |\mathbf{Q}|^{1/2} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \mathbf{Q}(\mathbf{x}-\boldsymbol{\mu})\right\}, \quad \mathbf{x} \in \mathbb{R}^p.$$

6.  $\text{Student}(\mu, \sigma, \nu)$  with location  $\mu \in \mathbb{R}$ , scale  $\sigma > 0$  and  $\nu$  degrees of freedom, has density

$$f(x; \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sigma \Gamma\left(\frac{\nu}{2}\right) \sqrt{\nu\pi}} \left(1 + \frac{(x-\mu)^2}{\sigma^2 \nu}\right)^{-\frac{\nu+1}{2}}, \quad x \in \mathbb{R}.$$

7.  $\text{gamma}(\alpha, \beta)$  with shape  $\alpha > 0$  and rate  $\beta > 0$ , has expectation  $\alpha/\beta$ , variance  $\alpha/\beta^2$  and density

$$f(x; \alpha, \beta) = \beta^\alpha / \Gamma(\alpha) x^{\alpha-1} \exp(-\beta x), \quad x > 0,$$

where  $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} \exp(-x) dx$  is the gamma function, and  $x\Gamma(x) = \Gamma(x+1)$  for  $x > 0$ . If  $\alpha = 1$ , we recover  $\text{expo}(\beta)$ .

8.  $\text{inv.gamma}(\alpha, \beta)$  with shape  $\alpha > 0$  and rate  $\beta > 0$ , has expectation  $\beta/(\alpha-1)$  for  $\alpha > 1$  and density

$$f(x; \alpha, \beta) = \beta^\alpha / \Gamma(\alpha) x^{-\alpha-1} \exp(-\beta/x), \quad x > 0.$$

9.  $\text{Laplace}(\mu, \sigma)$  with mean  $\mu \in \mathbb{R}$  and scale  $\sigma > 0$  with density  $f(x) = (2\sigma)^{-1} \exp(-|x-\mu|/\sigma)$  for  $x \in \mathbb{R}$ . It's variance is  $2\sigma^2$ .

**Jeffrey's prior:**  $p(\theta) \propto |I|^{1/2}$ , where  $I$  is the unit Fisher (expected) information.

**Change of variable formula:** consider an injective (one-to-one) differentiable function  $\mathbf{g} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , with inverse  $\mathbf{g}^{-1}$ . Then, if  $\mathbf{Y} = \mathbf{g}(\mathbf{X})$ , the density of  $\mathbf{Y}$  is

$$f_Y(\mathbf{y}) = f_X\{\mathbf{g}^{-1}(\mathbf{y})\} |\mathbf{J}_{\mathbf{g}^{-1}}(\mathbf{y})| = f_X(\mathbf{x}) |\mathbf{J}_{\mathbf{g}}(\mathbf{x})|^{-1},$$

where  $\mathbf{J}_{\mathbf{g}}(\mathbf{x})$  is the Jacobian matrix with  $(i, j)$ th element  $\partial[\mathbf{g}(\mathbf{x})]_i / \partial x_j$ .

**Metropolis–Hastings algorithm:** consider  $p(\boldsymbol{\theta} | \mathbf{y})$  and  $q(\boldsymbol{\theta} | \boldsymbol{\theta}_{t-1})$ , the proposal density evaluated at  $\boldsymbol{\theta}$ . The acceptance ratio for proposal  $\boldsymbol{\theta}_t^*$  given the current value  $\boldsymbol{\theta}_{t-1}$  is  $\min\{1, R\}$ , where

$$R = \frac{p(\boldsymbol{\theta}_t^* | \mathbf{y})}{p(\boldsymbol{\theta}_{t-1} | \mathbf{y})} \frac{q(\boldsymbol{\theta}_{t-1} | \boldsymbol{\theta}_t^*)}{q(\boldsymbol{\theta}_t^* | \boldsymbol{\theta}_{t-1})}$$

**Effective sample size:** for autocorrelation at lag  $t$  of  $\rho_t$  for a Markov chain of length  $B$ ,  $\text{ESS} = B / \{1 + 2 \sum_{t=1}^{\infty} \rho_t\}$ .

**Law of iterated mean and variance:**

$$\mathbb{E}_Y(Y) = \mathbb{E}_Z\{\mathbb{E}_{Y|Z}(Y)\}, \quad \text{Va}_Y(Y) = \mathbb{E}_Z\{\text{Va}_{Y|Z}(Y)\} + \text{Va}_Z\{\mathbb{E}_{Y|Z}(Y)\}.$$

**Kullback–Leibler divergence:**

$$\text{KL}(f \parallel g) = \int \{\log f(\mathbf{x}) - \log g(\mathbf{x})\} f(\mathbf{x}) d\mathbf{x}.$$

**Quadratic form completion:**

$$(\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a}) + (\mathbf{x} - \mathbf{b})^\top \mathbf{B}(\mathbf{x} - \mathbf{b}) \stackrel{x}{\propto} (\mathbf{x} - \mathbf{c})^\top \mathbf{C}(\mathbf{x} - \mathbf{c}),$$

where  $\mathbf{C} = \mathbf{A} + \mathbf{B}$  and  $\mathbf{c} = \mathbf{C}^{-1}(\mathbf{A}\mathbf{a} + \mathbf{B}\mathbf{b})$ .

**Laplace approximation:** Let  $h(\mathbf{x})$  be a twice-differentiable function which is concave in the vicinity of its mode  $\mathbf{x}_0$ , with  $\mathbf{H}(\mathbf{x})$  the Hessian matrix of  $-h(\mathbf{x})$ . If  $h(\mathbf{x}_0)$  is  $O(n)$ , then as  $n \rightarrow \infty$ ,

$$I_n = \int_{\mathbb{R}^d} \exp\{h(\mathbf{x})\} d\mathbf{x} \approx (2\pi)^{p/2} |\mathbf{H}(\mathbf{x}_0)|^{-1/2} \exp\{h(\mathbf{x}_0)\}.$$

**Evidence lower bound (ELBO):**  $\text{ELBO}(g) = \mathbb{E}_g\{\log p(\boldsymbol{\theta}, \mathbf{y})\} - \mathbb{E}_g\{\log g(\boldsymbol{\theta})\}$ .

**Question 1. True or false** .....**10**

Explain whether each of the following statement is true, false or uncertain. To get marks, you must briefly justify your reasoning by proving the statement or providing a counterexample if the statement is false. **Answers without justifications are worth zero mark.**

- 1.1 A Gaussian approximation  $g$  that minimizes the forward Kullback–Leibler divergence and a Laplace approximation to a univariate density  $p$  will have the same mean. [2]
- 1.2 We can readily obtain marginal posterior moments from variational inference procedures. [2]
- 1.3 The impact of the prior vanishes as the sample size gets larger. [2]
- 1.4 Marginalization in Markov chain Monte Carlo methods is always beneficial (fewer parameters, lower dependence between components). [2]
- 1.5 Consider a prediction  $\tilde{y}$  from a frequentist model with density  $f(\tilde{y}; \theta)$ , where  $\hat{\theta}$  is a point estimator. The posterior predictive distribution  $p(\tilde{y} | y)$  will be more variable than it's frequentist counterpart  $f(\tilde{y}; \hat{\theta})$ . [2]

**Question 2. Short questions** .....**9**

- 2.1 Zellner (1996) proposed the maximal data information (MDI) prior. The latter is proportional to the exponential of the negative entropy,  $p(\theta) \propto \exp\{E_g(\log g)\}$ . [3]
  - Show that the entropy of a Gaussian random variable  $\text{Gauss}(\mu, \sigma^2)$  is a function of the scale  $\sigma$  only.
  - Use this result to derive the corresponding MDI prior for  $(\mu, \sigma)$ .
  - Is the resulting prior proper? Justify your answer.
- 2.2 Define the concepts of **burn in**, **warmup**, and **thinning** for Markov chain Monte Carlo. Explain their relevance in the context of a sampling-based Bayesian analysis. [3]

- 2.3 The following **R** code implements a simple Markov chain Monte Carlo to sample from the posterior mean and std. deviation  $(\mu, \sigma)$ , where we assume independent and identically observations and a half-Cauchy prior on  $[0, \infty)$  for the scale, assumed independent apriori

[3]

$$Y_i | \mu, \sigma \sim \text{Gauss}(\mu, \sigma^2), i = 1, \dots, n; \quad p(\mu) \propto 1; \quad p(\sigma) \propto \text{Student}_+(1, 0, 1).$$

Find the error in the code and explain why this isn't sampling from the posterior distribution of interest.

```
sd_prop <- 0.05
# Unnormalized log posterior for mu, sigma
logpost <- function(pars, y){
  mu <- pars[1]; sigma <- pars[2]
  if(sigma < 0){ return(-Inf)}
  sum(dnorm(x = y, mean = mu, sd = sigma, log = TRUE)) + # log likelihood
  log(2) + dt(x = sigma, df = 1, log = TRUE) # log prior for scale
}

B <- 1e4L # number of simulations
mu_s <- sigma_s <- numeric(B) # containers
mu <- mean(y); sigma <- sd(y) # initial values
for(b in 1:B){
  # Gibbs step for mu
  mu <- mu_s[b] <- rnorm(n = 1, mean = mean(y), sd = sigma/sqrt(length(y)))
  # Metropolis random walk on the log scale
  sigma_prop <- exp(rnorm(n = 1, mean = log(sigma), sd = sd_prop))
  # Log of Metropolis acceptance ratio
  logR <- logpost(c(mu, sigma_prop), y = y) - logpost(c(mu, sigma), y = y)
  if(logR > log(runif(1))){
    sigma <- sigma_prop
  }
  sigma_s[b] <- sigma
}
```

**Question 3. Stochastic versus deterministic approximations .....****10**

- 3.1 Define the **marginal likelihood** as a function of the prior  $p(\boldsymbol{\theta})$  and the likelihood  $p(\mathbf{y} | \boldsymbol{\theta})$ . [2]
- 3.2 Explain in your own words why estimation of the marginal likelihood is difficult. [2]
- 3.3 Explain how to obtain a Laplace approximation of the marginal likelihood. Under which circumstances will it be a good approximation? [2]
- 3.4 Explain the relevance of the marginal likelihood in the context of [2]
1. calculation of posterior moments of the form  $E_{\boldsymbol{\theta} | \mathbf{Y}}\{g(\boldsymbol{\theta})\}$ .
  2. model comparison using Bayes factor.
- 3.5 How does Markov chain Monte Carlo (MCMC) methods get around estimation of the marginal likelihood  $p(\mathbf{y})$ ? *Hint*: consider the Metropolis–Hastings acceptance ratio. [2]

**Question 4. Probit regression .....****16**

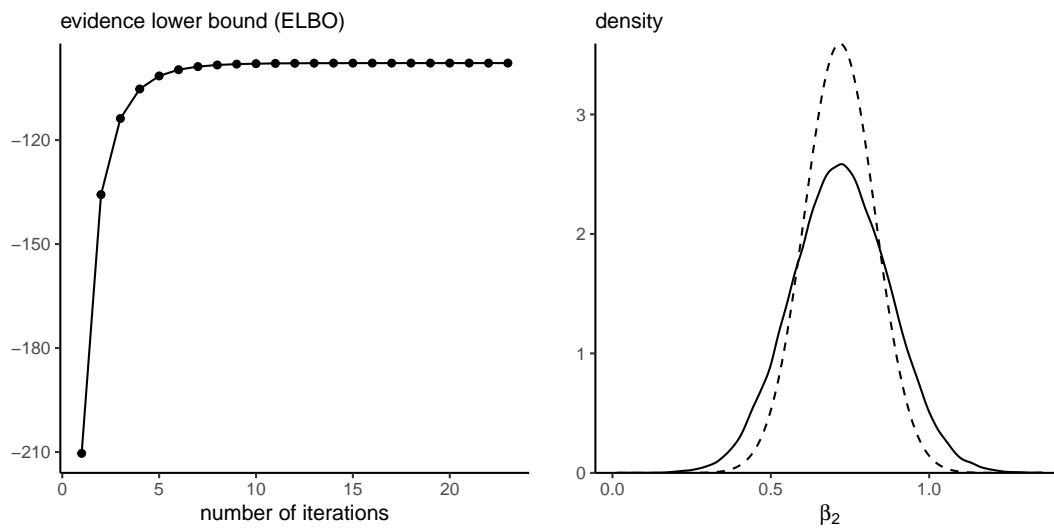
Experiment 2 of Duke and Amir (2023) consider the effect of sequential versus integrated decisions on customers decision to buy. Customers in an online experiment were exposed to products and decided whether to buy ( $Y_i = 1$ ) or not ( $Y_i = 0$ ). To model these, we consider a simple probit regression model with response  $Y_i \sim \text{binom}(1, p_i)$ , where

$$p_i = \Pr(Y_i = 1) = \Phi(\mathbf{x}_i \boldsymbol{\beta}),$$

with  $\Phi(\cdot)$  the distribution function of the standard Gaussian distribution. We set  $\boldsymbol{\beta} \sim \text{Gauss}_p(\mathbf{0}_p, c\mathbf{I}_p)$  for  $c > 0$  a known positive constant.

The model matrix include an intercept, a coefficient for age and a binary indicator equal to 1 if the participant was exposed to quantity-integrated decision, and zero for quantity-sequential (control group).

- 4.1 Consider the data augmentation scheme where  $Y_i = I(Z_i > 0)$ , where  $Z_i \sim \text{Gauss}(\mathbf{x}_i \boldsymbol{\beta}, 1)$ , with  $\mathbf{x}_i$  denoting the  $i$ th row of the  $n \times p$  design matrix. [2]
- Write down the expression for the joint distribution  $p(\mathbf{y}, \mathbf{z}, \boldsymbol{\beta}) = p(\mathbf{y} | \mathbf{z})p(\mathbf{z} | \boldsymbol{\beta})p(\boldsymbol{\beta})$ .
- 4.2 Derive the conditional distributions  $p(\boldsymbol{\beta} | \mathbf{z})$  and that of  $p(z_i | y_i, \boldsymbol{\beta})$  for  $i = 1, \dots, n$ . [4]
- 4.3 Based on the conditional distributions detail a Gibbs sampling algorithm for  $\boldsymbol{\beta}$  and  $\mathbf{z}$ . Explain the benefit of the latter over the marginal posterior  $p(\boldsymbol{\beta} | \mathbf{y})$ . [2]
- 4.4 Suppose that we instead used coordinate-ascent variational inference with a factorization of the posterior  $p_{\mathbf{Z}}(\mathbf{z})p_{\boldsymbol{\beta}}(\boldsymbol{\beta})$ . [4]
- Write down the optimal form of these distributions and the parameter updates. Explain how you would assess convergence.



**Figure 1:** Left: evidence lower bound (ELBO) as a function of iteration. Right: marginal density of  $\beta_2$  for quantity-integrated binary indicator, based on Monte Carlo samples (full), versus variational approximation (dashed).

*Hint:* if  $Y \sim \text{trunc.Gauss}(\mu, \sigma, a, b)$  a truncated Gaussian on  $[a, b]$  with location  $\mu$  and scale  $\sigma$  has expectation

$$E(Y) = \mu - \sigma \frac{\phi\{(b - \mu)/\sigma\} - \phi\{(a - \mu)/\sigma\}}{\Phi\{(b - \mu)/\sigma\} - \Phi\{(a - \mu)/\sigma\}},$$

where  $\phi$  and  $\Phi$  are the density and distribution functions of a standard Gaussian, respectively.

4.5 The right panel of Figure 1 shows the marginal density for the coefficient  $\beta_2$ . Explain why the two are not identical. [2]

4.6 What can we conclude from Figure 1 as to what is the most effective method? [2]