

# To-do list

## Contents

<b>1</b>	<b>Preceding chapters</b>	<b>3</b>
<b>A</b>	<b>Functional derivative of the entropy</b>	<b>5</b>
A.1	The usual functional derivative . . . . .	5
A.2	Fréchet differential of the entropy . . . . .	6
<b>B</b>	<b>Kronecker product and vectorisation</b>	<b>9</b>
<b>C</b>	<b>Statistical distributions and their properties</b>	<b>11</b>
C.1	Multivariate normal distribution . . . . .	11
C.2	Matrix normal distribution . . . . .	15
C.3	Truncated univariate normal distribution . . . . .	16
C.4	Truncated multivariate normal distribution . . . . .	17
C.5	Gamma distribution . . . . .	20
C.6	Inverse gamma distribution . . . . .	20
<b>D</b>	<b>Proofs related to the conically truncated independent multivariate normal distribution</b>	<b>21</b>
D.1	Proof of Lemma C.5: Pdf . . . . .	21
D.2	Proof of Lemma C.5: Moments . . . . .	22
D.3	Proof of Lemma C.5: Entropy . . . . .	26
<b>E</b>	<b>I-prior interpretation of the <math>g</math>-prior</b>	<b>27</b>
<b>F</b>	<b>Additional details for various I-prior regression models</b>	<b>29</b>
F.1	The I-prior for standard multilevel models . . . . .	29
F.2	The I-prior for naïve classification . . . . .	31
<b>G</b>	<b>Posterior distribution of the I-prior regression function</b>	<b>35</b>
G.1	Deriving the posterior distribution for $w$ . . . . .	35
G.2	Deriving the posterior predictive distribution . . . . .	36

<b>H</b>	<b>Variational EM algorithm for I-probit models</b>	<b>39</b>
H.1	Derivation of the variational densities . . . . .	39
H.1.1	Derivation of $\tilde{q}(\mathbf{y}^*)$ . . . . .	40
H.1.2	Derivation of $\tilde{q}(\mathbf{w})$ . . . . .	41
H.2	Deriving the ELBO expression . . . . .	44
H.2.1	Terms involving distributions of $\mathbf{y}^*$ . . . . .	44
H.2.2	Terms involving distributions of $\mathbf{w}$ . . . . .	45
<b>I</b>	<b>The Gibbs sampler for the I-prior Bayesian variable selection model</b>	<b>47</b>
I.1	Conditional posterior for $\beta$ . . . . .	48
I.2	Conditional posterior for $\gamma$ . . . . .	48
I.3	Conditional posterior for $\alpha$ . . . . .	49
I.4	Conditional posterior for $\sigma^2$ . . . . .	49
I.5	Conditional posterior for $\kappa$ . . . . .	50
I.6	Computational note . . . . .	50
	<b>Bibliography</b>	<b>51</b>
	<b>Figures</b>	<b>52</b>
	<b>Tables</b>	<b>53</b>
	<b>Theorems</b>	<b>54</b>
	<b>Definitions</b>	<b>55</b>
	<b>Nomenclature</b>	<b>60</b>
	<b>Abbreviations</b>	<b>61</b>

Haziq Jamil

*Department of Statistics*

*London School of Economics and Political Science*

PhD thesis: ‘Regression modelling using Fisher information covariance kernels (I-priors)’

## Chapter 1

# Preceding chapters



## Appendix A

# Functional derivative of the entropy

apx:funcder

We present the functional derivative of the entropy  $H(p)$  in equation 3.6 (p. 17). Typically, this is tackled using calculus of variations, but it can also be obtained using the Fréchet and Gâteaux differentials. Both methods are presented.

### A.1 The usual functional derivative

The functional derivative is defined as follows.

**Definition A.1** (Functional derivative). Given a manifold  $M$  representing continuous/smooth functions  $\rho$  with certain boundary conditions, and a functional  $F : M \rightarrow \mathbb{R}$ , the functional derivative of  $F(\rho)$  with respect to  $\rho$ , denoted  $\partial F / \partial \rho$ , is defined by

$$\begin{aligned} \int \frac{\partial F}{\partial \rho}(x) \phi(x) \, dx &= \lim_{\epsilon \rightarrow 0} \frac{F(\rho + \epsilon \phi) - F(\rho)}{\epsilon} \\ &= \left[ \frac{d}{d\epsilon} F(\rho + \epsilon \phi) \right]_{\epsilon=0}, \end{aligned}$$

where  $\phi$  is an arbitrary function. The function  $\partial F / \partial \rho$  as the gradient of  $F$  at the point  $\rho$ , and

$$\partial F(\rho, \phi) = \int \frac{\partial F}{\partial \rho}(x) \phi(x) \, dx$$

as the directional derivative at point  $\rho$  in the direction of  $\phi$ . Analogous to vector calculus, the inner product with the gradient gives the directional derivative.

Now let  $X$  be a discrete random variable with probability mass function  $p(x) \geq 0$ , for  $\forall x \in \Omega$ , a finite set. The entropy is a functional of  $p$ , namely

$$H(p) = - \sum_{x \in \Omega} p(x) \log p(x).$$

Equivalently, using the counting measure  $\nu$  on  $\Omega$ , we can write

$$H(p) = - \int_{\Omega} p(x) \log p(x) \, d\nu(x).$$

Using the definition of functional derivatives, we find that

$$\begin{aligned} \int_{\Omega} \frac{\partial H}{\partial p}(x) \phi(x) \, dx &= \left[ \frac{d}{d\epsilon} H(p + \epsilon \phi) \right]_{\epsilon=0} \\ &= \left[ - \frac{d}{d\epsilon} (p(x) + \epsilon \phi(x)) \log (p(x) + \epsilon \phi(x)) \right]_{\epsilon=0} \\ &= - \int_{\Omega} \left( \frac{p(x) \phi(x)}{p(x) + \epsilon \phi(x)} + \frac{\epsilon \phi(x)}{p(x) + \epsilon \phi(x)} + \phi(x) \log (p(x) + \epsilon \phi(x)) \right) dx \\ &= - \int_{\Omega} (1 + \log p(x)) \phi(x) \, dx. \end{aligned}$$

Thus,  $(\partial H / \partial p)(x) = -1 - \log p(x)$ .

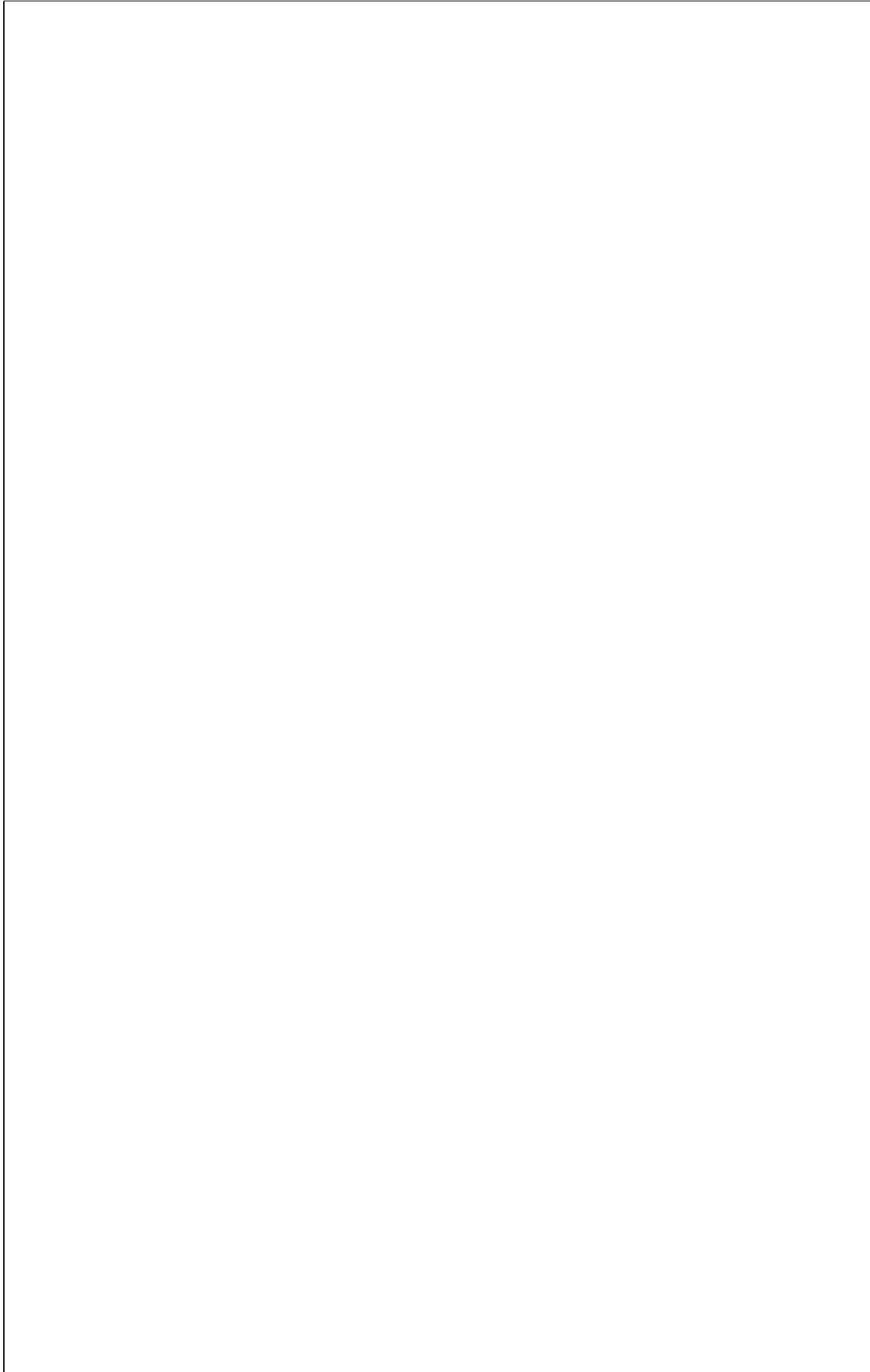
## A.2 Fréchet differential of the entropy

Since we have already introduced concepts of Fréchet and Gâteaux derivatives earlier, we shall use those instead. Assume that the entropy  $H$  is Fréchet differentiable at  $p$ , and that the probability densities  $p$  under consideration belong to the Hilbert space of square integrable functions  $L^2(\Theta, \nu)$  with inner product  $\langle p, p' \rangle_{L^2(\Theta, \nu)} = \int p p' \, d\nu$ . Now since the Fréchet derivative of  $H$  at  $p$  is assumed to exist, it is equal to the Gâteaux

derivative, which can be computed as follows:

$$\begin{aligned}
\partial_q H(p) &= \left. \frac{d}{dt} H(p + tq) \right|_{t=0} \\
&= \left. \frac{d}{dt} \left\{ - \int_{\Theta} (p(\theta) + tq(\theta)) \log (p(\theta) + tq(\theta)) \, d\nu(\theta) \right\} \right|_{t=0} \\
&= - \int_{\Theta} \left\{ \left. \frac{d}{dt} (p(\theta) + tq(\theta)) \log (p(\theta) + tq(\theta)) \right|_{t=0} \right\} d\nu(\theta) \\
&= - \int_{\Theta} \left( \frac{p(\theta)q(\theta)}{p(\theta) + tq(\theta)} + \frac{tq(\theta)^2}{p(\theta) + tq(\theta)} + q(\theta) \log (p(\theta) + tq(\theta)) \right) \Big|_{t=0} d\nu(\theta) \\
&= - \int_{\Theta} q(\theta) (1 + \log p(\theta)) \, d\nu(\theta) \\
&= \langle -(1 + \log p), q \rangle_{\Theta} \\
&= dH(p)(q).
\end{aligned}$$

By definition, the gradient of  $H$  at  $p$ , denoted  $\nabla H(p)$ , is equal to  $-1 - \log p$ . This agrees with the usual functional derivative of the entropy obtained via standard calculus of variations.





## Appendix B

# Kronecker product and vectorisation

The Kronecker product crops up in the definition of matrix normal distributions, which is used in [Chapter 5](#) for the I-probit model.

def:kroneckerprod

**Definition B.1** (Kronecker product). The Kronecker matrix product, denoted by  $\otimes$ , for two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{p \times q}$  is defined by

$$A \otimes B = \begin{pmatrix} A_{11}B & A_{12}B & \cdots & A_{1m}B \\ A_{21}B & A_{22}B & \cdots & A_{2m}B \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1}B & A_{n2}B & \cdots & A_{nm}B \end{pmatrix} \in \mathbb{R}^{np \times mq}.$$

The Kronecker product is a generalisation of the outer product for vectors to matrices. Of use will be these properties of the Kronecker product ([Zhang and Ding, 2013](#)):

- **Bilinearity and associativity.** For appropriately sized matrices  $A$ ,  $B$  and  $C$ , and a scalar  $\lambda$ ,

$$A \otimes (B + C) = A \otimes B + A \otimes C$$

$$(A + B) \otimes C = A \otimes C + B \otimes C$$

$$\lambda A \otimes B = A \otimes \lambda B = \lambda(A \otimes B)$$

$$(A \otimes B) \otimes C = A \otimes (B \otimes C)$$

- **Non-commutative.** In general,  $A \otimes B \neq B \otimes A$ , but they are *permutation equivalent*, i.e.  $A \otimes B \neq P(B \otimes A)Q$  for some permutation matrices  $P$  and  $Q$ .
- **The mixed product property.**  $(A \otimes B)(C \otimes D) = AC \otimes BD$ .

- **Inverse.**  $A \otimes B$  is invertible if and only if  $A$  and  $B$  are both invertible, and  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$ .
- **Transpose.**  $(A \otimes B)^\top = A^\top \otimes B^\top$ .
- **Determinant.** If  $A$  is  $n \times n$  and  $B$  is  $m \times m$ , then  $|A \otimes B| = |A|^m |B|^n$ . Note that the exponent of  $|A|$  is the order of  $B$  and vice versa.
- **Trace.** Suppose  $A$  and  $B$  are square matrices. Then  $\text{tr}(A \otimes B) = \text{tr } A \text{tr } B$ .
- **Rank.**  $\text{rank}(A \otimes B) = \text{rank } A \text{rank } B$ .
- **Matrix equations.**  $AXB = C \Leftrightarrow (B^\top \otimes A) \text{vec } X = \text{vec}(AXB) = \text{vec } C$ .

The equivalence between matrix normal and multivariate normal distributions are established making use of vectorisation for matrices. This is defined below.

def:vectori  
sation

**Definition B.2** (Vectorisation). The vectorisation operation ‘vec’ stacks the columns of the matrices into one long vector, for instance, for the matrix  $A \in \mathbb{R}^{n \times m}$

$$\text{vec } A = (A_{11}, \dots, A_{n1}, A_{12}, \dots, A_{n2}, \dots, A_{1m}, \dots, A_{nm})^\top \in \mathbb{R}^{nm}.$$

## Appendix C

# Statistical distributions and their properties

This is a reference relating to the multivariate normal, matrix normal, truncated univariate and multivariate normal, which are collated from various sources for convenience. Of interest are probability density functions, first and second moments, and entropy (Definition 3.5, page 17). Note that in this part of the appendix, boldface notation for matrix and vectors are not used.

### C.1 Multivariate normal distribution

apx:fisherm  
ultinormal

Let  $X \in \mathbb{R}^d$  be distributed according to a multivariate normal (Gaussian) distribution with mean  $\mu \in \mathbb{R}^d$  and covariance matrix  $\Sigma \in \mathbb{R}^d$  (a square, symmetric, positive-definite matrix). We say that  $X \sim N_d(\mu, \Sigma)$ . Then,

- **Pdf.**  $p(X|\mu, \Sigma) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp \left( -\frac{1}{2} (X - \mu)^\top \Sigma^{-1} (X - \mu) \right)$ .
- **Moments.**  $E X = \mu$ ,  $E[XX^\top] = \Sigma + \mu\mu^\top$ .
- **Entropy.**  $H(p) = \frac{1}{2} \log |2\pi e \Sigma| = \frac{d}{2} (1 + \log 2\pi) + \frac{1}{2} \log |\Sigma|$ .

For  $d = 1$ , i.e.  $X$  is univariate, then its pdf is  $p(X|\mu, \sigma^2) = \frac{1}{\sigma} \phi \left( \frac{X - \mu}{\sigma} \right)$ , and its cdf is  $F(X|\mu, \sigma^2) = \Phi \left( \frac{X - \mu}{\sigma} \right)$ , where  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the pdf and cdf of a univariate standard normal distribution. In the special case that  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ , then the components of  $X = (X_1, \dots, X_d)^\top$  are independently distributed according to  $X_i \sim N(\mu_i, \sigma_i^2)$ .

**Lemma C.1** (Properties of multivariate normal). Assume that  $X \sim N_d(\mu, \Sigma)$  and  $Y \sim N_d(\nu, \Psi)$ , where

$$X = \begin{pmatrix} X_a \\ X_b \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \quad \text{and} \quad \Sigma = \begin{pmatrix} \Sigma_a & \Sigma_{ab} \\ \Sigma_{ab}^\top & \Sigma_b \end{pmatrix}.$$

Then,

- **Marginal distributions.**

$$X_a \sim N_{\dim X_a}(\mu_a, \Sigma_a) \quad \text{and} \quad X_b \sim N_{\dim X_b}(\mu_b, \Sigma_b).$$

- **Conditional distributions.**

$$X_a | X_b \sim N_{\dim X_a}(\tilde{\mu}_a, \tilde{\Sigma}_a) \quad \text{and} \quad X_b \sim N_{\dim X_b}(\tilde{\mu}_b, \tilde{\Sigma}_b),$$

where

$$\begin{aligned} \tilde{\mu}_a &= \mu_a + \Sigma_{ab} \Sigma_b^{-1} (X_b - \mu_b) & \tilde{\mu}_b &= \mu_b + \Sigma_{ab}^\top \Sigma_a^{-1} (X_a - \mu_a) \\ \tilde{\Sigma}_a &= \Sigma_a - \Sigma_{ab} \Sigma_b^{-1} \Sigma_{ab}^\top & \tilde{\Sigma}_b &= \Sigma_b - \Sigma_{ab}^\top \Sigma_a^{-1} \Sigma_{ab} \end{aligned}$$

- **Linear combinations.**

$$AX + BY + C \sim N_d(A\mu + B\nu + C, A\Sigma A^\top + B\Psi B^\top)$$

where  $A$  and  $B$  are appropriately sized matrices, and  $C \in \mathbb{R}^d$ .

- **Product of Gaussian densities.**

$$p(X|\mu, \Sigma)p(Y|\nu, \Psi) \propto p(Z|m, S)$$

where  $p(Z)$  is a Gaussian density,  $m = S(\Sigma^{-1}\mu + \Psi^{-1}\nu)$  and  $S = (\Sigma^{-1} + \Psi^{-1})^{-1}$ .

The normalising constant is equal to the density of  $\mu \sim N(\nu, \Sigma + \Psi)$ .

*Proof.* Omitted—see Petersen and Pedersen (2008, §8). ■

Frequently, in Bayesian statistics especially, the following identities will be useful in deriving posterior distributions involving multivariate normals.

**Lemma C.2.** Let  $x, b \in \mathbb{R}^d$  be a vector,  $X, B \in \mathbb{R}^{n \times d}$  a matrix, and  $A \in \mathbb{R}^{d \times d}$  a symmetric, invertible matrix. Then,

$$\begin{aligned} -\frac{1}{2}x^\top Ax + b^\top x &= -\frac{1}{2}(x - A^{-1}b)^\top A(x - A^{-1}b) + \frac{1}{2}b^\top A^{-1}b \\ -\frac{1}{2}\text{tr}(X^\top AX) + \text{tr}(B^\top X) &= -\frac{1}{2}\text{tr}((X - A^{-1}B)^\top A(X - A^{-1}B)) + \frac{1}{2}\text{tr}(B^\top A^{-1}B). \end{aligned}$$

*Proof.* Omitted—see Petersen and Pedersen (2008, §8.1.6). ■

**Lemma C.3.** Let  $X \sim N_p(\mu_\theta, \Sigma_\theta)$ , that is, the mean vector  $\mu_\theta$  and covariance matrix  $\Sigma_\theta$  depends on a real,  $q$ -dimensional vector  $\theta$ . The Fisher information matrix  $U \in \mathbb{R}^{q \times q}$  for  $\theta$  has  $(i, j)$  entries given by

$$U_{ij} = \frac{\partial \mu_\theta^\top}{\partial \theta_i} \Sigma_\theta^{-1} \frac{\partial \mu_\theta}{\partial \theta_j} + \frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_j} \right) \quad (\text{C.1})$$

for  $i, j = 1, \dots, q$ .

*Proof.* Define the derivative of a matrix  $\Sigma \in \mathbb{R}^{p \times p}$  with respect to a scalar  $z$ , denoted  $\partial \Sigma / \partial z \in \mathbb{R}^{p \times p}$ , by  $(\partial \Sigma / \partial z)_{ij} = \partial \Sigma_{ij} / \partial z$ , i.e. derivatives are taken element-wise. The two identities below are useful:

$$\frac{\partial}{\partial z} \text{tr} \Sigma = \text{tr} \frac{\partial \Sigma}{\partial z} \quad (\text{C.2})$$

$$\frac{\partial}{\partial z} \log |\Sigma| = \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial z} \right) \quad (\text{C.3})$$

$$\frac{\partial \Sigma^{-1}}{\partial z} = -\Sigma^{-1} \frac{\partial \Sigma}{\partial z} \Sigma^{-1} \quad (\text{C.4})$$

A useful reference for these identities is Petersen and Pedersen (2008).

Differentiating the log-likelihood for  $\theta$  with respect to the  $i$ 'th component of  $\theta$  yields

$$\begin{aligned} \frac{\partial}{\partial \theta_i} L(\theta|X) &= -\frac{1}{2} \frac{\partial}{\partial \theta_i} \log |\Sigma_\theta| - \frac{1}{2} \frac{\partial}{\partial \theta_i} \text{tr}(\Sigma_\theta^{-1} (X - \mu_\theta)(X - \mu_\theta)^\top) \\ &= -\frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \right) - \frac{1}{2} \text{tr} \left( \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_i} (X - \mu_\theta)(X - \mu_\theta)^\top \right) \\ &\quad - \frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial}{\partial \theta_i} ((X - \mu_\theta)(X - \mu_\theta)^\top) \right) \\ &= \underbrace{-\frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \right)}_{(A)} - \underbrace{\frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} (X - \mu_\theta)(X - \mu_\theta)^\top \right)}_{(B)} \\ &\quad + \underbrace{\text{tr} \left( \Sigma_\theta^{-1} (X - \mu_\theta) \frac{\partial \mu_\theta^\top}{\partial \theta_i} \right)}_{(C)}. \end{aligned}$$

Taking derivatives again, this time with respect to  $\theta_j$ , of the three parts (A), (B) and (C) above, we get:

• (A)

$$\frac{1}{2} \frac{\partial}{\partial \theta_j} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \right) = \frac{1}{2} \text{tr} \left( \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} \frac{\partial \Sigma_\theta}{\partial \theta_i} + \Sigma_\theta^{-1} \frac{\partial^2 \Sigma_\theta}{\partial \theta_i \partial \theta_j} \right)$$

• (B)

$$\begin{aligned} & \frac{1}{2} \frac{\partial}{\partial \theta_j} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} (X - \mu_\theta)(X - \mu_\theta)^\top \right) \\ &= \frac{1}{2} \text{tr} \left( \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} (X - \mu_\theta)(X - \mu_\theta)^\top \right) \\ & \quad + \frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial^2 \Sigma_\theta}{\partial \theta_i \partial \theta_j} \Sigma_\theta^{-1} (X - \mu_\theta)(X - \mu_\theta)^\top \right) \\ & \quad + \frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} (X - \mu_\theta)(X - \mu_\theta)^\top \right) \\ & \quad - \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} \frac{\partial \mu_\theta}{\partial \theta_j} (X - \mu_\theta)^\top \right) \end{aligned}$$

• (C)

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \text{tr} \left( \Sigma_\theta^{-1} (X - \mu_\theta) \frac{\partial \mu_\theta^\top}{\partial \theta_i} \right) &= \text{tr} \left( \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} (X - \mu_\theta) \frac{\partial \mu_\theta^\top}{\partial \theta_i} - \Sigma_\theta^{-1} \frac{\partial \mu_\theta}{\partial \theta_j} \frac{\partial \mu_\theta^\top}{\partial \theta_i} \right. \\ & \quad \left. - \Sigma_\theta^{-1} (X - \mu_\theta) \frac{\partial^2 \mu_\theta}{\partial \theta_i \partial \theta_j} \right) \end{aligned}$$

The Fisher information matrix  $U$  contains  $(i, j)$  entries equal to the expectation of  $-\frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\theta|X)$ . Using the fact that 1)  $E[X - \mu_\theta] = 0$ ; 2)  $E[\text{tr} \Sigma] = \text{tr}(E \Sigma)$ ; 3)  $E[XX^\top] = \Sigma_\theta$ ; and 4) the trace is invariant under cyclic permutations, we get

$$\begin{aligned} U_{ij} &= \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \mu_\theta}{\partial \theta_j} \frac{\partial \mu_\theta^\top}{\partial \theta_i} \right) \\ & \quad + \frac{1}{2} \text{tr} \left( \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} \frac{\partial \Sigma_\theta}{\partial \theta_i} + \cancel{\Sigma_\theta^{-1} \frac{\partial^2 \Sigma_\theta}{\partial \theta_i \partial \theta_j}} - \cancel{\frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} \frac{\partial \Sigma_\theta}{\partial \theta_i}} - \cancel{\Sigma_\theta^{-1} \frac{\partial^2 \Sigma_\theta}{\partial \theta_i \partial \theta_j}} - \frac{\partial \Sigma_\theta}{\partial \theta_i} \frac{\partial \Sigma_\theta^{-1}}{\partial \theta_j} \right) \\ &= \frac{\partial \mu_\theta^\top}{\partial \theta_i} \Sigma_\theta^{-1} \frac{\partial \mu_\theta}{\partial \theta_j} + \frac{1}{2} \text{tr} \left( \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_i} \Sigma_\theta^{-1} \frac{\partial \Sigma_\theta}{\partial \theta_j} \right) \end{aligned}$$

as required. ■

apx:matrixn  
ormal

## C.2 Matrix normal distribution

The matrix normal distribution is an extension of the Gaussian distribution to matrices. Let  $X \in \mathbb{R}^{n \times m}$  matrix, and let  $X$  follow a matrix normal distribution with mean  $\mu \in \mathbb{R}^{n \times m}$  and row and column variances  $\Sigma \in \mathbb{R}^{n \times n}$  and  $\Psi \in \mathbb{R}^{m \times m}$  respectively, which we denote by  $X \sim \text{MN}_{n,m}(\mu, \Sigma, \Psi)$ . Then,

- **Pdf.**  $p(X|\mu, \Sigma, \Psi) = (2\pi)^{-nm/2} |\Sigma|^{-m/2} |\Psi|^{-n/2} e^{-\frac{1}{2} \text{tr}(\Psi^{-1}(X-\mu)^\top \Sigma^{-1}(X-\mu))}$ .
- **Moments.**  $\mathbb{E} X = \mu$ ,  $\text{Var}(X_{i.}) = \Psi$  for  $i = 1, \dots, n$ , and  $\text{Var}(X_{.j}) = \Sigma$  for  $j = 1, \dots, m$ .
- **Entropy.**  $H(p) = \frac{1}{2} \log |2\pi e(\Psi \otimes \Sigma)| = \frac{nm}{2} (1 + \log 2\pi) + \frac{1}{2} \log |\Sigma|^m |\Psi|^n$ .

**Lemma C.4** (Equivalence between matrix and multivariate normal).  $X \sim \text{MN}_{n,m}(\mu, \Sigma, \Psi)$  if and only if  $\text{vec } X \sim \text{N}_{nm}(\text{vec } \mu, \Psi \otimes \Sigma)$ .

*Proof.* In the exponent of the matrix normal pdf, we have

$$\begin{aligned} -\frac{1}{2} \text{tr}(\Psi^{-1}(X - \mu)^\top \Sigma^{-1}(X - \mu)) \\ &= -\frac{1}{2} \text{vec}(X - \mu)^\top \text{vec}(\Sigma^{-1}(X - \mu)\Psi^{-1}) \\ &= -\frac{1}{2} \text{vec}(X - \mu)^\top (\Psi^{-1} \otimes \Sigma^{-1}) \text{vec}(X - \mu) \\ &= -\frac{1}{2} (\text{vec } X - \text{vec } \mu)^\top (\Psi \otimes \Sigma)^{-1} (\text{vec } X - \text{vec } \mu). \end{aligned}$$

Also,  $|\Sigma|^{-m/2} |\Psi|^{-n/2} = |\Psi \otimes \Sigma|^{-1/2}$ . This converts the matrix normal pdf to that of a multivariate normal pdf. ■

Some useful properties of the matrix normal distribution are listed:

- **Expected values.**

$$\begin{aligned} \mathbb{E}(X - \mu)(X - \mu)^\top &= \text{tr}(\Psi)\Sigma \in \mathbb{R}^{n \times n} \\ \mathbb{E}(X - \mu)^\top (X - \mu) &= \text{tr}(\Sigma)\Psi \in \mathbb{R}^{m \times m} \\ \mathbb{E}[XAX^\top] &= \text{tr}(A^\top \Psi)\Sigma + \mu A \mu^\top \\ \mathbb{E}[X^\top BX] &= \text{tr}(\Sigma B^\top)\Psi + \mu^\top B \mu \\ \mathbb{E}[XCX] &= \Sigma C^\top \Psi + \mu C \mu \end{aligned}$$

- **Transpose.**  $X^\top \sim \text{MN}_{m,n}(\mu^\top, \Psi, \Sigma)$ .

- **Linear transformation.** Let  $A \in \mathbb{R}^{a \times n}$  be of full-rank  $a \leq n$  and  $B \in \mathbb{R}^{m \times b}$  be of full-rank  $b \leq m$ . Then  $AXB \sim \text{MN}_{a,b}(\mu^\top, A\Sigma A^\top, B^\top \Psi B)$ .
- **Iid.** If  $X_i \stackrel{\text{iid}}{\sim} N_m(\mu, \Psi)$  for  $i = 1, \dots, n$ , and we arranged these vectors row-wise into the matrix  $X = (X_1^\top, \dots, X_n^\top)^\top \in \mathbb{R}^{n \times m}$ , then  $X \sim \text{MN}(1_n \mu^\top, I_n, \Psi)$ .

### C.3 Truncated univariate normal distribution

apx:truncun  
inorm

Let  $X \sim N(\mu, \sigma^2)$  with the random variable  $X$  restricted to the interval  $(a, b) \subset \mathbb{R}$ . Then we say that  $X$  follows a truncated normal distribution, and we denote this by  $X \sim {}^tN(\mu, \sigma^2, a, b)$ . Let  $\alpha = (a - \mu)/\sigma$ ,  $\beta = (b - \mu)/\sigma$ , and  $C = \Phi(\beta) - \Phi(\alpha)$ . Then,

- **Pdf.**  $p(X|\mu, \sigma, a, b) = C^{-1}(2\pi\sigma^2)^{-1/2}e^{-\frac{1}{2\sigma^2}(X-\mu)^2} = \sigma C^{-1}\phi(\frac{X-\mu}{\sigma})$ .
- **Moments.**

$$\begin{aligned} \mathbb{E} X &= \mu + \sigma \frac{\phi(\alpha) - \phi(\beta)}{C} \\ \mathbb{E} X^2 &= \sigma^2 + \mu^2 + \sigma^2 \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{C} + 2\mu\sigma \frac{\phi(\alpha) - \phi(\beta)}{C} \\ \text{Var } X &= \sigma^2 \left[ 1 + \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{C} - \left( \frac{\phi(\alpha) - \phi(\beta)}{C} \right)^2 \right] \end{aligned}$$

- **Entropy.**

$$\begin{aligned} H(p) &= \frac{1}{2} \log 2\pi e\sigma^2 + \log C + \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{2C} \\ &= \frac{1}{2} \log 2\pi e\sigma^2 + \log C + \frac{1}{2\sigma^2} \cdot \overbrace{\sigma^2 \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{C}}^{\text{Var } X - \sigma^2 + (\mathbb{E} X - \mu)^2} \\ &= \frac{1}{2} \log 2\pi\sigma^2 + \log C + \frac{1}{2\sigma^2} \mathbb{E}[X - \mu]^2 \end{aligned}$$

$$\text{because } \text{Var } X + (\mathbb{E} X - \mu)^2 = \mathbb{E} X^2 - (\mathbb{E} X)^2 + (\mathbb{E} X)^2 + \mu^2 - 2\mu \mathbb{E} X.$$

For binary probit models, the distributions that come up are one-sided truncations at zero, i.e.  ${}^tN(\mu, \sigma^2, 0, +\infty)$  (upper tail/positive part) and  ${}^tN(\mu, \sigma^2, -\infty, 0)$  (lower tail/negative part), for which their moments are of interest. As an aside, if  $\mu = 0$  then the truncation  ${}^tN(0, \sigma^2, 0, +\infty) \equiv N_+(0, \sigma^2)$  is called the *folded-normal* distribution. For the positive one-sided truncation at zero,  $C = \Phi(+\infty) - \Phi(-\mu/\sigma) = 1 - \Phi(-\mu/\sigma) = \Phi(\mu/\sigma)$ , and for the negative one-sided truncation at zero,  $C = \Phi(-\mu/\sigma) - \Phi(-\infty) = 1 - \Phi(\mu/\sigma)$ . Additionally, if  $\sigma = 1$ , then  ${}^tN(0, 1, 0, +\infty) \equiv N_+(0, 1)$  is called the *half-normal* distribution.



One may simulate random draws from a truncated normal distribution by drawing from  $N(\mu, \sigma^2)$  and discarding samples that fall outside  $(a, b)$ . Alternatively, the inverse-transform method using

$$X = \mu + \sigma \Phi^{-1}(\Phi(\alpha) + UC)$$

with  $U \sim \text{Unif}(0, 1)$  will work too. Either of these methods will work reasonably well as long as the truncation region is not too far away from  $\mu$ , but neither is particularly fast. Efficient algorithms have been explored which are along the lines of either accept/reject algorithms (Robert, 1995), Gibbs sampling (Damien and Walker, 2001), or pseudo-random number generation algorithms (Chopin, 2011). The latter algorithm is inspired by the Ziggurat algorithm (Marsaglia and Tsang, 2000) which is considered to be the fastest Gaussian random number generator.

## C.4 Truncated multivariate normal distribution

apx:truncmu  
ltinorm

Consider the restriction of  $X \sim N_d(\mu, \Sigma)$  to a convex subset<sup>1</sup>  $\mathcal{A} \subset \mathbb{R}^d$ . Call this distribution the truncated multivariate normal distribution, and denote it  $X \sim {}^tN_d(\mu, \Sigma, \mathcal{A})$ . The pdf is  $p(X|\mu, \Sigma, \mathcal{A}) = C^{-1}\phi(X|\mu, \Sigma)\mathbb{1}[X \in \mathcal{A}]$ , where

$$C = \int_{\mathcal{A}} \phi(x|\mu, \Sigma) dx = P(X \in \mathcal{A}).$$

Generally speaking, there are no closed-form expressions for  $Eg(X)$  for any well-defined functions  $g$  on  $X$ . One strategy to obtain values such as  $E X$  (mean),  $E X^2$  (second moment) and  $E \log p(X)$  (entropy) would be Monte Carlo integration. If  $X^{(1)}, \dots, X^{(T)}$  are samples from  $X \sim {}^tN_d(\mu, \Sigma, \mathcal{A})$ , then  $\widehat{Eg(X)} = \frac{1}{T} \sum_{i=1}^T g(X^{(i)})$ .

Sampling from a truncated multivariate normal distribution is described by Robert (1995), who used a Gibbs-based approach, which we now describe. Assume that the one-dimensional slices of  $\mathcal{A}$

$$\mathcal{A}_k(X_{-j}) = \{X_j | (X_1, \dots, X_{j-1}, X_j, X_{j+1}, \dots, X_d) \in \mathcal{A}\}$$

are readily available so that the bounds or anti-truncation region of  $X_j$  given the rest of the components  $X_{-j}$  are known to be  $(x_j^-, x_j^+)$ . Using properties of the normal

<sup>1</sup>A convex subset is a subset of a space that is closed under convex combinations. In Euclidean space, for every pair of points in a convex set, all the points that lie on the straight line segment which joins the pair of points are also in the set.

distribution, the full conditionals of  $X_j$  given  $X_{-j}$  is

$$\begin{aligned} X_j|X_{-j} &\sim {}^t\text{N}(\tilde{\mu}_j, \tilde{\sigma}_j^2, x_j^-, x_j^+) \\ \tilde{\mu}_j &= \mu_j + \Sigma_{j,-j}^\top \Sigma_{-j,-j} (x_{-j} - \mu_{-j}) \\ \tilde{\sigma}_j^2 &= \Sigma_{11} - \Sigma_{j,-j}^\top \Sigma_{-j,-j} \Sigma_{j,-j}. \end{aligned}$$

According to Robert (1995), if  $\Psi = \Sigma^{-1}$ , then

$$\Sigma_{-j,-j}^{-1} = \Psi_{-j,-j} - \Psi_{j,-j} \Psi_{-j,-j}^\top / \Psi_{jj}$$

which means that we need only compute one global inverse  $\Sigma^{-1}$ . Therefore, the Gibbs sampler makes draws from truncated normal distributions in the following sequence, given initial values  $X^{(0)}$ :

- Draw  $X_1^{(t)}|X_2^{(t)}, \dots, X_d^{(t)} \sim {}^t\text{N}(\tilde{\mu}_1, \tilde{\sigma}_1^2, x_1^-, x_1^+)$ .
- Draw  $X_2^{(t)}|X_1^{(t+1)}, X_3^{(t)}, \dots, X_d^{(t)} \sim {}^t\text{N}(\tilde{\mu}_2, \tilde{\sigma}_2^2, x_2^-, x_2^+)$ .
- ...
- Draw  $X_d^{(t)}|X_1^{(t+1)}, \dots, X_{d-1}^{(t+1)} \sim {}^t\text{N}(\tilde{\mu}_d, \tilde{\sigma}_d^2, x_d^-, x_d^+)$ .

In a later work, Damien and Walker (2001) introduce a latent variable  $Y \in \mathbb{R}$  such that the joint pdf of  $X$  and  $Y$  is

$$p(X_1, \dots, X_d, Y) \propto \exp(-Y/2) \mathbb{1}(Y > (X - \mu)^\top \Sigma^{-1}(X - \mu)) \mathbb{1}(X \in \mathcal{A}).$$

Now, the Gibbs conditional densities for the  $X_k$ 's are given by

$$p(X_j|X_{-j}, Y) \propto \mathbb{1}(X_j \in \mathcal{B}_j)$$

where

$$\mathcal{B}_j \in (x_j^-, x_j^+) \cap \{X_j | (X - \mu)^\top \Sigma^{-1}(X - \mu) < Y\}.$$

Thus, given values for  $X_{-j}$  and  $Y$ , the bounds for  $X_j$  involves solving a quadratic equation in  $X_j$ . The Gibbs conditional density for  $Y|X$  is a shifted exponential distribution, which can be sampled using the inverse-transform method. Thus, both  $X$  and  $Y$  can be sampled directly from uniform variates.

For probit models, we are interested in the conical truncations  $\mathcal{C}_j = \{X_j > X_k | k \neq j, \text{ and } k = 1, \dots, m\}$  for which the  $j$ 'th component of  $X$  is largest. These truncations form cones in  $d$ -dimensional space such that  $\mathcal{C}_1 \cup \dots \cup \mathcal{C}_d = \mathbb{R}^d$ , and hence the name.

In the case where  $\Sigma$  is a diagonal matrix, the conically truncated multivariate normal distributions are easier to deal with due to the independence structure in the covariance matrix. In particular, most calculations of interest involve only a one dimensional inte-

gral of products of normal cdfs. We present some results that we have not previously seen before elsewhere.

thm:contrun  
cn

**Lemma C.5.** Let  $X \sim {}^t\text{N}_d(\mu, \Sigma, \mathcal{C}_j)$ , with  $\mu = (\mu_1, \dots, \mu_d)^\top$  and  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ , and  $\mathcal{C}_j = \{X_j > X_k | k \neq j, \text{ and } k = 1, \dots, m\}$  a conical truncation of  $\mathbb{R}^d$  such that the  $j$ 'th component is largest. Then,

(i) **Pdf.** The pdf of  $X$  has the following functional form:

$$p(X) = \frac{C^{-1}}{\sigma_1 \dots \sigma_d (2\pi)^{d/2}} \exp \left[ -\frac{1}{2} \sum_{i=1}^d \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

where  $\phi$  is the pdf of a standard normal distribution and

$$C = \mathbb{E}_Z \left[ \prod_{\substack{i=1 \\ i \neq j}}^d \Phi \left( \frac{\sigma_j}{\sigma_i} Z + \frac{\mu_j - \mu_i}{\sigma_i} \right) \right]$$

where  $Z \sim \text{N}(0, 1)$ .

(ii) **Moments.** The expectation  $\mathbb{E} X = (\mathbb{E} X_1, \dots, \mathbb{E} X_d)^\top$  is given by

$$\mathbb{E} X_i = \begin{cases} \mu_i - \sigma_i C^{-1} \mathbb{E}_Z \left[ \phi_i \prod_{k \neq i, j} \Phi_k \right] & \text{if } i \neq j \\ \mu_j - \sigma_j \sum_{i \neq j} (\mathbb{E} X_i - \mu_i) & \text{if } i = j \end{cases}$$

and the second moments  $\mathbb{E}[X - \mu]^2$  are given by

$$\mathbb{E}[X_i - \mu_i]^2 = \begin{cases} \sigma_i^2 + (\mu_j - \mu_i)(\mathbb{E} X_i - \mu_i) + \sigma_i \sigma_j C^{-1} \mathbb{E}_Z \left[ Z \phi_i \prod_{k \neq i, j} \Phi_k \right] & \text{if } i \neq j \\ C^{-1} \sigma_j^2 \mathbb{E}_Z \left[ Z^2 \prod_{k \neq j} \Phi_k \right] & \text{if } i = j \end{cases}$$

where we had defined

$$\phi_i = \phi_i(Z) = \phi \left( \frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i} \right), \text{ and} \\ \Phi_i = \Phi_i(Z) = \Phi \left( \frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i} \right).$$

(iii) **Entropy.** The entropy is given by

$$H(p) = \log C + \frac{d}{2} \log 2\pi + \frac{1}{2} \sum_{i=1}^d \log \sigma_i^2 + \frac{1}{2} \sum_{i=1}^d \frac{1}{\sigma_i^2} \mathbb{E}[x_i - \mu_i]^2.$$

*Proof.* See [Appendix D](#) for the proof. ■

## C.5 Gamma distribution

apx:gamma  
st

For  $X \in \mathbb{R}_{\geq 0}$ , let  $X$  be distributed according to the gamma distribution with shape  $s$  and rate  $r$ , denoted  $X \sim \Gamma(s, r)$ . Then,

- **Pdf.**  $p(X) = \Gamma(s)^{-1} r^s X^{s-1} e^{-rX}$ .
- **Moments.**  $E X = s/r$ ,  $\text{Var } X = s/r^2$ .
- **Entropy.**  $H(p) = s - \log r + \log \Gamma(s) + (1-s)\psi(s)$ .

In the above,  $\Gamma(\cdot) = \Gamma_1(\cdot)$  and  $\psi(\cdot) = \psi_1(\cdot)$  are the gamma and digamma functions, defined by

$$\Gamma(a) = \begin{cases} (a-1)! & \text{if } a \in \mathbb{Z}^+ \\ \int_0^\infty u^{a-1} e^{-u} \mathrm{d}u & \text{otherwise} \end{cases}$$

and

$$\psi(a) = \frac{\partial}{\partial a} \log \Gamma(a) = \frac{\partial \Gamma(a)/\partial a}{\Gamma(a)}.$$

Often, the gamma distribution is parameterised according to shape  $s$  and scale  $\sigma = 1/r$  parameters,  $X \sim \Gamma(s, \sigma)$ .

## C.6 Inverse gamma distribution

def:invgam

For  $X \in \mathbb{R}_{\geq 0}$ , a random variable  $X$  distributed according to an inverse gamma distribution with parameters  $s$  (shape) and  $\sigma$  (scale) is denoted by  $X \sim \Gamma^{-1}(s, \sigma)$ . Then,

- **Pdf.**  $p(X) = \Gamma(s)^{-1} \sigma^s X^{-(s+1)} e^{-\sigma/X}$ .
- **Moments.**  $E X = \sigma/(s-1)$ ,  $\text{Var } X = \sigma^2((s-1)^2(s-2))^{-1}$ .
- **Entropy.**  $H(p) = s + \log(\sigma \Gamma(s)) - (1+s)\psi(s)$ .

with  $\Gamma(\cdot)$  and  $\psi(\cdot)$  representing the gamma and digamma functions respectively, as defined in [Appendix C.5](#).

**Lemma C.6.** *If  $X$  has a Gamma distribution with shape and rate  $s$  and  $r$ , then  $1/X \sim \Gamma^{-1}(s, r)$ .*

*Proof.* Let  $Y = 1/X$ . Then the pdf of  $Y$  is

$$\begin{aligned} p_Y(Y) &= p_X(1/Y) \left| \frac{\partial}{\partial Y}(1/Y) \right| \\ &= \Gamma(s)^{-1} r^s (1/Y)^{s-1} e^{-r/Y} (1/Y^2) \\ &= \Gamma(s)^{-1} r^s Y^{-(s+1)} e^{-r/Y} \end{aligned}$$

which is the pdf of an inverse gamma with shape  $s$  and scale  $r$ . ■

## Appendix D

# Proofs related to the conically truncated independent multivariate normal distribution

apx:contrun  
proof

We present the proof for [Lemma C.5](#) related to the conically truncated independent multivariate normal distribution, which we had not encountered in the literature.

### D.1 Proof of [Lemma C.5](#): Pdf

Using the fact that  $\int p(x) dx = 1$ , and that

$$\begin{aligned}
 & \int \cdots \int [x_i < x_j, \forall i \neq j] \cdot \prod_{i=1}^d \phi(x_i | \mu_i, \sigma_i^2) dx_1 \cdots dx_d \\
 &= \int \cdots \int \mathbb{1}[x_i < x_j, \forall i \neq j] \prod_{i=1}^d \left[ \frac{1}{\sigma_i} \phi\left(\frac{x_i - \mu_i}{\sigma_i}\right) \right] dx_1 \cdots dx_d \\
 &= \int \cdots \int \mathbb{1}[x_i < x_j, \forall i \neq j] \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) \prod_{\substack{i=1 \\ i \neq j}}^d \left[ \frac{1}{\sigma_i} \phi\left(\frac{x_i - \mu_i}{\sigma_i}\right) \right] dx_1 \cdots dx_d \\
 &= \int \prod_{\substack{i=1 \\ i \neq j}}^d \Phi\left(\frac{x_j - \mu_i}{\sigma_i}\right) \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) dx_j \\
 &= \int \prod_{\substack{i=1 \\ i \neq j}}^d \Phi\left(\frac{\sigma_j z + \mu_j - \mu_i}{\sigma_i}\right) \phi(z) dz \\
 & \quad \text{(by using the standardisation } z = (x_j - \mu_j)/\sigma_j)
 \end{aligned}$$

$$= \mathbb{E}_Z \left[ \prod_{\substack{i=1 \\ i \neq j}}^d \Phi \left( \frac{\sigma_j}{\sigma_i} Z + \frac{\mu_j - \mu_i}{\sigma_i} \right) \right]$$

the proof follows directly.

## D.2 Proof of Lemma C.5: Moments

Recall that for  $Y \sim {}^t\mathcal{N}(\mu, \sigma^2, -\infty, b)$ , for some function  $g$  of  $Y$ , we have that

$$\mathbb{E} g(Y) = \Phi(\beta)^{-1} \int [y < b] \cdot g(y) \phi(y|\mu, \sigma^2) dy,$$

and in particular, we have

$$\mathbb{E}[Y - \mu] = -\sigma \frac{\phi(\beta)}{\Phi(\beta)} \quad (\text{D.1})$$

$$\mathbb{E}[Y - \mu]^2 - \sigma^2 = -\sigma^2 \frac{\beta \phi(\beta)}{\Phi(\beta)} \quad (\text{D.2})$$

where  $\beta = (b - \mu)/\sigma$ . For the conically truncated multivariate normal distribution  $X \sim {}^t\mathcal{N}_d(\mu, \Sigma, \mathcal{A}_j)$ , where  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ , the independence structure of  $\Sigma$  makes it possible to consider the expectations of each of the components separately by marginalising out the rest of the components. For simplicity, denote  $p(x_k) = \phi(x_k|\mu_k, \sigma_k) = \sigma_k^{-1} \phi(\frac{x_k - \mu_k}{\sigma_k})$ . For  $i \neq j$ , we have

$$\begin{aligned} \mathbb{E} g(X_i) &= C^{-1} \int \cdots \int [x_k < x_j, \forall k \neq j] \cdot g(x_i) \prod_{k=1}^d p(x_k) dx_1 \cdots dx_d \\ &= C^{-1} \frac{\Phi((x_j - \mu_j)/\sigma_j)}{\Phi((x_j - \mu_j)/\sigma_j)} \iint [x_i < x_j] \cdot g(x_i) p(x_i) p(x_j) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) dx_i dx_j \\ &= C^{-1} \int \mathbb{E}_{X_i \sim {}^t\mathcal{N}(\mu_i, \sigma_i^2, -\infty, x_j)} [g(X_i)] \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) p(x_j) dx_j \end{aligned} \quad (\text{D.3})$$

where  $C$  is the normalising constant for  $X$ , while for the  $j$ 'th component we have

$$\begin{aligned} \mathbb{E} g(X_j) &= C^{-1} \int \cdots \int [x_k < x_j, \forall k \neq j] \cdot g(x_j) \prod_{k=1}^d p(x_k) dx_1 \cdots dx_d \\ &= C^{-1} \int g(x_j) \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) p(x_j) dx_d. \end{aligned} \quad (\text{D.4})$$

Plugging in (D.1) for  $g(X_i) = X_i - \mu_i$  in (D.3) we get

$$\begin{aligned}
\mathbb{E} X_i - \mu_i &= -C^{-1} \int \left( \sigma_i \phi \left( \frac{x_j - \mu_i}{\sigma_i} \right) / \Phi \left( \frac{x_j - \mu_i}{\sigma_i} \right) \right) \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) p(x_j) dx_j \\
&= -\sigma_i C^{-1} \int \phi \left( \frac{x_j - \mu_i}{\sigma_i} \right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) p(x_j) dx_j \\
&= -\sigma_i C^{-1} \int \phi \left( \frac{\sigma_j z + \mu_j - \mu_i}{\sigma_i} \right) \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{\sigma_j z + \mu_j - \mu_k}{\sigma_k} \right) \phi(z) dz \\
&= -\sigma_i C^{-1} \mathbb{E}_Z \left[ \phi \left( \frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i} \right) \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{\sigma_j Z + \mu_j - \mu_k}{\sigma_k} \right) \right]
\end{aligned}$$

where  $Z$  is the distribution of  $N(0, 1)$ , and we had used a change of variable  $x_j = \sigma_j z + \mu_j$ , so that  $p(x_j) = \sigma_j^{-1} \phi(z)$  and  $dx_j = \sigma_j dz$ . For the  $j$ 'th component, substitute  $g(x_j) = x_j - \mu_j$  in (D.4) to get

$$\begin{aligned}
\mathbb{E} X_j - \mu_j &= C^{-1} \int (x_j - \mu_j) \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) p(x_j) dx_j \\
&= C^{-1} \sigma_j \int z \prod_{\substack{k=1 \\ k \neq j}}^d \Phi \left( \frac{\sigma_j z + \mu_j - \mu_k}{\sigma_k} \right) \phi(z) dz \\
&= \sigma_j \sum_{\substack{i=1 \\ i \neq j}}^d \sigma_i C^{-1} \mathbb{E} \left[ \phi \left( \frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i} \right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi \left( \frac{\sigma_j Z + \mu_j - \mu_k}{\sigma_k} \right) \right] \\
&= -\sigma_j \sum_{\substack{i=1 \\ i \neq j}}^d (\mathbb{E} X_i - \mu_i),
\end{aligned}$$

where we have made use of Lemma D.1 in the second last step.

For the second moments, plug in (D.2) for  $g(X_i) = (X_i - \mu_i)^2 - \sigma_i^2$  in (D.3) to get

$$\begin{aligned}
 \mathbb{E}[X_i - \mu_i]^2 - \sigma_i^2 &= -\sigma_i^2 C^{-1} \int \overbrace{\frac{x_j - \mu_i}{\sigma_i}}^{x_j - \mu_i - \mu_j + \mu_j} \cdot \frac{\phi((x_j - \mu_i)/\sigma_i)}{\Phi((x_j - \mu_i)/\sigma_i)} \prod_{\substack{k=1 \\ k \neq j}}^d \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) p(x_j) dx_j \\
 &= -\sigma_i C^{-1} \int (x_j - \mu_j) \phi\left(\frac{x_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) p(x_j) dx_j \\
 &\quad + \overbrace{(\mu_j - \mu_i) \cdot -\sigma_i C^{-1} \int \phi\left(\frac{x_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) p(x_j) dx_j}^{\mathbb{E} X_i - \mu_i} \\
 &= (\mu_j - \mu_i)(\mathbb{E} X_i - \mu_i) \\
 &\quad + \sigma_i C^{-1} \int \sigma_j z \phi\left(\frac{x_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{\sigma_j z + \mu_j - \mu_k}{\sigma_k}\right) \phi(z) dz \\
 &= (\mu_j - \mu_i)(\mathbb{E} X_i - \mu_i) \\
 &\quad + \sigma_i \sigma_j C^{-1} \mathbb{E} \left[ Z \phi\left(\frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{\sigma_j Z + \mu_j - \mu_k}{\sigma_k}\right) \right]
 \end{aligned}$$

And similarly, for the  $j$ 'th component

$$\begin{aligned}
 \mathbb{E}[X_j - \mu_j]^2 &= C^{-1} \int (x_j - \mu_j)^2 \prod_{\substack{k=1 \\ k \neq j}}^d \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) p(x_j) dx_j \\
 &= C^{-1} \sigma_j^2 \int z^2 \prod_{\substack{k=1 \\ k \neq j}}^d \Phi\left(\frac{z \sigma_j + \mu_j - \mu_k}{\sigma_k}\right) p(x_j) dz \\
 &= C^{-1} \sigma_j^2 \mathbb{E}_Z \left[ Z^2 \prod_{\substack{k=1 \\ k \neq j}}^d \Phi\left(\frac{Z \sigma_j + \mu_j - \mu_k}{\sigma_k}\right) \right].
 \end{aligned}$$

Lastly, we use the following result in the derivation above.

**Lemma D.1.** *Let  $Z \sim \mathcal{N}(0, 1)$ . Then for all  $m \in \{\mathbb{N} \mid m > 1\}$  and  $(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+$ ,*

$$\mathbb{E} \left[ Z \prod_{\substack{k=1 \\ k \neq j}}^m \Phi(\sigma_k Z + \mu_k) \right] = \sum_{\substack{i=1 \\ i \neq j}}^m \mathbb{E} \left[ \sigma_i \phi(\sigma_i Z + \mu_i) \prod_{\substack{k=1 \\ k \neq i, j}}^m \Phi(\sigma_k Z + \mu_k) \right]$$

for some  $j \in \{1, \dots, m\}$ .



*Proof.* Use the fact that for any differentiable function  $g$ ,  $E[Zg(Z)] = E[g'(Z)]$ , and apply the result with the function  $g_m : z \mapsto \prod_{k \neq j} \Phi(\sigma_k z + \mu_k)$ . All that is left is to derive the derivative of  $g$ , and we use an inductive proof to do this. Introduce the following notation for convenience:

$$\begin{aligned}\phi_i &= \phi(\sigma_i z + \mu_i) \\ \Phi_i &= \Phi(\sigma_i z + \mu_i)\end{aligned}$$

The simplest case is when  $m = 2$ , which can be trivially shown to be true. Without loss of generality, let  $j = 1$ . Then

$$\begin{aligned}g_2(z) &= \Phi_2 \\ \Rightarrow \dot{g}_2(z) &= \sigma_2 \phi_2 = \sum_{\substack{i=1 \\ i \neq 1}}^2 \left[ \sigma_i \phi_i \sum_{\substack{k=1 \\ k \neq 1,2}}^2 \Phi_k \right].\end{aligned}$$

Now assume that the inductive hypothesis holds for some  $m \in \{\mathbb{N} \mid m > 1\}$ . That is, the derivative of  $g_m(z) = \prod_{k \neq j} \Phi_k$ ,

$$\dot{g}_m(z) = \sum_{\substack{i=1 \\ i \neq j}}^m \left[ \sigma_i \phi_i \prod_{\substack{k=1 \\ k \neq i,j}}^m \Phi_k \right],$$

is assumed to be true. Also assume that, without loss of generality,  $j \neq m + 1$ . Then, the derivative of

$$g_{m+1}(z) = \prod_{\substack{k=1 \\ k \neq j}}^{m+1} \Phi_k = g_m(z) \Phi_{m+1}$$

is found to be

$$\begin{aligned}\dot{g}_{m+1}(z) &= \sigma_{m+1} \phi_{m+1} g_m(z) + \dot{g}_m(z) \Phi_{m+1} \\ &= \sigma_{m+1} \phi_{m+1} \prod_{\substack{k=1 \\ k \neq j}}^m \Phi_k + \sum_{\substack{i=1 \\ i \neq j}}^m \left[ \sigma_i \phi_i \prod_{\substack{k=1 \\ k \neq i,j}}^m \Phi_k \right] \Phi_{m+1} \\ &= \sigma_{m+1} \phi_{m+1} \prod_{\substack{k=1 \\ k \neq j, m+1}}^{m+1} \Phi_k + \sum_{\substack{i=1 \\ i \neq j}}^m \left[ \sigma_i \phi_i \prod_{\substack{k=1 \\ k \neq i,j}}^{m+1} \Phi_k \right] \\ &= \sum_{\substack{i=1 \\ i \neq j}}^{m+1} \left[ \sigma_i \phi_i \prod_{\substack{k=1 \\ k \neq i,j}}^{m+1} \Phi_k \right],\end{aligned}$$

as required for the inductive proof. Using linearity of expectations, the proof is complete. ■

### D.3 Proof of Lemma C.5: Entropy

As a direct consequence of the definition of entropy,

$$\begin{aligned}
 H(p) &= -\mathbb{E}[\log p(X)] \\
 &= -\mathbb{E} \left[ -\log C - \frac{d}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^d \log \sigma_i^2 - \frac{1}{2} \sum_{i=1}^d \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2 \right] \\
 &= \log C + \frac{d}{2} \log 2\pi + \frac{1}{2} \sum_{i=1}^d \log \sigma_i^2 + \frac{1}{2} \sum_{i=1}^d \frac{1}{\sigma_i^2} \mathbb{E}[x_i - \mu_i]^2.
 \end{aligned}$$

## Appendix E

# I-prior interpretation of the $g$ -prior

misc:gprior

The I-prior for  $\beta$  in a standard linear model resembles the objective  $g$ -prior (Zellner, 1986) for regression coefficients,

$$\beta \sim N_p(\mathbf{0}, g(\mathbf{X}^\top \Psi \mathbf{X})^{-1}),$$

although they are quite different objects. The  $g$ -prior for  $\beta$  has the *inverse* (scaled) Fisher information matrix as its covariance matrix. This, in itself, has a much different and arguably counterintuitive meaning: large amounts of Fisher information about  $\beta$  corresponds to a small prior variance, and hence less deviation away from the prior mean of zero in estimating  $\beta$ . The choice of the hyperparameter  $g$  has been the subject of much debate, with choices ranging from fixing  $g = n$  (corresponding to the concept of *unit Fisher information*), to fully Bayesian and empirical Bayesian methods of estimating  $g$  from the data.

On the other hand, we note that the  $g$ -prior has an I-prior interpretation when argued as follows. Assume that the regression function  $f$  lies in the continual dual space of  $\mathbb{R}^p$  equipped with the inner product  $\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}} = \mathbf{x}^\top (\mathbf{X}^\top \Psi \mathbf{X})^{-1} \mathbf{x}$ . With this inner product and from (3.3) (p. 11), the Fisher information for  $\beta$  is

$$\begin{aligned} \mathcal{I}_g(\beta) &= \sum_{i=1}^n \sum_{j=1}^n \psi_{ij} (\mathbf{X}^\top \Psi \mathbf{X})^{-1} \mathbf{x}_i \otimes (\mathbf{X}^\top \Psi \mathbf{X})^{-1} \mathbf{x}_j \\ &= (\mathbf{X}^\top \Psi \mathbf{X})^{-1} \cancel{(\mathbf{X}^\top \Psi \mathbf{X})} (\mathbf{X}^\top \Psi \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \Psi \mathbf{X})^{-1}, \end{aligned}$$

and this, rather than the usual  $\mathbf{X}^\top \Psi \mathbf{X}$  as the prior covariance matrix for  $\beta$ , means that the I-prior is in fact the standard  $g$ -prior.

The metric induced by the inner product is actually the *Mahalanobis distance*, a scale-invariant natural distance if the covariates are measured on different scales. To expand on this idea, circle back to the regression function and write it as  $f(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle_{\mathcal{X}}$ . In usual least squares regression, the choice of inner product is irrelevant, so the usual dot product is commonly used (however, as we have seen above, the choice of inner product determines the form of the Fisher information for  $\boldsymbol{\beta}$ ). In particular, suppose that all the  $x_{ik}$ 's,  $k = 1, \dots, p$  for each unit  $i = 1, \dots, n$  are measured on the same scale; for instance, these could be measurements in centimetres. In this case, the dot product is reasonable, because  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \sum_{k=1}^p x_{ik}x_{jk}$  and the inner product has a coherent unit, namely the squared unit of the  $x_{ik}$ 's. However, if they were a mix of various scaled measurements, then obviously the inner product's unit is incoherent—one would be resorted to adding measurements in different units, for example,  $\text{cm}^2$  and  $\text{kg}^2$  and so on. In such a case, a unitless inner product is appropriate, like the Mahalanobis inner product, which technically rescales the  $x_{ik}$ 's to unity. In summary, if the covariates are all measured on the same scale, then the I-prior is appropriate, and if not, the  $g$ -prior is appropriate.

## Appendix F

# Additional details for various I-prior regression models

These are additional details relating to discussion on various I-prior regression models in [Section 4.1 of Chapter 4](#) (p. 4). These details relate to the standard linear multilevel model and the naïve classification model.

### F.1 The I-prior for standard multilevel models

misc:multilevelmodels

We show the corresponding I-prior for the regression coefficients of the standard linear multilevel model [\(4.3\)](#). Write  $\alpha = \beta_0$ , and for simplicity, assume iid errors, i.e.,  $\Psi = \psi \mathbf{I}_n$ . The form of  $f \in \mathcal{F}$  is now  $f(\mathbf{x}_i^{(j)}, j) = \sum_{i'=1}^{n_{j'}} \sum_{j'=1}^m h_\lambda((\mathbf{x}_i^{(j)}, j), (\mathbf{x}_{i'}^{(j')}, j')) w_{i'j'}$ , where each  $w_{i'j'} \sim N(0, \psi^{-1})$ .

Now, functions in the scaled RKHS  $\mathcal{F}_2$  have the form

$$\begin{aligned} f_2(j) &= \sum_{i=1}^{n_{j'}} \sum_{j'=1}^m \lambda_2 \left( \frac{\delta_{jj'}}{p_j} - 1 \right) w_{ij'} \\ &= \lambda_2 \left( \frac{w_{+j}}{p_j} - w_{++} \right), \end{aligned}$$

where a '+' in the index of  $w_{ik}$  indicates a summation over that index, and  $p_j$  is the empirical distribution over  $\mathcal{M}$ , i.e.  $p_j = n_j/n$ . Clearly  $f_2(j)$  is a variable depending on

$j$ , so write  $f_2(j) = \beta_{0j}$ . The distribution of  $\beta_{0j}$  is normal with zero mean and variance

$$\begin{aligned}\text{Var } \beta_{0j} &= \lambda_2^2 \left( \frac{n_j \psi}{n_j^2/n^2} + n\psi \right) \\ &= n\psi \lambda_2^2 \left( \frac{1}{p_j} + 1 \right).\end{aligned}$$

The covariance between any two random intercepts  $\beta_{0j}$  and  $\beta_{0j'}$  is

$$\begin{aligned}\text{Cov}(\beta_{0j}, \beta_{0j'}) &= \text{Cov} \left( \lambda_2 \left( \frac{w_{+j}}{p_j} - w_{++} \right), \lambda_2 \left( \frac{w_{+j'}}{p_{j'}} - w_{++} \right) \right) \\ &= \frac{\lambda_2^2}{p_j p_{j'}} \text{Cov}(w_{+j}, w_{+j'}) - \frac{\lambda_2^2}{p_j} \text{Cov}(w_{+j}, w_{++}) - \frac{\lambda_2^2}{p_{j'}} \text{Cov}(w_{++}, w_{+j'}) \\ &\quad + \lambda_2^2 \text{Cov}(w_{++}, w_{++}) \\ &= -\frac{\lambda_2^2}{n_j/n} n_j \psi - \frac{\lambda_2^2}{n_{j'}/n} n_{j'} \psi + \lambda_2^2 n\psi \\ &= -n\psi \lambda_2^2.\end{aligned}$$

Functions in  $\mathcal{F}_{12}$ , on the other hand, have the form

$$\begin{aligned}f_{12}(\mathbf{x}_i, j) &= \sum_{i'=1}^{n_{j'}} \sum_{j'=1}^m \lambda_1 \lambda_2 \cdot \tilde{\mathbf{x}}_i^{(j)\top} \tilde{\mathbf{x}}_{i'}^{(j')} \cdot \left( \frac{\delta_{jj'}}{p_j} - 1 \right) w_{i'j'} \\ &= \tilde{\mathbf{x}}_i^{(j)\top} \underbrace{\left( \frac{\lambda_1 \lambda_2}{p_j} \sum_{i'=1}^{n_j} \tilde{\mathbf{x}}_{i'}^{(j)} w_{i'j} - \lambda_1 \lambda_2 \sum_{i'=1}^{n_{j'}} \sum_{j'=1}^m \tilde{\mathbf{x}}_{i'}^{(j')} w_{i'j'} \right)}_{\beta_{1j}},\end{aligned}$$

and this is, as expected, a linear form dependent on cluster  $j$ . We can calculate the variance for  $\beta_{1j}$  to be

$$\begin{aligned}\text{Var } \beta_{1j} &= \lambda_1^2 \lambda_2^2 \text{Var} \left( \frac{1}{p_j} \tilde{\mathbf{X}}_j^\top \mathbf{w}_j - \tilde{\mathbf{X}}^\top \mathbf{w} \right) \\ &= \lambda_1^2 \lambda_2^2 \left( \frac{\psi}{n_j^2/n^2} \tilde{\mathbf{X}}_j^\top \tilde{\mathbf{X}}_j + \psi \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \frac{1}{p_j} \tilde{\mathbf{X}}_j^\top \text{Cov}(\mathbf{w}_j, \mathbf{w}) \tilde{\mathbf{X}}^\top \right) \\ &= n\psi \lambda_1^2 \lambda_2^2 \left( \frac{1}{p_j} \mathbf{S}_j + \mathbf{S} - \mathbf{S}_j \right) \\ &= n\psi \lambda_1^2 \lambda_2^2 \left( \left( \frac{1}{p_j} - 1 \right) \mathbf{S}_j + \mathbf{S} \right)\end{aligned}$$

where  $\mathbf{S}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}})^\top (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}})$ ,  $\mathbf{S} = \frac{1}{n} \sum_{i=1}^{n_j} \sum_{j=1}^m (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}})^\top (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}})$ , and  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n_j} \sum_{j=1}^m \mathbf{x}_i^{(j)}$ . The covariance between two vectors of the random slopes is

$$\begin{aligned} \text{Cov}(\beta_{1j}, \beta_{1j'}) &= \lambda_1^2 \lambda_2^2 \text{Cov} \left( \frac{1}{p_j} \tilde{\mathbf{X}}_j^\top \mathbf{w}_j - \tilde{\mathbf{X}}^\top \mathbf{w}, \frac{1}{p_{j'}} \tilde{\mathbf{X}}_{j'}^\top \mathbf{w}_{j'} - \tilde{\mathbf{X}}^\top \mathbf{w} \right) \\ &= \psi \lambda_1^2 \lambda_2^2 \left( \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} - \frac{1}{p_j} \tilde{\mathbf{X}}_j^\top \tilde{\mathbf{X}}_j - \frac{1}{p_{j'}} \tilde{\mathbf{X}}_{j'}^\top \tilde{\mathbf{X}}_{j'} \right) \\ &= n \psi \lambda_1^2 \lambda_2^2 (\mathbf{S} - \mathbf{S}_j - \mathbf{S}_{j'}). \end{aligned}$$

Another quantity of interest is the covariance between the random intercepts and random slopes:

$$\begin{aligned} \text{Cov}(\beta_{0j}, \beta_{1j}) &= \lambda_1 \lambda_2^2 \text{Cov} \left( \frac{1}{p_j} \mathbf{1}_{n_j}^\top \mathbf{w}_j - \mathbf{1}_n^\top \mathbf{w}, \frac{1}{p_j} \tilde{\mathbf{X}}_j^\top \mathbf{w}_j - \tilde{\mathbf{X}}^\top \mathbf{w} \right) \\ &= \psi \lambda_1 \lambda_2^2 \left( \mathbf{1}_n^\top \tilde{\mathbf{X}}_j^0 + \frac{1}{p_j^2} \mathbf{1}_{n_j}^\top \tilde{\mathbf{X}}_j - \frac{2}{p_j} \mathbf{1}_{n_j}^\top \tilde{\mathbf{X}}_j \right) \\ &= n \psi \lambda_1 \lambda_2^2 \left( \left( \frac{1}{p_j} - 2 \right) \frac{1}{n_j} \sum_{i=1}^{n_j} (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}}) \right) \\ &= n \psi \lambda_1 \lambda_2^2 \left( \frac{1}{p_j} - 2 \right) (\bar{\mathbf{x}}^{(j)} - \bar{\mathbf{x}}) \end{aligned}$$

and

$$\begin{aligned} \text{Cov}(\beta_{0j}, \beta_{1j'}) &= \lambda_1 \lambda_2^2 \text{Cov} \left( \frac{1}{p_j} \mathbf{1}_{n_j}^\top \mathbf{w}_j - \mathbf{1}_n^\top \mathbf{w}, \frac{1}{p_{j'}} \tilde{\mathbf{X}}_{j'}^\top \mathbf{w}_{j'} - \tilde{\mathbf{X}}^\top \mathbf{w} \right) \\ &= \psi \lambda_1 \lambda_2^2 \left( \mathbf{1}_n^\top \tilde{\mathbf{X}}_{j'}^0 + \frac{1}{p_j p_{j'}} \mathbf{1}_{n_j}^\top \text{Cov}(\mathbf{w}_j, \mathbf{w}_{j'}) \tilde{\mathbf{X}}_{j'}^0 - \frac{1}{p_j} \mathbf{1}_{n_j}^\top \tilde{\mathbf{X}}_j - \frac{1}{p_{j'}} \mathbf{1}_{n_{j'}}^\top \tilde{\mathbf{X}}_{j'} \right) \\ &= n \psi \lambda_1 \lambda_2^2 \left( -\frac{1}{n_j} \sum_{i=1}^{n_j} (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}}) - \frac{1}{n_{j'}} \sum_{i=1}^{n_{j'}} (\mathbf{x}_i^{(j')} - \bar{\mathbf{x}}) \right) \\ &= n \psi \lambda_1 \lambda_2^2 (2\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(j)} - \bar{\mathbf{x}}^{(j')}). \end{aligned}$$

## F.2 The I-prior for naïve classification

For the naïve I-prior classification model (4.7), the I-prior is derived as follows. Firstly, the functions in  $\mathcal{F}_{\mathcal{M}}$  and  $\mathcal{F}_{\mathcal{X}}$  need necessarily be zero-mean functions (as per the functional ANOVA definition in Definition 2.36, but also, as per the definition of the Pearson RKHS and centred identity kernel RKHS). What this means is that  $\sum_{j=1}^m \alpha_j = 0$ ,

$\sum_{j=1}^m f_j(x_i) = 0$ , and  $\sum_{i=1}^n f_j(x_i) = 0$ . In particular,

$$\begin{aligned} \mathbb{E} \left[ \sum_{j=1}^m y_{ij} \right] &= \sum_{j=1}^m (\alpha + \alpha_j + f_j(x_i)) \\ &= m\alpha + \sum_{j=1}^m \alpha_j + \sum_{j=1}^m f_j(x_i) \end{aligned}$$

and since  $\sum_{j=1}^m y_{ij} = 1$ , we get the ML estimate  $\hat{\alpha} = 1/m$ , and thus the grand intercept can be fixed to resolve identification.

It is much more convenient to work in vector and matrix form, so let us introduce some notation. Let  $\mathbf{w}$  (c.f.  $\mathbf{y}$ ,  $\mathbf{f}$  and  $\boldsymbol{\epsilon}$ ) be an  $n \times m$  matrix whose  $(i, j)$  entries contain  $w_{ij}$  (c.f.  $y_{ij}$ ,  $f(x_i, j)$ , and  $\epsilon_{ij}$ ). The row-wise entries of  $\mathbf{w}$  are independent of each other (independence assumption of the  $n$  observations), while any two of their columns have covariance as specified in  $\boldsymbol{\Psi}$ . This means that  $\mathbf{w}$  follows a matrix normal distribution  $\text{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \boldsymbol{\Psi})$ , which implies  $\text{vec } \mathbf{w} \sim \text{N}_{nm}(\mathbf{0}, \boldsymbol{\Psi} \otimes \mathbf{I}_n)$ , and similarly,  $\boldsymbol{\epsilon} \sim \text{N}_{nm}(\mathbf{0}, \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n)$ . Denote by  $\mathbf{B}_\eta$  the  $n \times n$  kernel matrix with entries supplied by kernel  $1 + b_\eta$  over  $\mathcal{X} \times \mathcal{X}$ , and  $\mathbf{A}$  the  $m \times m$  matrix with entries supplied by  $a$  over  $\mathcal{M} \times \mathcal{M}$ . From (4.7), we have that

$$\mathbf{f} = \mathbf{B}_\eta \mathbf{w} \mathbf{A} \in \mathbb{R}^{n \times m},$$

and thus  $\text{vec } \mathbf{f} \sim \text{N}_{nm}(\mathbf{0}, \mathbf{A} \boldsymbol{\Psi} \mathbf{A} \otimes \mathbf{B}_\eta^2)$ . As  $\mathbf{y} = \mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{f} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\alpha} \in \mathbb{R}^m$  with  $j$ 'th component  $\alpha + \alpha_j = 1/m + \alpha_j$ , by linearity we have that

$$\text{vec } \mathbf{y} \sim \text{N}_{nm}(\text{vec } \boldsymbol{\alpha}, \mathbf{A} \boldsymbol{\Psi} \mathbf{A} \otimes \mathbf{B}_\eta^2 + \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \quad (\text{F.1})$$

and

$$\text{vec } \mathbf{y} | \mathbf{w} \sim \text{N}_{nm}(\text{vec}(\boldsymbol{\alpha} + \mathbf{B}_\eta \mathbf{w} \mathbf{A}), \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n). \quad (\text{F.2})$$

By the results of Chapter 4, the posterior distribution of the I-prior random effects is  $\text{vec } \mathbf{w} | \mathbf{y} \sim \text{N}(\text{vec } \tilde{\mathbf{w}}, \tilde{\mathbf{V}}_w)$ , where

$$\text{vec } \tilde{\mathbf{w}} = \tilde{\mathbf{V}}_w (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta) \text{vec}(\mathbf{y} - \mathbf{1}_n \boldsymbol{\alpha}^\top) \quad \text{and} \quad \tilde{\mathbf{V}}_w^{-1} = \mathbf{A} \boldsymbol{\Psi} \mathbf{A} \otimes \mathbf{B}_\eta^2 + \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n = \mathbf{V}_y. \quad (\text{F.3})$$

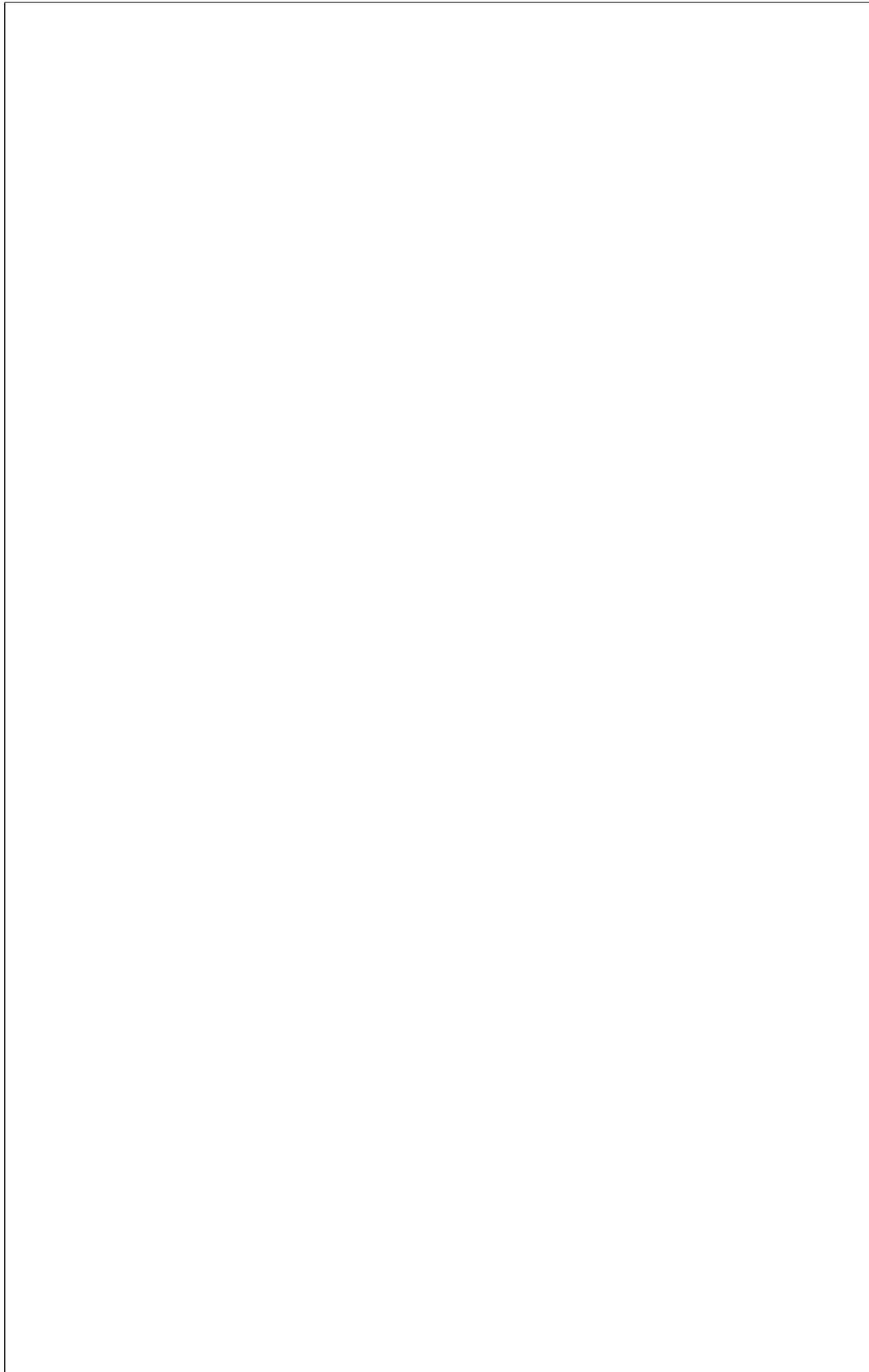
Suppose hypothetically, one uses the uncentered identity kernel  $a(j, j') = \delta_{jj'}$ , in which case centring of the intercepts  $\alpha_j$  must be handled separately. In conjunction with an assumption of iid errors ( $\boldsymbol{\Psi} = \psi \mathbf{I}_n$ ), the above distributions simplify further.



Specifically, the variance in the marginal distribution becomes

$$\begin{aligned}
 \text{Var}[\text{vec } \mathbf{y}] &= (\psi \mathbf{I}_m \otimes \mathbf{B}_\eta^2) + (\psi^{-1} \mathbf{I}_m \otimes \mathbf{I}_n) \\
 &= (\mathbf{I}_m \otimes \psi \mathbf{B}_\eta^2) + (\mathbf{I}_m \otimes \psi^{-1} \mathbf{I}_n) \\
 &= \mathbf{I}_m \otimes \overbrace{(\psi \mathbf{B}_\eta^2 + \psi^{-1} \mathbf{I}_n)}^{\tilde{\mathbf{V}}_y}.
 \end{aligned}$$

which implies independence and identical variances  $\tilde{\mathbf{V}}_y$  for the vectors  $(y_{1j}, \dots, y_{nj})^\top$  for each class  $j = 1, \dots, m$ . Evidently, this stems from the implied independence structure of the prior on  $f$  too, since now  $\text{Var}(\text{vec } \mathbf{f}) = \text{diag}(\psi \mathbf{B}_\eta^2, \dots, \psi \mathbf{B}_\eta^2)$ , which could be interpreted as having independent and identical I-priors on the regression functions for each class  $\mathbf{f}_{\cdot j} = (f(x_1, j), \dots, f(x_n, j))^\top$ .



## Appendix G

# Posterior distribution of the I-prior regression function

We derive the posterior distribution for the I-prior random effects  $\mathbf{w} = (w_1, \dots, w_n)^\top$ , which is related to the I-prior regression function via  $f(x_i) = \sum_{k=1}^n h_\eta(x_i, x_k) w_k$ , or in matrix terms,  $\mathbf{f} := (f(x_1), \dots, f(x_n))^\top = \mathbf{H}_\eta \mathbf{w}$ , and  $f \in \mathcal{F}$  an RKHS with kernel  $h_\eta$ . A closely related distribution of interest is the posterior predictive distribution of  $y_{\text{new}}$ , the prediction at a new data point  $x_{\text{new}}$ . We note the similarity of these results with the posterior distributions of Gaussian process regressions (Rasmussen and Williams, 2006).

### G.1 Deriving the posterior distribution for $\mathbf{w}$

apx:posteriorw

In the following derivation, we implicitly assume the dependence on  $\mathbf{f}_0$  and  $\theta$ . The distribution of  $\mathbf{y}|\mathbf{w}$  is  $N_n(\boldsymbol{\alpha} + \mathbf{f}_0 + \mathbf{H}_\eta \mathbf{w}, \boldsymbol{\Psi}^{-1})$ , where  $\boldsymbol{\alpha} = \alpha \mathbf{1}_n$ , while the prior distribution for  $\mathbf{w}$  is  $N_n(\mathbf{0}, \boldsymbol{\Psi})$ . Since  $p(\mathbf{w}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$ , we have that

$$\begin{aligned} \log p(\mathbf{w}|\mathbf{y}) &= \log p(\mathbf{y}|\mathbf{w}) + \log p(\mathbf{w}) \\ &= \text{const.} + \frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} (\mathbf{y} - \boldsymbol{\alpha} - \mathbf{f}_0 - \mathbf{H}_\eta \mathbf{w})^\top \boldsymbol{\Psi} (\mathbf{y} - \boldsymbol{\alpha} - \mathbf{f}_0 - \mathbf{H}_\eta \mathbf{w}) \\ &\quad - \frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \mathbf{w}^\top \boldsymbol{\Psi}^{-1} \mathbf{w} \\ &= \text{const.} - \frac{1}{2} \mathbf{w}^\top (\mathbf{H}_\eta \boldsymbol{\Psi} \mathbf{H}_\eta + \boldsymbol{\Psi}^{-1}) \mathbf{w} + (\mathbf{y} - \boldsymbol{\alpha} - \mathbf{f}_0)^\top \boldsymbol{\Psi} \mathbf{H}_\eta \mathbf{w}. \end{aligned}$$

Setting  $\mathbf{A} = \mathbf{H}_\eta \boldsymbol{\Psi} \mathbf{H}_\eta + \boldsymbol{\Psi}^{-1}$ ,  $\mathbf{a}^\top = (\mathbf{y} - \boldsymbol{\alpha} - \mathbf{f}_0)^\top \boldsymbol{\Psi} \mathbf{H}_\eta$ , and using the fact that

$$\mathbf{w}^\top \mathbf{A} \mathbf{w} - 2\mathbf{a}^\top \mathbf{w} = (\mathbf{w} - \mathbf{A}^{-1} \mathbf{a})^\top \mathbf{A} (\mathbf{w} - \mathbf{A}^{-1} \mathbf{a}),$$

we have that  $\mathbf{w}|\mathbf{y}$  is normally distributed with the required mean and variance.

Alternatively, one could have shown this using standard results of multivariate normal distributions. Noting that the covariance between  $\mathbf{y}$  and  $\mathbf{w}$  is

$$\begin{aligned}\text{Cov}[\mathbf{y}, \mathbf{w}] &= \text{Cov}[\boldsymbol{\alpha} + \mathbf{f}_0 + \mathbf{H}_\eta \mathbf{w} + \boldsymbol{\epsilon}, \mathbf{w}] \\ &= \mathbf{H}_\eta \text{Cov}[\mathbf{w}, \mathbf{w}] \\ &= \mathbf{H}_\eta \boldsymbol{\Psi}\end{aligned}$$

and that  $\text{Cov}[\mathbf{w}, \mathbf{y}] = \boldsymbol{\Psi} \mathbf{H}_\eta = \mathbf{H}_\eta \boldsymbol{\Psi} = \text{Cov}[\mathbf{y}, \mathbf{w}]$  by symmetry, the joint distribution  $(\mathbf{y}, \mathbf{w})$  is

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{w} \end{pmatrix} \sim N_{n+n} \left( \begin{pmatrix} \boldsymbol{\alpha} + \mathbf{f}_0 \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V}_y & \mathbf{H}_\eta \boldsymbol{\Psi} \\ \boldsymbol{\Psi} \mathbf{H}_\eta & \boldsymbol{\Psi} \end{pmatrix} \right).$$

Thus,

$$\begin{aligned}\mathbb{E}[\mathbf{w}|\mathbf{y}] &= \mathbb{E} \mathbf{w} + \text{Cov}(\mathbf{w}, \mathbf{y})(\text{Var} \mathbf{y})^{-1}(\mathbf{y} - \mathbb{E} \mathbf{y}) \\ &= \boldsymbol{\Psi} \mathbf{H}_\eta \mathbf{V}_y^{-1}(\mathbf{y} - \boldsymbol{\alpha} - \mathbf{f}_0),\end{aligned}$$

and

$$\begin{aligned}\text{Var}[\mathbf{w}|\mathbf{y}] &= \text{Var} \mathbf{w} - \text{Cov}(\mathbf{w}, \mathbf{y})(\text{Var} \mathbf{y})^{-1} \text{Cov}(\mathbf{y}, \mathbf{w}) \\ &= \boldsymbol{\Psi} - \mathbf{H}_\eta \boldsymbol{\Psi} \mathbf{V}_y^{-1} \mathbf{H}_\eta \boldsymbol{\Psi} \\ &= \boldsymbol{\Psi} - \boldsymbol{\Psi} \mathbf{H}_\eta (\boldsymbol{\Psi}^{-1} + \mathbf{H}_\eta \boldsymbol{\Psi} \mathbf{H}_\eta)^{-1} \mathbf{H}_\eta \boldsymbol{\Psi} \\ &= (\boldsymbol{\Psi}^{-1} + \mathbf{H}_\eta \boldsymbol{\Psi} \mathbf{H}_\eta)^{-1} \\ &= \mathbf{V}_y^{-1}\end{aligned}$$

as a direct consequence of the Woodbury matrix identity (Petersen and Pedersen, 2008, eq. 156, §3.2.2).

## G.2 Deriving the posterior predictive distribution

The posterior predictive distribution is obtained in an empirical Bayesian manner, in which the parameters of the model are replaced with their ML estimates (denoted with hats).

A priori, assume that  $y_{\text{new}} \sim N(\hat{\alpha}, v_{\text{new}})$ , where  $v_{\text{new}} = \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{h}_{\hat{\eta}}(x_{\text{new}}) + \psi_{\text{new}}^{-1}$ . Consider the joint distribution of  $(y_{\text{new}}, \mathbf{y}^\top)^\top$ , which is multivariate normal (since both  $y_{\text{new}}$  and  $\mathbf{y}$  are. Write

$$\begin{pmatrix} y_{\text{new}} \\ \mathbf{y} \end{pmatrix} \sim N_{n+1} \left( \begin{pmatrix} \hat{\alpha} \\ \hat{\alpha} \mathbf{1}_n \end{pmatrix}, \begin{pmatrix} v_{\text{new}} & \text{Cov}[y_{\text{new}}, \mathbf{y}] \\ \text{Cov}[y_{\text{new}}, \mathbf{y}]^\top & \hat{\mathbf{V}}_y \end{pmatrix} \right),$$

where

$$\begin{aligned}
\text{Cov}[y_{\text{new}}, \mathbf{y}] &= \text{Cov}[f_{\text{new}} + \epsilon_{\text{new}}, \mathbf{f} + \boldsymbol{\epsilon}] \\
&= \text{Cov}[f_{\text{new}}, \mathbf{f}] + \text{Cov}(\epsilon_{\text{new}}, \boldsymbol{\epsilon}) \\
&= \text{Cov} \left[ \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \tilde{\mathbf{w}}, \mathbf{H}_{\hat{\eta}} \tilde{\mathbf{w}} \right] + (\sigma_{\text{new},1}, \dots, \sigma_{\text{new},n}) \\
&= \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} + \boldsymbol{\sigma}_{\text{new}}.
\end{aligned}$$

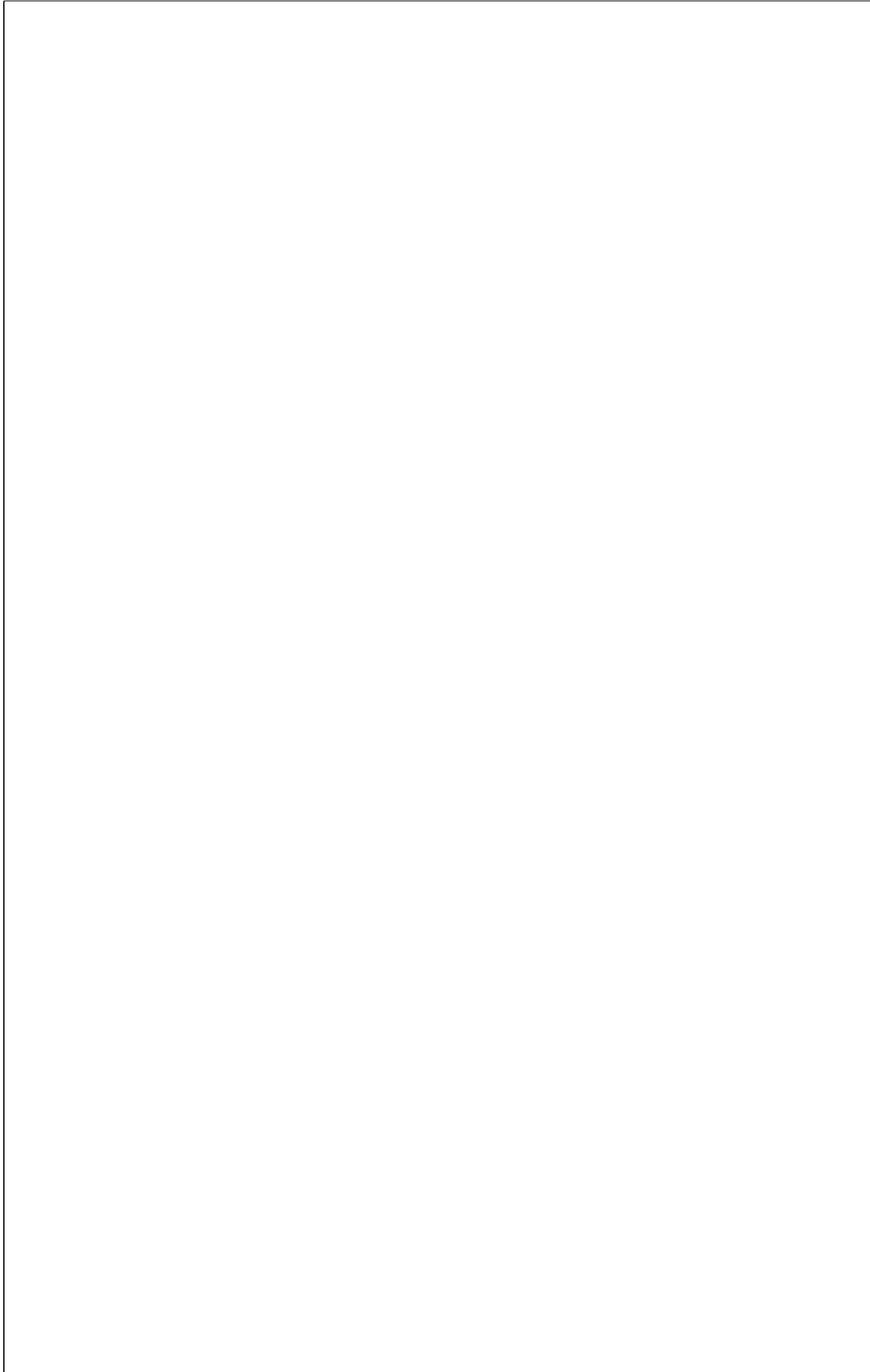
The vector of covariances  $\boldsymbol{\sigma}_{\text{new}}$  between observations  $y_1, \dots, y_n$  and the predicted point  $y_{\text{new}}$  would need to be prescribed a priori (treated as extra parameters), or estimated again, which seems excessive. Under an iid assumption of the error precisions, then  $\boldsymbol{\sigma}_{\text{new}} = \mathbf{0}$  would be acceptable.

In any case, using standard multivariate normal results, we get that  $y_{\text{new}}|\mathbf{y}$  is also normally distributed with mean

$$\begin{aligned}
E[y_{\text{new}}|\mathbf{y}] &= \hat{\alpha} + (\mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} + \boldsymbol{\sigma}_{\text{new}}) \hat{\mathbf{V}}_y^{-1} \tilde{\mathbf{y}} \\
&= \hat{\alpha} + \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \overbrace{\mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} \hat{\mathbf{V}}_y^{-1} \tilde{\mathbf{y}}}^{\hat{\mathbf{w}}} + \boldsymbol{\sigma}_{\text{new}} \hat{\mathbf{V}}_y^{-1} \tilde{\mathbf{y}} \\
&= \hat{\alpha} + E[f(x_{\text{new}})|\mathbf{y}] + \text{mean correction term}
\end{aligned}$$

and variance

$$\begin{aligned}
\text{Var}[y_{\text{new}}|\mathbf{y}] &= v_{\text{new}} - (\mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} + \boldsymbol{\sigma}_{\text{new}}) \hat{\mathbf{V}}_y^{-1} (\mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} + \boldsymbol{\sigma}_{\text{new}})^\top \\
&= \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{h}_{\hat{\eta}}(x_{\text{new}}) + \psi_{\text{new}}^{-1} - \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} \hat{\mathbf{V}}_y^{-1} \mathbf{H}_{\hat{\eta}} \hat{\boldsymbol{\Psi}} \mathbf{h}_{\hat{\eta}}(x_{\text{new}}) \\
&\quad + \text{variance correction term} \\
&= \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top (\hat{\boldsymbol{\Psi}} - \hat{\boldsymbol{\Psi}} \mathbf{H}_{\hat{\eta}} \hat{\mathbf{V}}_y^{-1} \mathbf{H}_{\hat{\eta}} \hat{\boldsymbol{\Psi}}) \mathbf{h}_{\hat{\eta}}(x_{\text{new}}) + \psi_{\text{new}}^{-1} \\
&\quad + \text{variance correction term} \\
&= \mathbf{h}_{\hat{\eta}}(x_{\text{new}})^\top \hat{\mathbf{V}}_y^{-1} \mathbf{h}_{\hat{\eta}}(x_{\text{new}}) + \psi_{\text{new}}^{-1} + \text{variance correction term} \\
&= \text{Var}[f(x_{\text{new}})|\mathbf{y}] + \psi_{\text{new}}^{-1} + \text{variance correction term}.
\end{aligned}$$



## Appendix H

# Variational EM algorithm for I-probit models

apx:varemip  
robit

The two sections that follow detail the derivation of the variational densities used in the E-step of the variational EM algorithm, and also the lower bound (ELBO) used to monitor convergence.

### H.1 Derivation of the variational densities

In what follows, the implicit dependence of the densities on the parameters of the model  $\theta$  are dropped. We derive a mean-field variational approximation of

$$\begin{aligned} p(\mathbf{y}^*, \mathbf{w} | \mathbf{y}) &\approx q(\mathbf{y}^*)q(\mathbf{w}) \\ &= \prod_{i=1}^n q(\mathbf{y}_i^*)q(\mathbf{w}). \end{aligned}$$

The first line is by assumption, while the second line follows from an induced factorisation on the latent propensities, as we will see later. Recall that the optimal mean-field variational density  $\tilde{q}$  satisfy

$$\log \tilde{q}(\mathbf{y}^*) = \mathbb{E}_{\mathbf{w} \sim \tilde{q}} [\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w})] + \text{const.} \quad (\text{from 5.13})$$

$$\log \tilde{q}(\mathbf{w}) = \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} [\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w})] + \text{const.} \quad (\text{from 5.14})$$

The joint likelihood is given by

$$p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}) = p(\mathbf{y} | \mathbf{y}^*)p(\mathbf{y}^* | \mathbf{w})p(\mathbf{w}).$$

For reference, the three relevant distributions are listed below.

- $p(\mathbf{y}|\mathbf{y}^*)$ . For each observation  $i \in \{1, \dots, n\}$ , given the corresponding latent propensities  $\mathbf{y}_i^* = (y_{i1}^*, \dots, y_{im}^*)$ , the distribution for  $y_i$  is a degenerate distribution which depends on the  $j$ 'th component of  $\mathbf{y}_i^*$  being largest, where the value observed for  $y_i$  was  $j$ . Since each of the  $y_i$ 's are independent, everything is multiplicative.

$$p(\mathbf{y}|\mathbf{y}^*) = \prod_{i=1}^n \prod_{j=1}^m p_{ij}^{[y_i=j]} = \prod_{i=1}^n \prod_{j=1}^m \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \mathbb{1}[y_i=j].$$

- $p(\mathbf{y}^*|\mathbf{w})$ . Given values for the parameters and I-prior random effects, the distribution of the latent propensities is matrix normal

$$\mathbf{y}^*|\mathbf{w} \sim \text{MN}_{n,m}(\mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{H}_\eta \mathbf{w}, \mathbf{I}_n, \boldsymbol{\Psi}^{-1}).$$

Write  $\boldsymbol{\mu} = \mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{H}_\eta \mathbf{w}$ . Its pdf is

$$\begin{aligned} p(\mathbf{y}^*|\mathbf{w}) &= \exp \left[ -\frac{nm}{2} \log 2\pi + \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \text{tr}((\mathbf{y}^* - \boldsymbol{\mu}) \boldsymbol{\Psi} (\mathbf{y}^* - \boldsymbol{\mu})^\top) \right] \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi + \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot}) \right], \end{aligned}$$

where  $\mathbf{y}_{i\cdot}^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}_{i\cdot} \in \mathbb{R}^m$  are the rows of  $\mathbf{y}^*$  and  $\boldsymbol{\mu}$  respectively. The second line follows directly from the definition of the trace, but also emanates from the fact that  $\mathbf{y}_{i\cdot}^*$  are independent multivariate normal with mean  $\boldsymbol{\mu}_{i\cdot}$  and variance  $\boldsymbol{\Psi}^{-1}$ .

- $p(\mathbf{w})$ . The  $\mathbf{w}$ 's are normal random matrices  $\mathbf{w} \sim \text{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \boldsymbol{\Psi})$  with pdf

$$\begin{aligned} p(\mathbf{w}) &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \text{tr}(\mathbf{w} \boldsymbol{\Psi}^{-1} \mathbf{w}^\top) \right] \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \sum_{i=1}^n \mathbf{w}_{i\cdot}^\top \boldsymbol{\Psi}^{-1} \mathbf{w}_{i\cdot} \right]. \end{aligned}$$

### H.1.1 Derivation of $\tilde{q}(\mathbf{y}^*)$

The rows of  $\mathbf{y}^*$  are independent, and thus we can consider the variational density for each  $\mathbf{y}_i^*$  separately. Consider the case where  $y_i$  takes one particular value  $j \in \{1, \dots, m\}$ . In such cases, we have that  $y_{ij}^* > y_{ik}$  for all  $k \neq j$ , and that

$$\begin{aligned} \log \tilde{q}(\mathbf{y}_{i\cdot}^*) &= \mathbb{E}_{\mathbf{w} \sim \tilde{q}} \left[ -\frac{1}{2} (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_i)^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_i) \right] + \text{const.} \\ &= \left[ -\frac{1}{2} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_i)^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_i) \right] + \text{const.} \end{aligned} \quad (\star)$$



where  $\tilde{\boldsymbol{\mu}}_{i.} = \boldsymbol{\alpha} + \tilde{\mathbf{w}}\mathbf{h}_\eta(x_i)$ ,  $\tilde{\mathbf{w}} = \mathbb{E}_{\mathbf{w} \sim \tilde{q}}[\mathbf{w}]$ . This is recognised as the logarithm of a multivariate normal pdf with mean  $\tilde{\boldsymbol{\mu}}_{i.}$  and variance  $\boldsymbol{\Psi}^{-1}$ . On the other hand, when  $y_i \neq j$ , the pdf is zero. Thus,

$$\tilde{q}(\mathbf{y}_{i.}^*) = \begin{cases} \phi(\mathbf{y}_{i.}^* | \tilde{\boldsymbol{\mu}}_{i.}, \boldsymbol{\Psi}^{-1}) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise,} \end{cases}$$

implying a truncated multivariate normal distribution for  $\mathbf{y}_{i.}^*$ . The required moments from the truncated multivariate normal distribution can be obtained using the methods described in [Appendix C.4](#) (p. 17).

*Remark H.1.* In the above derivation, we needn't consider the second order terms in the expectations because they do not involve  $\mathbf{y}_{i.}^*$ , and thus, these terms can be absorbed into the constant. To see this,

$$\begin{aligned} \mathbb{E}[(\mathbf{y}_{i.}^* - \boldsymbol{\mu}_{i.})^\top \boldsymbol{\Psi}(\mathbf{y}_{i.}^* - \boldsymbol{\mu}_{i.})] &= \mathbb{E}[\mathbf{y}_{i.}^{*\top} \boldsymbol{\Psi} \mathbf{y}_{i.}^* + \boldsymbol{\mu}_{i.}^\top \boldsymbol{\Psi} \boldsymbol{\mu}_{i.} - 2\boldsymbol{\mu}_{i.}^\top \boldsymbol{\Psi} \mathbf{y}_{i.}^*] \\ &= \mathbf{y}_{i.}^{*\top} \boldsymbol{\Psi} \mathbf{y}_{i.}^* - 2\mathbb{E}[\boldsymbol{\mu}_{i.}^\top] \boldsymbol{\Psi} \mathbf{y}_{i.}^* + \text{const.} \\ &= \mathbf{y}_{i.}^{*\top} \boldsymbol{\Psi} \mathbf{y}_{i.}^* - 2\tilde{\boldsymbol{\mu}}_{i.}^\top \boldsymbol{\Psi} \mathbf{y}_{i.}^* + \text{const.} \\ &= (\mathbf{y}_{i.}^* - \tilde{\boldsymbol{\mu}}_{i.})^\top \boldsymbol{\Psi}(\mathbf{y}_{i.}^* - \tilde{\boldsymbol{\mu}}_{i.}) + \text{const.} \end{aligned}$$

The square is then completed to get the final line, which is the expression for the term  $(\star)$  multiplied by a half.

### H.1.2 Derivation of $\tilde{q}(\mathbf{w})$

apx:qw

The terms involving  $\mathbf{w}$  in the joint likelihood (5.14) are the  $p(\mathbf{y}^*|\mathbf{w})$  and  $p(\mathbf{w})$  terms, so the rest are absorbed into the constant. The easiest way to derive  $\tilde{q}(\mathbf{w})$  is to vectorise  $\mathbf{y}^*$  and  $\mathbf{w}$ . We know that

$$\begin{aligned} \text{vec } \mathbf{y}^* | \boldsymbol{\alpha}, \mathbf{w}, \eta, \boldsymbol{\Psi} &\sim N_{nm}(\text{vec}(\mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{H}_\eta \mathbf{w}), \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \\ &\text{and} \\ \text{vec } \mathbf{w} | \boldsymbol{\Psi} &\sim N_{nm}(\mathbf{0}, \boldsymbol{\Psi} \otimes \mathbf{I}_n) \end{aligned}$$

using properties of matrix normal distributions.

We also use the fact that  $\text{vec}(\mathbf{H}_\eta \mathbf{w}) = (\mathbf{I}_m \otimes \mathbf{H}_\eta) \text{vec } \mathbf{w}$ . For simplicity, write  $\bar{\mathbf{y}}^* = \text{vec}(\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top)$ , and  $\mathbf{M} = (\mathbf{I}_m \otimes \mathbf{H}_\eta)$ . Thus,

$$\begin{aligned} \log \tilde{q}(\mathbf{w}) &= \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ -\frac{1}{2} (\bar{\mathbf{y}}^* - \mathbf{M} \text{vec } \mathbf{w})^\top (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n)^{-1} (\bar{\mathbf{y}}^* - \mathbf{M} \text{vec } \mathbf{w}) \right] \\ &\quad + \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ -\frac{1}{2} (\text{vec } \mathbf{w})^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n)^{-1} \text{vec } (\mathbf{w}) \right] + \text{const.} \\ &= -\frac{1}{2} \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ (\text{vec } \mathbf{w})^\top \overbrace{(\mathbf{M}^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n) \mathbf{M} + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n))}^{\mathbf{A}} \text{vec } (\mathbf{w}) \right] \\ &\quad + \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ \overbrace{\bar{\mathbf{y}}^{*\top} (\boldsymbol{\Psi} \otimes \mathbf{I}_n) \mathbf{M}}^{\mathbf{a}^\top} \text{vec } (\mathbf{w}) \right] + \text{const.} \\ &= -\frac{1}{2} \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ (\text{vec } \mathbf{w} - \mathbf{A}^{-1} \mathbf{a})^\top \mathbf{A} (\text{vec } \mathbf{w} - \mathbf{A}^{-1} \mathbf{a}) \right] + \text{const.} \end{aligned}$$

This is recognised as a multivariate normal of dimension  $nm$  with mean and precision given by  $\text{vec } \tilde{\mathbf{w}} = \mathbb{E}[\mathbf{A}^{-1} \mathbf{a}]$  and  $\tilde{\mathbf{V}}_w^{-1} = \mathbb{E}[\mathbf{A}]$  respectively. With a little algebra, we find that

$$\begin{aligned} \tilde{\mathbf{V}}_w &= \{ \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}}[\mathbf{A}] \}^{-1} \\ &= \left\{ \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} \left[ (\mathbf{I}_m \otimes \mathbf{H}_\eta)^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n) (\mathbf{I}_m \otimes \mathbf{H}_\eta) + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \right] \right\}^{-1} \\ &= (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta^2 + \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n)^{-1} \end{aligned}$$

and

$$\begin{aligned} \text{vec } \tilde{\mathbf{w}} &= \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}}[\mathbf{A}^{-1} \mathbf{a}] \\ &= \tilde{\mathbf{V}}_w \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} [(\mathbf{I}_m \otimes \mathbf{H}_\eta) (\boldsymbol{\Psi} \otimes \mathbf{I}_n) \text{vec}(\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top)] \\ &= \tilde{\mathbf{V}}_w (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta) \mathbb{E}_{\mathbf{y}^* \sim \tilde{q}} [\text{vec}(\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top)] \\ &= \tilde{\mathbf{V}}_w (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta) \text{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top). \end{aligned}$$

We will often refer to  $\tilde{\mathbf{w}}$  as the  $n \times m$  matrix constructed by filling in its entries with  $\text{vec } \tilde{\mathbf{w}}$  column-wise (akin to the opposite of vectorisation). This way, the  $\tilde{\mathbf{w}}$  contains posterior mean values arranged by class  $j = 1, \dots, m$  column-wise, and by observations  $i = 1, \dots, n$  row-wise. Ideally, we do not want to work with the  $nm \times nm$  matrix  $\mathbf{V}_w$ , since its inverse is expensive to compute. Refer to [Section 5.6.2](#) for details.

In the case of the I-probit model, where  $\Psi = \text{diag}(\psi_1, \dots, \psi_m)$ , then the covariance matrix takes a simpler form. Specifically, it has the block diagonal structure:

$$\begin{aligned}\tilde{\mathbf{V}}_w &= (\text{diag}(\psi_1, \dots, \psi_m) \otimes \mathbf{H}_\eta^2 + \text{diag}(\psi_1, \dots, \psi_m) \otimes \mathbf{I}_n)^{-1} \\ &= \text{diag}\left((\psi_1 \mathbf{H}_\eta^2 + \psi_1^{-1} \mathbf{I}_n)^{-1}, \dots, (\psi_m \mathbf{H}_\eta^2 + \psi_m^{-1} \mathbf{I}_n)^{-1}\right) \\ &=: \text{diag}(\tilde{\mathbf{V}}_{w_1}, \dots, \tilde{\mathbf{V}}_{w_m}).\end{aligned}$$

The mean  $\text{vec } \tilde{\mathbf{w}}$  is

$$\begin{aligned}\text{vec } \tilde{\mathbf{w}} &= \tilde{\mathbf{V}}_w (\text{diag}(\psi_1, \dots, \psi_m) \otimes \tilde{\mathbf{H}}_\eta) \text{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top) \\ &= \text{diag}(\tilde{\mathbf{V}}_{w_1}, \dots, \tilde{\mathbf{V}}_{w_m}) \text{diag}(\psi_1 \mathbf{H}_\eta, \dots, \psi_m \mathbf{H}_\eta) \text{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top) \\ &= \text{diag}(\psi_1 \tilde{\mathbf{V}}_{w_1} \mathbf{H}_\eta, \dots, \psi_m \tilde{\mathbf{V}}_{w_m} \mathbf{H}_\eta) (\tilde{\mathbf{y}}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top) \\ &= \begin{pmatrix} \tilde{\mathbf{w}}_{\cdot,1}^\top & \dots & \tilde{\mathbf{w}}_{\cdot,m}^\top \end{pmatrix}^\top \\ &= \left( (\psi_1 \tilde{\mathbf{V}}_{w_1} \mathbf{H}_\eta (\tilde{\mathbf{y}}_{\cdot,1}^* - \alpha_1 \mathbf{1}_n))^\top \quad \dots \quad (\psi_m \tilde{\mathbf{V}}_{w_m} \mathbf{H}_\eta (\tilde{\mathbf{y}}_{\cdot,m}^* - \alpha_m \mathbf{1}_n))^\top \right)^\top.\end{aligned}$$

Therefore, we can consider the distribution of  $\mathbf{w} = (\mathbf{w}_{\cdot,1}, \dots, \mathbf{w}_{\cdot,m})$  columnwise, and each are normally distributed with mean and variance

$$\tilde{\mathbf{w}}_{\cdot,j} = \psi_j \tilde{\mathbf{V}}_{w_j} \mathbf{H}_\eta (\tilde{\mathbf{y}}_{\cdot,j}^* - \alpha_j \mathbf{1}_n) \quad \text{and} \quad \tilde{\mathbf{V}}_{w_j} = (\psi_j \mathbf{H}_\eta^2 + \psi_j^{-1} \mathbf{I}_n)^{-1}.$$

A quantity that we will be requiring time and again will be  $\text{tr}(\mathbf{C} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}])$ , where  $\mathbf{C} \in \mathbb{R}^{m \times m}$  and  $\mathbf{D} \in \mathbb{R}^{n \times n}$  are both square and symmetric matrices. Using the definition of the trace directly, we get

$$\begin{aligned}\text{tr}(\mathbf{C} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]) &= \sum_{i,j=1}^m \mathbf{C}_{ij} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]_{ij} \\ &= \sum_{i,j=1}^m \mathbf{C}_{ij} \mathbf{E}[\mathbf{w}_{\cdot,i}^\top \mathbf{D} \mathbf{w}_{\cdot,j}].\end{aligned} \tag{H.1}$$

The expectation of the univariate quantity  $\mathbf{w}_{\cdot,i}^\top \mathbf{D} \mathbf{w}_{\cdot,j}$  is inspected below:

$$\begin{aligned}\mathbf{E}[\mathbf{w}_{\cdot,i}^\top \mathbf{D} \mathbf{w}_{\cdot,j}] &= \text{tr}(\mathbf{D} \mathbf{E}[\mathbf{w}_{\cdot,j} \mathbf{w}_{\cdot,i}^\top]) \\ &= \text{tr}(\mathbf{D} (\text{Cov}(\mathbf{w}_{\cdot,j}, \mathbf{w}_{\cdot,i}) + \mathbf{E}[\mathbf{w}_{\cdot,j}] \mathbf{E}[\mathbf{w}_{\cdot,i}]^\top)) \\ &= \text{tr}(\mathbf{D} (\mathbf{V}_w[i, j] + \tilde{\mathbf{w}}_{\cdot,j} \tilde{\mathbf{w}}_{\cdot,i}^\top)).\end{aligned}$$

where  $\mathbf{V}_w[i, j] \in \mathbb{R}^{n \times n}$  refers to the  $(i, j)$ 'th submatrix block of  $\mathbf{V}_w$ . Of course, in the independent the I-probit model, this is equal to

$$\mathbf{V}_w[i, j] = \delta_{ij} (\psi_j \mathbf{H}_\eta^2 + \psi_j^{-1} \mathbf{I}_n)^{-1}$$

where  $\delta$  is the Kronecker delta. Continuing on (H.1) leads us to

$$\text{tr}(\mathbf{C} \mathbb{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]) = \sum_{i,j=1}^m \mathbf{C}_{ij} \left( \text{tr}(\mathbf{D}(\delta_{ij} \mathbf{V}_{w_j} + \tilde{\mathbf{w}}_{\cdot j} \tilde{\mathbf{w}}_{\cdot i}^\top)) \right).$$

If  $\mathbf{C} = \text{diag}(c_1, \dots, c_m)$ , then

$$\begin{aligned} \text{tr}(\mathbf{C} \mathbb{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]) &= \sum_{j=1}^m c_j \left( \text{tr}(\mathbf{D} \tilde{\mathbf{V}}_{w_j}) + \tilde{\mathbf{w}}_{\cdot j}^\top \mathbf{D} \tilde{\mathbf{w}}_{\cdot j} \right) \\ &= \sum_{j=1}^m c_j \text{tr}(\mathbf{D}(\tilde{\mathbf{V}}_{w_j} + \tilde{\mathbf{w}}_{\cdot j} \tilde{\mathbf{w}}_{\cdot j}^\top)). \end{aligned}$$

## H.2 Deriving the ELBO expression

The evidence lower bound (ELBO) expression involves the following calculation:

$$\begin{aligned} \mathcal{L}_q(\theta) &= \int \cdots \int q(\mathbf{y}^*, \mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{y}^*, \mathbf{w} | \theta)}{q(\mathbf{y}^*, \mathbf{w})} d\mathbf{y}^* d\mathbf{w} d\theta \\ &= \mathbb{E} \left[ \overbrace{\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w} | \theta)}^{\text{joint likelihood}} \right] + \overbrace{(-\mathbb{E} [\log q(\mathbf{y}^*, \mathbf{w})])}^{\text{entropy}} \\ &= \mathbb{E} \left[ \sum_{i=1}^n \sum_{j=1}^m \log p(\overline{y_i} | y_{ij}^*) + \sum_{i=1}^n \log p(y_i^* | \boldsymbol{\alpha}, \mathbf{w}, \boldsymbol{\Psi}, \eta) + \log p(\mathbf{w} | \boldsymbol{\Psi}) \right] \\ &\quad + \sum_{i=1}^n H[q(\mathbf{y}_{i\cdot}^*)] + H[q(\mathbf{w})]. \end{aligned}$$

As discussed, given the latent propensities  $\mathbf{y}^*$ , the pdf of  $\mathbf{y}$  is degenerate and hence can be disregarded.

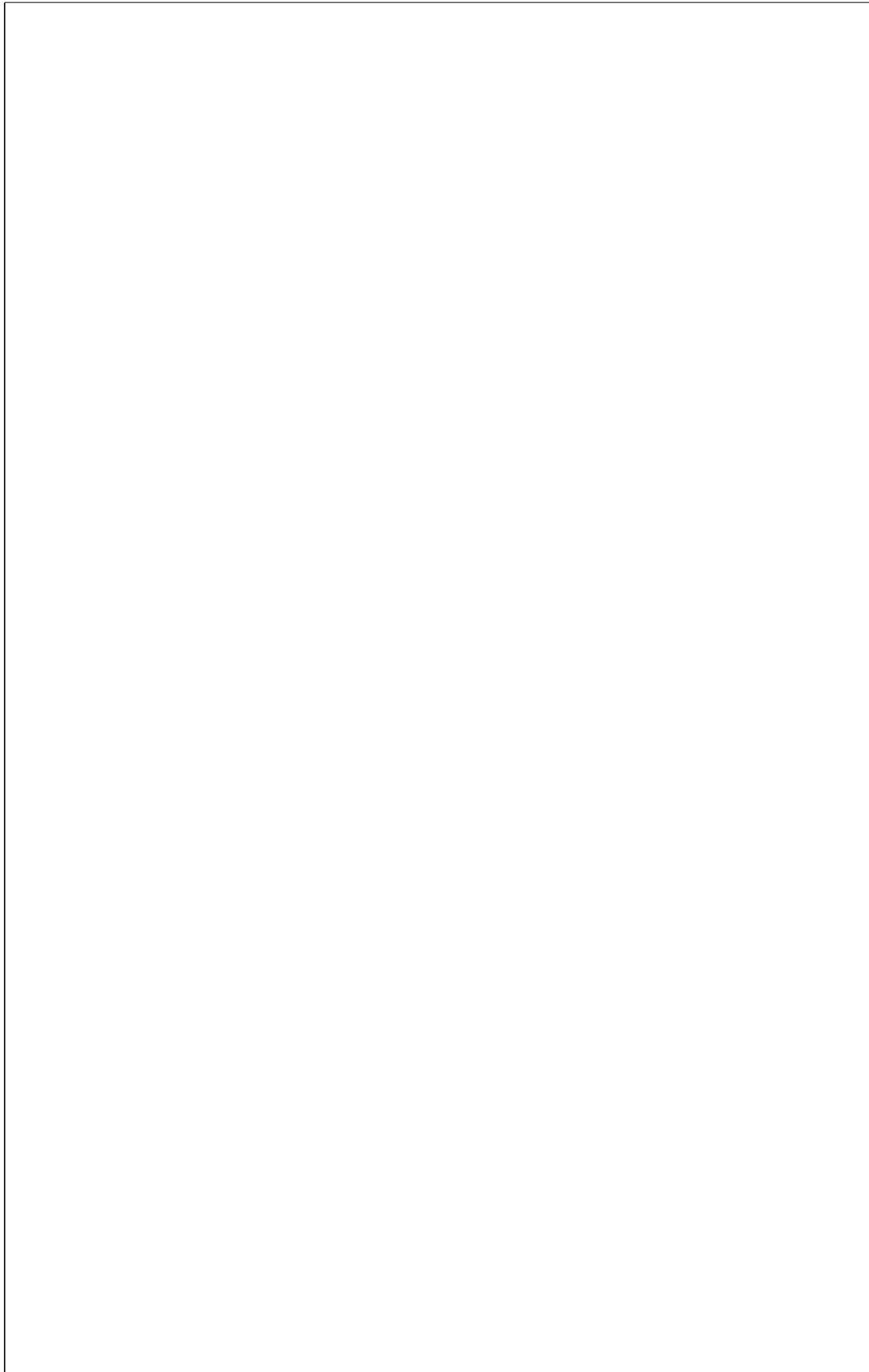
### H.2.1 Terms involving distributions of $\mathbf{y}^*$

$$\begin{aligned} &\sum_{i=1}^n \left\{ \mathbb{E} [\log p(\mathbf{y}_{i\cdot}^* | \boldsymbol{\alpha}, \mathbf{w}, \boldsymbol{\Psi}, \eta)] + H[q(\mathbf{y}_{i\cdot}^*)] \right\} \\ &= -\frac{nm}{2} \log 2\pi + \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \mathbb{E} \left[ \sum_{i=1}^n (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot}) \right] \\ &\quad + \frac{nm}{2} \log 2\pi - \frac{n}{2} \log |\boldsymbol{\Psi}| + \frac{1}{2} \mathbb{E} \left[ \sum_{i=1}^n (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot}) \right] + \log C_i \\ &= \sum_{i=1}^n \log C_i \end{aligned}$$

where  $C_i$  is the normalising constant for the distribution of multivariate truncated normal  $\mathbf{y}_{i.}^* \sim {}^t\text{N}(\tilde{\boldsymbol{\mu}}(x_i), \boldsymbol{\Psi}^{-1}, \mathcal{C}_{y_i})$ , with  $\tilde{\boldsymbol{\mu}}(x_i) = \boldsymbol{\alpha} + \tilde{\mathbf{w}}\mathbf{h}_\eta(x_i)$ .

## H.2.2 Terms involving distributions of $\mathbf{w}$

$$\begin{aligned} \mathbb{E} \log p(\mathbf{w}|\boldsymbol{\Psi}) + H[q(\mathbf{w})] &= -\frac{nm}{2} \log 2\pi - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \mathbb{E} \text{tr} (\mathbf{w} \boldsymbol{\Psi}^{-1} \mathbf{w}^\top) \\ &\quad + \frac{nm}{2} (1 + \log 2\pi) + \frac{1}{2} \log |\tilde{\mathbf{V}}_w| \\ &= \frac{nm}{2} - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \sum_{i,j=1}^m \boldsymbol{\Psi}_{ij}^{-1} \text{tr} \mathbb{E} [\tilde{\mathbf{w}}_{.j} \tilde{\mathbf{w}}_{.j}^\top] + \frac{1}{2} \log |\tilde{\mathbf{V}}_w| \end{aligned}$$



## Appendix I

# The Gibbs sampler for the I-prior Bayesian variable selection model

apx:gibbsbv  
s

The I-prior Bayesian variable selection model has the following hierarchical form:

$$\begin{aligned} \mathbf{y}|\alpha, \boldsymbol{\beta}, \gamma, \sigma^2, \kappa &\sim \text{N}_n(\alpha \mathbf{1}_n + \mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n) \\ \boldsymbol{\theta} &= (\gamma_1 \beta_1, \dots, \gamma_p \beta_p)^\top \\ \boldsymbol{\beta}|\sigma^2, \kappa &\sim \text{N}_p(\mathbf{0}, \sigma^2 \kappa \mathbf{X}^\top \mathbf{X}) \\ \alpha|\sigma^2 &\sim \text{N}(0, \sigma^2 A) \\ \sigma^2, \kappa &\sim \Gamma^{-1}(c, d) \\ \gamma_j &\sim \text{Bern}(\pi_j) \quad j = 1, \dots, p \end{aligned}$$

In the simulations and real-data examples, we used  $\pi_j = 0.5, \forall j$ ,  $A = 100$ , and  $c = d = 0.001$ , and the columns of the matrix  $\mathbf{X}$  are standardised.

The first line of the set of equations above is the likelihood, while the joint prior density is given by

$$p(\alpha, \boldsymbol{\beta}, \gamma, \sigma^2, \kappa) = p(\boldsymbol{\beta}|\sigma^2)p(\alpha|\sigma^2)p(\sigma^2)p(\kappa)p(\gamma_1) \cdots p(\gamma_p).$$

For simplicity, in the following subsections we shall denote by  $\Theta$  the entire set of parameters, while  $\Theta_{-\xi}$  implies the set of parameters excluding the parameter  $\xi$ .

## I.1 Conditional posterior for $\beta$

$$\begin{aligned}
\log p(\beta|\mathbf{y}, \Theta_{-\beta}) &= \text{const.} + \log p(\mathbf{y}|\Theta) + \log p(\beta|\sigma^2) \\
&= \text{const.} - \frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}_\gamma \beta\|^2 - \frac{1}{2\sigma^2} \beta^\top (\kappa \mathbf{X}^\top \mathbf{X})^{-1} \beta \\
&= \text{const.} - \frac{1}{2\sigma^2} \left( \beta^\top (\mathbf{X}_\gamma^\top \mathbf{X}_\gamma + (\kappa \mathbf{X}^\top \mathbf{X})^{-1}) \beta - 2(\mathbf{y} - \alpha \mathbf{1}_n)^\top \mathbf{X}_\gamma \beta \right) \\
&= \text{const.} - \frac{1}{2\sigma^2} (\beta - \tilde{\mathbf{B}}(\mathbf{y} - \alpha \mathbf{1}_n))^\top \tilde{\mathbf{B}}^{-1} (\beta - \tilde{\mathbf{B}}(\mathbf{y} - \alpha \mathbf{1}_n))
\end{aligned}$$

where  $\tilde{\mathbf{B}} = \mathbf{X}_\gamma^\top \mathbf{X}_\gamma + (\kappa \mathbf{X}^\top \mathbf{X})^{-1}$ , and  $\mathbf{X}_\gamma = (\gamma_1 X_1 \cdots \gamma_p X_p)$  is the  $n \times p$  design matrix  $\mathbf{X}$  with each of the  $p$  columns multiplied by the indicator variable  $\gamma$ . This is of course recognised as the log density of a  $p$ -variate normal distribution with mean and variance

$$\mathbb{E}[\beta|\Theta_{-\beta}] = \tilde{\mathbf{B}}(\mathbf{y} - \alpha \mathbf{1}_n) \text{ and } \text{Var}[\beta|\Theta_{-\beta}] = \sigma^2 \tilde{\mathbf{B}}.$$

## I.2 Conditional posterior for $\gamma$

Consider each  $\gamma_j$  in turn. For  $j \in \{1, \dots, p\}$ ,

$$\begin{aligned}
p(\gamma_j|\mathbf{y}, \Theta_{-\gamma_j}) &\propto p(\mathbf{y}|\Theta) p(\gamma_j) \\
&\propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}\|^2\right) \pi_j^{\gamma_j} (1 - \pi_j)^{1-\gamma_j}
\end{aligned}$$

Since the support of  $\gamma_j$  is  $\{0, 1\}$ , the above is a probability mass function which can be normalised easily. When  $\gamma_j = 1$ , we have

$$p(\gamma_j|\mathbf{y}, \Theta_{-\gamma_j}) \propto \pi_j \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}_j^{[1]}\|^2\right) := u_j$$

while for  $\gamma_j = 0$ , we have

$$p(\gamma_j|\mathbf{y}, \Theta_{-\gamma_j}) \propto (1 - \pi_j) \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}_j^{[0]}\|^2\right) := v_j.$$

For  $j = 1, \dots, p$ , we have used the notation  $\boldsymbol{\theta}_j^{[\omega]}$  to mean

$$\boldsymbol{\theta}_j^{[\omega]} = \begin{cases} (\theta_1, \dots, \theta_{j-1}, \beta_j, \theta_{j+1}, \dots, \theta_p) & \omega = 1 \\ (\theta_1, \dots, \theta_{j-1}, 0, \theta_{j+1}, \dots, \theta_p) & \omega = 0. \end{cases}$$



Therefore, the conditions distribution for  $\gamma_j$  is Bernoulli with success probability

$$\tilde{\pi}_j = \frac{u_j}{u_j + v_j}.$$

### I.3 Conditional posterior for $\alpha$

We can obtain the conditional posterior for  $\alpha$  in a similar fashion we obtained the conditional posterior for  $\beta$ . That is,

$$\begin{aligned} \log p(\alpha|\mathbf{y}, \Theta_{-\alpha}) &= \text{const.} + \log p(\mathbf{y}|\Theta) + \log p(\alpha|\sigma^2) \\ &= \text{const.} - \frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}\|^2 - \frac{\alpha^2}{2\sigma^2 A} \\ &= \text{const.} - \frac{1}{2\sigma^2} \left( (n + A^{-1})\alpha^2 - 2\alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\theta}) \right) \\ &= \text{const.} - \frac{1}{2\sigma^2(n + A^{-1})} \left( \alpha - \frac{\sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\theta})}{n + A^{-1}} \right)^2. \end{aligned}$$

Thus, the conditional posterior for  $\alpha$  is normal with mean and variance which can be easily read off the final line above.

### I.4 Conditional posterior for $\sigma^2$

The conditional density for  $\sigma^2$  is

$$\begin{aligned} \log p(\sigma^2|\mathbf{y}, \Theta_{-\sigma^2}) &= \text{const.} + \log p(\mathbf{y}|\Theta) + \log p(\sigma^2) \\ &= \text{const.} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}\|^2 - (c + 1) \log \sigma^2 - d/\sigma^2 \\ &= \text{const.} - (n/2 + c + 1) \log \sigma^2 - \frac{\|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}\|^2/2 + d}{\sigma^2} \end{aligned}$$

which is an inverse gamma distribution with shape  $\tilde{c} = n/2 + c + 1$  and scale  $\tilde{d} = \|\mathbf{y} - \alpha \mathbf{1}_n - \mathbf{X}\boldsymbol{\theta}\|^2/2 + d$ .

## I.5 Conditional posterior for $\kappa$

Interestingly, since  $\kappa$  is a hyperparameter to be estimated, it does not actually make use of any data, apart from the appearance of  $\mathbf{X}$  in the covariance matrix for  $\boldsymbol{\beta}$ .

$$\begin{aligned}\log p(\kappa|\mathbf{y}, \Theta_{-\kappa}) &= \text{const.} + \log p(\boldsymbol{\beta}|\sigma^2, \kappa) + \log p(\kappa) \\ &= \text{const.} - \frac{p}{2} \log \kappa - \frac{1}{\kappa} \cdot \frac{1}{2\sigma^2} \boldsymbol{\beta}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\beta} - (c+1) \log \kappa - d/\kappa \\ &= \text{const.} - (p/2 + c + 1) \log \kappa - \frac{\boldsymbol{\beta}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\beta} / \sigma^2 + d}{\kappa}\end{aligned}$$

This is an inverse gamma distribution with shape  $\tilde{c} = p/2 + c + 1$  and scale  $\tilde{d} = \boldsymbol{\beta}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \boldsymbol{\beta} / \sigma^2 + d$ .

## I.6 Computational note

From the above, we see that all of the Gibbs conditionals are of recognisable form, making Gibbs sampling a straightforward MCMC method to implement. We built an R package **ipriorBVS** that uses JAGS (Plummer, 2003), a variation of WinBUGS, internally for the Gibbs sampling, and wrote a wrapper function which takes formula based inputs for convenience. The **ipriorBVS** also performs two-stage BVS, and supported priors are the I-prior,  $g$ -prior, and independent prior, as used in this thesis. Although a Gibbs sampler could be coded from scratch, JAGS has the advantage of being tried and tested and has simple controls for tuning (burn-in, adaptation, thinning, etc.). Furthermore, the output from JAGS can be inspected using a myriad of multipurpose MCMC tools to diagnose convergence problems. The **ipriorBVS** package is available at <https://github.com/haziqj/ipriorBVS>.

In all examples, a default setting of 4,000 burn-in samples, 1,000 adaptation size, and 10,000 samples with no thinning seemed adequate. There were no major convergence issues encountered.

Computational complexity is dominated by the inversion of a  $p \times p$  matrix, and matrix multiplications of order  $O(np^2)$ . These occur in the conditional posterior for  $\boldsymbol{\beta}$ . Overall, if  $n \gg p$ , then time complexity is  $O(np^2)$ . Storage requirements are  $O(np)$ .

# Bibliography

- chopin2011fast Chopin, Nicolas (2011). “Fast simulation of truncated Gaussian distributions”. In: *Statistics and Computing* 21.2, pp. 275–288.
- damien2001sampling Damien, Paul and Stephen G Walker (2001). “Sampling truncated normal, beta, and gamma densities”. In: *Journal of Computational and Graphical Statistics* 10.2, pp. 206–215.
- marsaglia2000ziggurat Marsaglia, George and Wai Wan Tsang (2000). “The ziggurat method for generating random variables”. In: *Journal of statistical software* 5.8, pp. 1–7.
- petersen2008matrix Petersen, Kaare Brandt and Michael Syskind Pedersen (2008). “The matrix cookbook”. In: *Technical University of Denmark* 7.15, p. 510.
- plummer2003jags Plummer, Martyn (2003). “JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling”. In: *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*. Vol. 124. Vienna, Austria, p. 125.
- rasmussen2006gaussian Rasmussen, Carl Edward and Christopher K I Williams (2006). *Gaussian Processes for Machine Learning*. The MIT Press.
- robert1995simulation Robert, Christian P (1995). “Simulation of truncated normal variables”. In: *Statistics and computing* 5.2, pp. 121–125.
- zellner1986assessing Zellner, Arnold (1986). “On assessing prior distributions and Bayesian regression analysis with g-prior distributions”. In: *Bayesian inference and decision techniques*.
- zhang2013kronecker Zhang, Huamin and Feng Ding (2013). “On the Kronecker products and their applications”. In: *Journal of Applied Mathematics* 2013.

# Figures

# Tables

# Theorems

C.1	Lemma (Properties of multivariate normal) . . . . .	11
C.4	Lemma (Equivalence between matrix and multivariate normal) . . . . .	15

# Definitions

A.1	Definition (Functional derivative) . . . . .	5
B.1	Definition (Kronecker product) . . . . .	9
B.2	Definition (Vectorisation) . . . . .	10





# Nomenclature

As much as possible, and unless otherwise stated, the following conventions are used throughout this thesis.

## Conventions

$\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots$	Boldface lower case letters denote real vectors
$\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$	Boldface upper case letters denote real matrices
$\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$	Calligraphic upper case letters denote sets
$x'$	Primes are used to distinguish elements (not indicate derivatives)
$\hat{\theta}$	Hats are used to denote estimators of parameters

## Indexing

$\mathbf{A}_{ij}, A_{ij}, a_{ij}$	The $(i, j)$ 'th element of the matrix $\mathbf{A}$
$\mathbf{A}_i.$	The $i$ 'th row of the matrix $\mathbf{A}$ as a tall vector (transposed row vector)
$\mathbf{A}.j$	The $j$ 'th column vector of the matrix $\mathbf{A}$

## Symbols

$\mathbb{N}$	The set of natural numbers (excluding zero)
$\mathbb{Z}$	The set of integers
$\mathbb{R}$	The set of real numbers
$\mathbb{R}_{>0}$	The set of positive real numbers, $\{x \in \mathbb{R}   x > 0\}$
$\mathbb{R}_{\geq 0}$	The set of non-negative real numbers, $\{x \in \mathbb{R}   x \geq 0\}$
$\mathbb{R}^d$	The $d$ -dimensional Euclidean space
$\mathcal{A}^c$	The complement of a set $\mathcal{A}$
$\mathcal{P}(\mathcal{A})$	The power set of the set $\mathcal{A}$
$\{\}, \emptyset$	The empty set
$\mathbf{0}$	A vector of zeroes
$\mathbf{1}_n$	A length $n$ vector of ones
$\mathbf{I}_n$	The $n \times n$ identity matrix
$\exists$	(short hand) There exists
$\forall$	(short hand) For all
$\lim_{n \rightarrow \infty}$	The limit as $n$ tends to infinity
$\xrightarrow{\text{dist.}}$	Convergence in distribution
$\mathcal{O}(n)$	Computational complexity (time or storage)
$\Delta x$	A quantity representing a change in $x$

**Relations**

$a \approx b$	$a$ is approximately or almost equal to $b$
$a \propto b$	$a$ is equivalent to $b$ up to a constant of proportionality
$a \equiv b$	$a$ is identical to $b$
$A \Rightarrow B$	The statement $B$ being true is predicated on $A$ being true
$A \Leftrightarrow B$	The statement $A$ is true if and only if $B$ is true
$a \in \mathcal{A}$	$a$ is an element of the set $\mathcal{A}$
$\mathcal{A} \subseteq \mathcal{B}$	$\mathcal{A}$ is a subset of $\mathcal{B}$ which may include itself
$\mathcal{A} \subset \mathcal{B}$	$\mathcal{A}$ is a subset of $\mathcal{B}$ which does not include itself
$a := b, a \leftarrow b$	$a$ is assigned the value $b$
$X \sim p(X)$	The random variable $X$ is distributed according to the pdf $p(X)$
$X \sim D$	The random variable $X$ is distributed according to the pdf specified by the distribution $D$ , e.g. $D \equiv \mathcal{N}(0, 1)$
$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} D$	Each random variable $X_i, i = 1, \dots, n$ is independently and identically distributed according to the pdf specified by the distribution $D$
$X Y$	The (random) variable $X$ given/conditional on $Y$

**Functions**

$\inf \mathcal{A}$	The infimum of a set $\mathcal{A}$
$\sup \mathcal{A}$	The supremum of a set $\mathcal{A}$
$\min \mathcal{A}$	The minimum value of a set $\mathcal{A}$
$\max \mathcal{A}$	The maximum value of a set $\mathcal{A}$
$\arg \min_x f(x)$	The value of $x$ which minimises the function $f(x)$
$\arg \max_x f(x)$	The value of $x$ which maximises the function $f(x)$
$ a $ with $a \in \mathbb{R}$	The absolute value of $a$ ; $ a  = a$ if $a$ is positive, and $-a$ if $a$ is negative, and $ 0  = 0$
$\delta_{xx'}$	The Kronecker delta; $\delta_{xx'} = 1$ if $x = x'$ , and 0 otherwise
$[A]$	The Iverson bracket; $[A] = 1$ if the logical proposition $A$ is true, and 0 otherwise
$\mathbb{1}_{\mathcal{A}}(x)$	The indicator function; $\mathbb{1}_{\mathcal{A}}(x) = 1$ if $x \in \mathcal{A}$ , and 0 otherwise
$e^x, \exp(x)$	The natural exponential function
$\log(x)$	The natural logarithmic function
$\frac{d}{dx} f(x), \dot{f}(x)$	The derivative of $f$ with respect to $x$
$f \circ g$	Composition of functions, i.e. $g$ following $f$

**Abstract vector space operations and notations**

$\mathcal{V}^\perp$	The orthogonal complement of the space $\mathcal{V}$
$\mathcal{V}^\vee$	The algebraic dual space of $\mathcal{V}$
$\mathcal{V}^*$	The continuous dual space of $\mathcal{V}$
$\overline{\mathcal{V}}$	The closure of the space $\mathcal{V}$
$\mathcal{B}(\mathcal{V})$	The Borel $\sigma$ -algebra of $\mathcal{V}$
$L^p(\mathcal{X}, \nu)$	The set of $p$ -integrable functions over the space $\mathcal{X}$ with measure $\nu$
$L(\mathcal{V}; \mathcal{W})$	The set of bounded, linear operators from $\mathcal{V}$ to $\mathcal{W}$
$\dim(\mathcal{V})$	The dimensions of the vector space $\mathcal{V}$
$\langle x, y \rangle_{\mathcal{V}}$	The inner product between $x$ and $y$ in the vector space $\mathcal{V}$

$\ x\ _{\mathcal{V}}$	The norm of $x$ in the vector space $\mathcal{V}$
$D(x, y)$	The distance between $x$ and $y$
$x \otimes y$	The tensor product of $x$ and $y$ which are elements of a vector space
$\mathcal{F} \otimes \mathcal{G}$	The tensor product space of two vector spaces
$\mathcal{F} \oplus \mathcal{G}$	The direct sum (or tensor sum) of two vector spaces
$df(x), d^2f(x)$	The first and second Fréchet differentials of $f$ at $x$
$\partial_v f(x), \partial_v^2 f(x)$	The first and second Gâteaux differentials of $f$ at $x$ in the direction $v$
$\nabla f(x), \nabla^2 f(x)$	The gradient and Hessian of $f$ at $x$ in the direction $v$ ( $f$ is a mapping of a Hilbert space)

### Matrix and vector operations

$\mathbf{a}^\top, \mathbf{A}^\top$	The transpose of a vector $\mathbf{a}$ or matrix $\mathbf{A}$
$\mathbf{A}^{-1}$	The inverse of a square matrix $\mathbf{A}$
$\ \mathbf{a}\ ^2$	The squared 2-norm the vector $\mathbf{a}$ , equivalent to $\mathbf{a}^\top \mathbf{a}$
$ \mathbf{A} $	The determinant of a matrix $\mathbf{A}$
$\text{tr}(\mathbf{A})$	The trace of a square matrix $\mathbf{A}$
$\text{diag}(\mathbf{A})$	The diagonal elements of a square matrix $\mathbf{A}$
$\text{rank}(\mathbf{A})$	The rank of a matrix $\mathbf{A}$
$\text{vec}(\mathbf{A})$	The column-wise vectorisation of a matrix $\mathbf{A}$
$\mathbf{a} \otimes \mathbf{b}$	The outer product of two vectors $\mathbf{a}$ and $\mathbf{b}$
$\mathbf{A} \otimes \mathbf{B}$	The Kronecker product of matrix $\mathbf{A}$ with matrix $\mathbf{B}$
$\mathbf{A} \circ \mathbf{B}$	The Hadamard product two matrices $\mathbf{A}$ and $\mathbf{B}$

### Statistical functions

$P(A)$	The probability of event $A$ occurring
$p(X \theta)$	The probability density function of $X$ given parameters $\theta$
$L(\theta X)$	The log-likelihood of $\theta$ given data $X$ , sometimes simply $L(\theta)$
$\text{BF}(M, M')$	Bayes factor for comparing two models $M$ and $M'$
$\mathcal{I}(\theta)$	The Fisher information for $\theta$
$\mathbb{E}[X], \mathbb{E} X$	The expectation <sup>1</sup> of the random element $X$
$\text{Var}[X], \text{Var} X$	The variance <sup>1</sup> of the random element $X$
$\text{Cov}[X, Y]$	The covariance <sup>1</sup> between two random elements $X$ and $Y$
$H(p)$	The entropy of the distribution $p(X)$
$\text{D}_{\text{KL}}(q(x)  p(x))$	The Kullback-Leibler divergence from $p(x)$ to $q(x)$ , denoted also by $\text{D}_{\text{KL}}(q  p)$ for short

### Statistical distributions

$N(\mu, \sigma^2)$	Univariate normal distribution with mean $\mu$ and variance $\sigma^2$
$N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	$d$ -dimensional multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$
$\phi(z)$	The standard normal pdf
$\Phi(z)$	The standard normal cdf

<sup>1</sup>When there is ambiguity as to which random element the expectation or variance is taken under or what its distribution is, this is explicated by means of subscripting, e.g.  $\mathbb{E}_{X \sim N(0,1)} X$  to denote the expectation of a standard normal random variable.

$\phi(x \mu, \sigma^2)$	The pdf of $N(\mu, \sigma^2)$
$\phi(\mathbf{x} \boldsymbol{\mu}, \boldsymbol{\Sigma})$	The pdf of $N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
$MN_{n,m}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\Psi})$	Matrix normal distribution with mean $\boldsymbol{\mu}$ and row variances $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ and column variances $\boldsymbol{\Psi} \in \mathbb{R}^{m \times m}$
${}^tN(\mu, \sigma^2, a, b)$	Truncated univariate normal distribution with mean $\mu$ and variance $\sigma^2$ restricted to the interval $(a, b)$
$N_+(0, 1)$	The half-normal distribution with variance $\sigma^2$
$N_+(0, \sigma^2)$	The folded-normal distribution with variance $\sigma^2$
${}^tN_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathcal{A})$	Truncated $d$ -dimensional multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ restricted to the set $\mathcal{A}$
$\Gamma(s, r)$	Gamma distribution with shape $s$ and rate $r$ parameters
$\Gamma^{-1}(s, \sigma)$	Inverse gamma distribution with shape $s$ and scale $\sigma$ parameters
$\chi_d^2$	Chi-squared distribution with $d$ degrees of freedom
$\text{Bern}(p)$	Bernoulli distribution with probability of success $p$
$\text{Cat}(p_1, \dots, p_m)$	Categorical distribution with $m$ categories, and each category has probability of success $p_j$

# Abbreviations

This document is incomplete. The external file associated with the glossary ‘abbreviations’ (which should be called `appendix.gls-abr`) hasn’t been created.

Check the contents of the file `appendix.gls-abr`. If it’s empty, that means you haven’t indexed any of your entries in this glossary (using commands like `\gls` or `\glsadd`) so this list can’t be generated. If the file isn’t empty, the document build process hasn’t been completed.

Try one of the following:

- Add `automake` to your package option list when you load `glossaries-extra.sty`. For example:

```
\usepackage[automake]{glossaries-extra}
```

- Run the external (Lua) application:

```
makeglossaries-lite.lua "appendix"
```

- Run the external (Perl) application:

```
makeglossaries "appendix"
```

Then rerun  $\text{\LaTeX}$  on this document.

This message will be removed once the problem has been fixed.