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.

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

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

Crib sheet

- 1. binom(n, p) with *n* fixed and $p \in [0, 1]$, has mass function $\binom{n}{p} p^y (1-p)^{n-y}$ for $y \in \{0, ..., n\}$. The expectation is np and the variance np(1-p).
- 2. 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. Poisson(λ) with expectation $\lambda > 0$ has mass function $f(x) = \lambda^x / x! \exp(-\lambda), x \in 0, 1, ...$
- 4. Gauss(μ, σ^2) with expectation $\mu \in \mathbb{R}$ and variance $\sigma^2 > 0$, has density

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

5. Gauss_{*p*}($\boldsymbol{\mu}, \boldsymbol{Q}^{-1}$) with expectation $\boldsymbol{\mu} \in \mathbb{R}^{p}$ and precision (reciprocal variance) $\boldsymbol{Q} \geq 0$, has density

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

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

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

7. gamma(α , β) with shape $\alpha > 0$ and rate $\beta > 0$, has expectation α / β , variance α / β^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 $\exp(\beta)$.

8. inv.gamma(α , β) 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. Laplace(μ, σ) 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:** $p(\theta) \propto |\iota|^{1/2}$, where ι is the unit Fisher (expected) information.

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

$$f_{\boldsymbol{Y}}(\boldsymbol{y}) = f_{\boldsymbol{X}} \left\{ \boldsymbol{g}^{-1}(\boldsymbol{y}) \right\} \left| \boldsymbol{J}_{\boldsymbol{g}^{-1}}(\boldsymbol{y}) \right| = f_{\boldsymbol{X}}(\boldsymbol{x}) \left| \boldsymbol{J}_{\boldsymbol{g}}(\boldsymbol{x}) \right|^{-1},$$

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

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

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

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

$$\mathsf{E}_{Y}(Y) = \mathsf{E}_{Z}\left\{\mathsf{E}_{Y|Z}(Y)\right\}, \qquad \mathsf{Va}_{Y}(Y) = \mathsf{E}_{Z}\left\{\mathsf{Va}_{Y|Z}(Y)\right\} + \mathsf{Va}_{Z}\left\{\mathsf{E}_{Y|Z}(Y)\right\}.$$

Kullback-Leibler divergence:

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

Quadratic form completion:

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

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 it's 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 $n \to \infty$,

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

Evidence lower bound (ELBO): ELBO(g) = E_g{log $p(\theta, y)$ } – E_g{log $g(\theta)$ }.

10

[2]

[2]

Question 1. True or false

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 [2] Laplace approximation to a univariate density *p* will have the same mean.

Solution: False; the Laplace approximation will be centered at the mode, while the other model will be centered at the mean of the true distribution. They coincide for symmetric unimodal distributions.

1.2 We can readily obtain marginal posterior moments from variational inference procedures.

Solution: Uncertain; the optimal decomposition might lead to a multivariate family, which may be unnormalized. Gaussian approximations will of course be marginalizable, and we could resort to Monte Carlo methods if we can simulate.

1.3 The impact of the prior vanishes as the sample size gets larger.

Solution: False; this depends largely on what level the parameter or prior is located; for example, a random effect model with different variance per group would vanish only if the number of groups increases (in addition to the number of observations within each group). Other counterexample: any prior that restricts the support would not vanish.

1.4 Marginalization in Markov chain Monte Carlo methods is always beneficial (fewer parameters, [2] lower dependence between components).

Solution: False; we have seen examples where parameter expansion improves mixing for the eight school example by adding redundancy. Generally, it complexifies the dependence structure, while reducing the dimension (thinks data augmentation in Gibbs sampling).

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 $f(\tilde{y}; \hat{\theta})$.

Solution: True, follows from the law of iterated variance as

$$p(\widetilde{y} \mid \mathbf{y}) = \int f(\widetilde{y}; \boldsymbol{\theta}, \mathbf{y}) p(\boldsymbol{\theta} \mid \mathbf{y}) \mathrm{d}\boldsymbol{\theta}$$

page 1 of 8

so $Va_{\widetilde{y}|y}(\widetilde{Y}) = E_{\theta|y} \{ Va_{\widetilde{y}|\theta}(\widetilde{Y}) \} + Va_{\theta|y} \{ E_{\widetilde{y}|\theta}(\widetilde{Y}) \}$. The frequentist version amounts to taking a point mass at $p(\theta | y) = I(\theta = \widehat{\theta})$, so the second term vanishes as there is no variability.

Question 2. Short questions

- 2.1 Zellner (1996) proposed the maximal data information (MDI) prior. The latter is proportional to the exponential of the negative entropy, $p(\theta) \propto \exp\{E_g(\log g)\}$.
 - Show that the entropy of a Gaussian random variable $Gauss(\mu, \sigma^2)$ is a function of the scale σ only.
 - Use this result to derive the corresponding MDI prior for (μ, σ) .
 - Is the resulting prior proper? Justify your answer.

Solution:

$$-\mathsf{E}_{g}(\log g) = \frac{1}{2}\log(2\pi) + \log(\sigma) + \frac{1}{2\sigma^{2}}\mathsf{E}_{X}\{(X-\mu)^{2}\} = \frac{1+\log(2\pi)}{2} + \log(\sigma)$$

so the MDI prior is proportional to $1/\sigma$, hence is improper.

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.

Solution: Burn in refers to discarding the initial transient runs of the Markov chain, to ensure that it has converged to the stationnary distribution. Warmup is used sometimes interchangedly, but also serves for tuning parameters of the algorithm. Both of these initial runs are discarded after sampling is performed, to avoid biasing the results.

Thining refers to the practice of keeping only a fraction of the draws. This is only useful when there is strong autocorrelation to reduce storage costs.

9

[3]

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

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

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){</pre>
  mu <- pars[1]; sigma <- pars[2]</pre>
  if(sigma < 0){ return(-Inf)}</pre>
  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</pre>
mu <- mean(y); sigma <- sd(y) # initial values</pre>
for(b in 1:B){
  # Gibbs step for mu
 mu <- mu_s[b] <- rnorm(n = 1, mean = mean(y), sd = sigma/sqrt(length(y)))</pre>
  # Metropolis random walk on the log scale
  sigma_prop <- exp(rnorm(n = 1, mean = log(sigma), sd = sd_prop))</pre>
  # Log of Metropolis acceptance ratio
  logR <- logpost(c(mu, sigma_prop), y = y) - logpost(c(mu, sigma), y = y)</pre>
  if(logR > log(runif(1))){
    sigma <- sigma_prop
  }
  sigma_s[b] <- sigma</pre>
}
```

Solution: The Jacobian of the transformation is missing; the ratio of proposal densities isn't zero since the lognormal is not symmetric.

Question 3. Stochastic versus deterministic approximations

3.1 Define the **marginal likelihood** as a function of the prior $p(\theta)$ and the likelihood $p(y | \theta)$.

Solution: The marginal likelihood is the normalizing constant for the posterior, defined as

$$p(\mathbf{y}) = \int p(\mathbf{y} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

3.2 Explain in your own words why estimation of the marginal likelihood is difficult.

Solution: The marginal likelihood is often small (think discrete components with probability) and decreases with *n*. We have seen examples where numerical integral will fail because of numerical overflow. Likewise, Monte Carlo does not readily work because sampling from the prior unless the latter aligns with the likelihood model.

3.3 Explain how to obtain a Laplace approximation of the marginal likelihood. Under which circumstances will it be a good approximation? [2]

Solution: Write $\hat{\theta}$ for the maximum a posteriori (MAP) and $\mathbf{H}(\hat{\theta})$ for the Hessian of the negative log posterior evaluate at the MAP. Then a straightforward application of Laplace approximation gives

$$p(\mathbf{y}) \approx (2\pi)^{p/2} |\mathbf{H}(\widehat{\boldsymbol{\theta}})|^{-1/2} p(\mathbf{y} | \widehat{\boldsymbol{\theta}}) p(\widehat{\boldsymbol{\theta}})$$

It is a good approximation if the sample size n is large and the posterior is approximately symmetric in the parameterization considered.

- 3.4 Explain the relevance of the marginal likelihood in the context of
 - 1. calculation of posterior moments of the form $E_{\Theta|Y}\{g(\theta)\}$.
 - 2. model comparison using Bayes factor.

Solution: Given the joint $p(y, \theta)$, any expectation is of the form

$$\mathsf{E}_{\boldsymbol{\Theta}|\boldsymbol{Y}}\{g(\boldsymbol{\theta})\} = p(\boldsymbol{y})^{-1} \int g(\boldsymbol{\theta}) p(\boldsymbol{y}, \boldsymbol{\theta}) \mathrm{d}\boldsymbol{\theta}$$

we need the normalizing constant to compute the moments.

[2]

10

Model comparison in Bayesian is based on the posterior model odds, which are

$$p(\mathcal{M}_i \mid \mathbf{y}) \propto p(\mathbf{y} \mid \mathcal{M}_i) p(\mathcal{M}_i).$$

For two models \mathcal{M}_1 and \mathcal{M}_2 , the Bayes factor involve the ratio of marginal likelihoods $p(\mathbf{y} | \mathcal{M}_1) / p(\mathbf{y} | \mathcal{M}_2)$. These require necessarily proper priors.

3.5 How does Markov chain Monte Carlo (MCMC) methods get around estimation of the marginal likelihood p(y)? *Hint:* consider the Metropolis–Hastings acceptance ratio.

Solution: We compare current value with the proposal; since the marginal likelihood is a constant, taking ratios of posterior leads to it's cancellation in the Metropolis–Hastings ratio.

Once we have posterior samples, we can evaluate any functional of interest via Monte Carlo, circumventing the evaluation of the normalizing constant. We however cannot use this for model comparison.

Question 4. Probit regression

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 where 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 $\beta \sim \text{Gauss}_p(\mathbf{0}_p, c\mathbf{I}_p)$ 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 $Z_i \sim Gauss(\mathbf{x}_i \boldsymbol{\beta}, 1)$, with \mathbf{x}_i denoting the *i*th row of the $n \times p$ design matrix.

Write down the expression for the joint distribution $p(y, z, \beta) = p(y | z) p(z | \beta) p(\beta)$.

Solution:

$$p(\boldsymbol{z}, \boldsymbol{y}, \boldsymbol{\beta}) \propto \exp\left\{-\frac{1}{2c}\boldsymbol{\beta}^{\top}\boldsymbol{\beta} - \frac{1}{2}(\boldsymbol{z} - \mathbf{X}\boldsymbol{\beta})^{\top}(\boldsymbol{z} - \mathbf{X}\boldsymbol{\beta})\right\} \times \prod_{i=1}^{n} |(z_i > 0)^{y_i}| (z_i \le 0)^{1-y_i}$$

4.2 Derive the conditional distributions $p(\boldsymbol{\beta} | \boldsymbol{z})$ and that of $p(z_i | y_i, \boldsymbol{\beta})$ for i = 1, ..., n.

[2]

Solution: A simple Gaussian completion of squares gives

$$\boldsymbol{\beta} \mid \boldsymbol{z}, \boldsymbol{y} \sim \text{Gauss}_p \left\{ (\mathbf{X}^\top \mathbf{X} + c^{-1} \mathbf{I}_p)^{-1} \mathbf{X}^\top \boldsymbol{z}, (\mathbf{X}^\top \mathbf{X} + c \mathbf{I}_p)^{-1} \right\}.$$

The augmented variables Z_i are conditionally independent and truncated Gaussian with unit variance

$$Z_i \mid y_i, \boldsymbol{\beta} \sim \begin{cases} \text{trunc.Gauss}(\mathbf{x}_i \boldsymbol{\beta}, 1, -\infty, 0) & y_i = 0\\ \text{trunc.Gauss}(\mathbf{x}_i \boldsymbol{\beta}, 1, 0, \infty) & y_i = 1. \end{cases}$$

4.3 Based on the conditional distributions detail a Gibbs sampling algorithm for β and z. Explain the benefit of the latter over the marginal posterior $p(\beta | y)$.

Solution: We obtain conditional conjugacy, which simplifies the likelihood. The algorithm is also automatic since the acceptance rate of Gibbs sampling is one; there are no tuning parameters.

4.4 Suppose that we instead used coordinate-ascent variational inference with a factorization of [4] the posterior $p_Z(z)p_{\beta}(\beta)$.

Write down the optimal form of these distributions and the parameter updates. Explain how you would assess convergence.

Hint: if *Y* ~ trunc.Gauss(μ , σ , a, b) a truncated Gaussian on [a, b] with location μ and scale σ has expectation

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

where ϕ and Φ are the density and distribution functions of a standard Gaussian, respectively.

Solution: If we consider a factorization of the form $g_Z(z)g_\beta(\beta)$, then we exploit the conditionals in the same way as for Gibbs sampling, but substituting unknown parameter functionals by their expectations. Furthermore, the optimal form of the density further factorizes as $g_Z(z) = \prod_{i=1}^n g_{Z_i}(z_i)$.

The model depends on the mean parameter of Z, say μ_Z , and that of β , say μ_β . To see this, consider the terms in the posterior proportional to Z_i , where

$$p(z_i \mid \boldsymbol{\beta}, y_i) \propto -\frac{z_i^2 - 2z_i \mathbf{x}_i \boldsymbol{\beta}}{2} \times \mathbf{I}(z_i > 0)^{y_i} \mathbf{I}(z_i < 0)^{1 - y_i}$$

which is linear in $\boldsymbol{\beta}$. The expectation of a univariate truncated Gaussian $Z \sim \text{trunc.Gauss}(\mu, \sigma^2, l, u)$ is

$$\mathsf{E}(Z) = \mu - \sigma \frac{\phi\{(u - \mu/\sigma)\} - \phi\{(l - \mu/\sigma)\}}{\Phi\{(u - \mu/\sigma)\} - \Phi\{(l - \mu/\sigma)\}}.$$

If we replace $\mu = \mathbf{x}_i \mu_{\beta}$ in this expression, we get the update

$$\mu_{Z_i}(z_i) = \begin{cases} \mathbf{x}_i \mu_{\boldsymbol{\beta}} - \frac{\phi(\mathbf{x}_i \mu_{\boldsymbol{\beta}})}{1 - \Phi(\mathbf{x}_i \mu_{\boldsymbol{\beta}})} & y_i = 0; \\ \mathbf{x}_i \mu_{\boldsymbol{\beta}} + \frac{\phi(\mathbf{x}_i \mu_{\boldsymbol{\beta}})}{\Phi(\mathbf{x}_i \mu_{\boldsymbol{\beta}})} & y_i = 1, \end{cases}$$

since $\phi(x) = \phi(-x)$.

The optimal form for β is Gaussian and the only unknown parameter is μ_{β} ; the log density only involves a linear form for Z, so follows from the same principle. We get regrouping the vector of means for the latent variables

$$\boldsymbol{\mu}_{\boldsymbol{\beta}} = (\mathbf{X}^{\top}\mathbf{X} + \mathbf{Q}_{\mathbf{0}})^{-1} (\mathbf{X}\boldsymbol{\mu}_{\boldsymbol{Z}} + \mathbf{Q}_{0}\boldsymbol{\mu}_{0})$$

Starting with an initial vector for either parameters, the CAVI algorithm alternatives between updates of μ_{β} and μ_{Z} until the value of the evidence lower bound stabilizes.

4.5 The right panel of Figure 1 shows the marginal density for the coefficient β_2 . Explain why the two are not identical.

Solution: CAVI solves an approximate problem by maximizing the ELBO (or equivalently minimizing the reverse KL divergence), so it needs not coincide. We see that the Gaussian approximation does not fully capture the scale of the marginal.

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

[2]



Figure 1: Left: evidence lower bound (ELBO) as a function of iteration. Right: marginal density of β_2 for quantity-integrated binary indicator, based on Monte Carlo samples (full), versus variational approximation (dashed).

Solution: The posterior for the effect of quantity-integrated (versus control for the intercept has marginal posterior probability $Pr(\beta_2 > 0) \approx 1$. The latter thus indicates that, for the same age, it leads to more sales. This finding is contingent on the model being correct.