

# To-do list

1. Exponential family for $y$ not really necessary, it just follows nicely from the latent variable motivation. . . . .	2
2. Section X . . . . .	10
3. Section X . . . . .	14
4. Compare: Laplace, variational and HMC. . . . .	17
5. Section 4.3.3 . . . . .	22
6. Section 4.4 . . . . .	24
7. Chapter 3 . . . . .	34
8. Lemma X . . . . .	44
9. Section 4.3.3 . . . . .	48
10. expression XXX . . . . .	50
11. page . . . . .	52

# Contents

<b>5 I-priors for categorical responses</b>	<b>2</b>
5.1 A naïve model . . . . .	5
5.2 A latent variable motivation: the I-probit model . . . . .	8
5.3 Identifiability and IIA . . . . .	11
5.4 Estimation . . . . .	14
5.4.1 Laplace approximation . . . . .	14
5.4.2 Variational inference . . . . .	16
5.4.3 Markov chain Monte Carlo methods . . . . .	17
5.4.4 Comparison of estimation methods . . . . .	17
5.5 A variational algorithm . . . . .	18
5.5.1 Latent propensities $\mathbf{y}^*$ . . . . .	19
5.5.2 I-prior random effects $\mathbf{w}$ . . . . .	21

5.5.3	Kernel parameters $\eta$ . . . . .	21
5.5.4	Intercepts $\alpha$ . . . . .	22
5.5.5	The CAVI algorithm . . . . .	23
5.6	Post-estimation . . . . .	23
5.7	Computational consideration . . . . .	26
5.7.1	Efficient computation of class probabilities . . . . .	27
5.7.2	Computational complexity of the CAVI algorithm . . . . .	29
5.7.3	Difficulties faced with estimating $\Psi$ . . . . .	30
5.8	Examples . . . . .	30
5.9	Discussion . . . . .	30
5.10	Miscellanea . . . . .	31
5.10.1	A brief introduction to variational inference . . . . .	31
5.10.2	The EM algorithm is intractable—variational Bayes EM . . . . .	33
5.11	Some distributions and their properties . . . . .	33
5.11.1	Multivariate normal distribution . . . . .	34
5.11.2	Matrix normal distribution . . . . .	35
5.11.3	Truncated univariate normal distribution . . . . .	38
5.11.4	Truncated multivariate normal distribution . . . . .	39
5.12	Derivation of the CAVI algorithm . . . . .	42
5.12.1	Derivation of $\tilde{q}(\mathbf{y}^*)$ . . . . .	44
5.12.2	Derivation of $\tilde{q}(\mathbf{w})$ . . . . .	45
5.12.3	Derivation of $\tilde{q}(\eta)$ . . . . .	48
5.12.4	Derivation of $\tilde{q}(\Psi)$ . . . . .	50
5.12.5	Derivation of $\tilde{q}(\alpha)$ . . . . .	52
5.13	Deriving the ELBO expression . . . . .	53
5.13.1	Terms involving distributions of $\mathbf{y}^*$ . . . . .	53
5.13.2	Terms involving distributions of $\mathbf{w}$ . . . . .	55
5.13.3	Terms involving distributions of $\eta$ . . . . .	55
5.13.4	Terms involving distribution of $\alpha$ . . . . .	56
5.13.5	ELBO summarised . . . . .	56
<b>Bibliography</b>		<b>59</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 5

# I-priors for categorical responses

In a regression setting, consider polytomous response variables  $y_1, \dots, y_n$ , where each  $y_i$  takes on exactly one of the values  $\{1, \dots, m\}$  from a set of  $m$  possible choices. Modelling categorical response variables is of profound interest in statistics, econometrics and machine learning, with applications aplenty. In the social sciences, categorical variables often arise from survey responses, and one may be interested in studying correlations between explanatory variables and the categorical response of interest. Economists are often interested in discrete choice models to explain and predict choices between several alternatives, such as consumer choice of goods or modes of transport. In this age of big data, machine learning algorithms are used for classification of observations based on what is usually a large set of variables or features.

As an extension to the I-prior methodology, we propose a flexible modelling framework suitable for regression of categorical response variables.

In the spirit of generalised linear models (McCullagh and Nelder, 1989), we relate the class probabilities of the observations to the I-prior regression model via a link function. Perhaps though, it is more intuitive to view it as machine learners do: Since the regression function is ranged on the entire real line, it is necessary to “squash” it through some sigmoid function to conform it to the interval  $[0, 1]$  suitable for probability measures. As in GLMs, the  $y_i$ ’s are assumed to follow an exponential family distribution, and in this case, the categorical distribution. We denote this by

$$y_i \sim \text{Cat}(p_{i1}, \dots, p_{im}),$$

1. Exponential family for  $y$  not really necessary, it just follows nicely from the latent variable motivation.

with the class probabilities satisfying  $p_{ij} \geq 0, \forall j = 1, \dots, m$  and  $\sum_{j=1}^m p_{ij} = 1$ . The probability mass function (PMF) of  $y_i$  is given by

$$p(y_i) = p_{i1}^{[y_i=1]} \dots p_{im}^{[y_i=m]} \quad (5.1)$$

where the notation  $[\cdot]$  refers to the Iverson bracket<sup>1</sup>. The dependence of the class probabilities on the covariates is specified through the relationship

$$g(p_{ij}) = (\alpha_j + f_j(x_i))_{j=1}^m$$

where  $g : [0, 1] \rightarrow \mathbb{R}^m$  is some specified link function. As we will see later, the normality assumption of the errors naturally implies a *probit* link function, i.e.,  $g$  is the inverse cumulative distribution function (CDF) of a standard normal distribution (or more precisely, a function that *involves* the standard normal CDF). Normality is also a required assumption for I-priors to be specified on the regression functions. We call this method of probit regression using I-priors the *I-probit* regression model.

Note that the probabilities are modelled per class  $j \in \{1, \dots, m\}$  by individual regression curves  $f_j$ , and in the most general setting,  $m$  sets of intercepts  $\alpha_j$  and kernel hyperparameters  $\eta_j$  must be estimated. The dependence of these  $m$  curves are specified through covariances  $\sigma_{jk} := \text{Cov}[\epsilon_{ij}, \epsilon_{ik}]$ , for each  $j, k \in \{1, \dots, m\}$  and  $j \neq k$ . While it may be of interest to estimate these covariances, this paper considers cases where the regression functions are class independent, i.e.  $\sigma_{jk} = 0, \forall j \neq k$ . This violates the independence of irrelevant alternatives (IIA) assumption (see Section 5.3 for details) crucial in choice models, but not so much necessary for classification when the alternatives are distinctively different.

The many advantages of the I-prior methodology of [Jamil and Bergsma, 2017](#) transfer over quite well to the I-probit model for classification and inference. In particular, by choosing appropriate RKHSs for the regression functions, we are able to fit a multitude of binary and multinomial models, including multilevel or random-effects models, linear and non-linear classification models, and even spatio-temporal models. Examples of these models applied to real-world data is shown in Section ???. Working in a Bayesian setting together with variational inference allows us to estimate the model much faster than traditional MCMC sampling methods, yet provides us with the conveniences that come with posterior estimates. For example, inferences around log-odds is usually cumbersome

---

<sup>1</sup> $[A]$  returns 1 if the proposition  $A$  is true, and 0 otherwise. The Iverson bracket is a generalisation of the Kronecker delta.

for probit models, but a credibility interval can easily be obtained by resampling methods from the relevant posteriors, which are normally exponential family distributions in the I-probit model.

## 5.1 A naïve model

The I-prior methodology can be used naïvely to fit a categorical regression model. Suppose, as before, we observe data  $\{(y_1, x_1), \dots, (y_n, x_n)\}$  where each  $x_i \in \mathcal{X}$ , for  $i = 1, \dots, n$ . Here, the responses are categorical  $y_i \in \{1, \dots, m\} =: \mathcal{M}$ , and additionally, write  $y_i = (y_{i1}, \dots, y_{im})$  where the class responses  $y_{ij}$  equal one if individual  $i$ 's response category is  $y_i = j$ , and zero otherwise. In other words, there is exactly a single ‘1’ at the  $j$ 'th position in the vector  $y_i = (y_{i1}, \dots, y_{im})$ , and zeroes everywhere else. For  $j = 1, \dots, m$ , we model

$$\begin{aligned} y_{ij} &= \alpha + \alpha_j + f_j(x_i) + \epsilon_{ij} \\ (\epsilon_{i1}, \dots, \epsilon_{im})^\top &\stackrel{\text{iid}}{\sim} N_m(\mathbf{0}, \Psi^{-1}). \end{aligned} \tag{5.2}$$

{eq:naiveclassmod}

The idea here being that we attempt to model the class responses  $y_{ij}$  using class-specific regression functions  $f_j$ , and the class responses are assumed to be independent among individuals, but may or may not be correlated among classes for each individual. The class correlations are manifest themselves in the variance of the errors  $\Psi^{-1}$ , which is an  $m \times m$  matrix.

Denote the regression function  $f$  in (5.2) on the set  $\mathcal{X} \times \mathcal{M}$  as  $f(x_i, j) = \alpha_j + f_j(x_i)$ . This regression function can be seen as an ANOVA decomposition of the spaces  $\mathcal{F}_{\mathcal{M}}$  and  $\mathcal{F}_{\mathcal{X}}$  of functions over  $\mathcal{M}$  and  $\mathcal{X}$  respectively. That is,  $\mathcal{F} = \mathcal{F}_{\mathcal{M}} \oplus (\mathcal{F}_{\mathcal{M}} \otimes \mathcal{F}_{\mathcal{X}})$  is a decomposition into the main effects of ‘class’, and an interaction effect of the covariates for each class. Let  $\mathcal{F}_{\mathcal{M}}$  and  $\mathcal{F}_{\mathcal{X}}$  be RKHSs respectively with kernels  $a : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$  and  $b : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Then, the ANOVA RKKS  $\mathcal{F}$  possesses the reproducing kernel  $h : (\mathcal{X} \times \mathcal{M})^2 \rightarrow \mathbb{R}$  as defined by

$$b_\eta((x, j), (x', j')) = a(j, j') + a(j, j')h_\eta(x, x'). \tag{5.3}$$

{eq:anovaclass}

The kernel  $h_\eta$  may be any of the kernels described in this thesis, ranging from the linear kernel, to the fBm kernel, or even an ANOVA kernel. Choices for  $a : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$  include

1. **The Pearson kernel** (as defined in Definition 2.34). With  $J \sim P$ , a probability measure over  $\mathcal{M}$ ,

$$a(j, j') = \frac{\delta_{jj'}}{P(J = j)} - 1.$$

2. **The identity kernel.** With  $\delta$  denoting the Kronecker delta function,

$$a(j, j') = \delta_{jj'}.$$

The purpose of either of these kernels is to contribute to the class intercepts  $\alpha_j$ , and to associate a regression function in each class. We have a slight preference for the identity kernel, which lends itself as being easy to handle computationally. The only difference between the two is the inverse probability weighting per class that is applied in the Pearson kernel, but not in the identity kernel.

As a remark, 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 2.37](#)). 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,

$$\begin{aligned} \sum_{j=1}^m y_{ij} &= \sum_{j=1}^m (\alpha + \alpha_j + f_j(x_i)) \\ &= m\alpha + \sum_{j=1}^m \alpha_j + \sum_{j=1}^m f_j(x_i) \end{aligned}$$

and since  $\sum_{j=1}^m y_{ij} = 1$ , we have that  $\alpha = 1/m$  and can thus be fixed to resolve identification. The Pearson RKHS will contain zero mean functions, but the RKHS of constant functions induced by the identity kernel may not. If this is the case, then it should be ensured that  $\sum_{j=1}^m \alpha_j = 0$  in other ways; perhaps during the estimation process.

With  $f \in \mathcal{F}$  the RKKS with kernel  $h_\eta$ , it is straightforward to assign an I-prior on  $f$ . It is in fact

$$\begin{aligned} f(x_i, j) &= \sum_{j'=1}^m \sum_{i'=1}^n a(j, j') (1 + h_\eta(x_i, x_{i'})) w_{i'j'} \\ (w_{i'1}, \dots, w_{i'm})^\top &\sim N_m(\mathbf{0}, \Psi) \end{aligned} \tag{5.4}$$

{eq:naivecl  
assiprior}

assuming a zero prior mean  $f_0(x, j) = 0$ . It is much 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  $\Psi$ . This means that  $\mathbf{w} \sim \text{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \Psi)$  which implies  $\text{vec } \mathbf{w} \sim N_{nm}(\mathbf{0}, \Psi \otimes \mathbf{I}_n)$ , and similarly,

$\epsilon \sim N_{nm}(\mathbf{0}, \Psi^{-1} \otimes \mathbf{I}_n)$ . Denote by  $\mathbf{H}_\eta$  the  $n \times n$  kernel matrix with entries supplied by  $1 + h_\eta$ , and  $\mathbf{A}$  the  $m \times m$  matrix with entries supplied by  $a$ . From (5.4), we have that

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

and thus  $\text{vec } \mathbf{f} \sim N_{nm}(\mathbf{0}, \mathbf{A} \Psi \mathbf{A} \otimes \mathbf{H}_\eta^2)$ . As  $\mathbf{y} = \boldsymbol{\alpha} + \mathbf{f} + \epsilon$ , where  $\boldsymbol{\alpha} \in \mathbb{R}^{n \times m}$  with  $(i, j)$  entries given by  $\alpha + \alpha_j = \alpha_j + 1/m$ , by linearity we have that

$$\text{vec } \mathbf{y} \sim N_{nm}(\text{vec } \boldsymbol{\alpha}, (\mathbf{A} \Psi \mathbf{A} \otimes \mathbf{H}_\eta^2) + (\Psi^{-1} \otimes \mathbf{I}_n)) \quad (5.5)$$

and

$$\text{vec } \mathbf{y} | \text{vec } \mathbf{w} \sim N_{nm}(\text{vec}(\boldsymbol{\alpha} + \mathbf{H}_\eta \mathbf{w} \mathbf{A}), (\Psi^{-1} \otimes \mathbf{I}_n)). \quad (5.6)$$

which can then be estimated using the methods described in Chapter 4.

When using the identity kernel in conjunction with an assumption of iid errors ( $\Psi = \psi \mathbf{I}_n$ ), the above distributions simplify further. Specifically, the variance in the marginal distribution becomes

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

which implies independence and identical variances  $\mathbf{V}_y$  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  $\text{Var}(\text{vec } \mathbf{f}) = \text{diag}(\psi \mathbf{H}_\eta^2, \dots, \psi \mathbf{H}_\eta^2)$ , which could be interpreted as having independent and identical I-priors on the regression functions for each class  $\mathbf{f}_j = (f(x_1, j), \dots, f(x_n, j))^\top$ .

There are several downfalls to using the model described above. Unlike in the case of continuous response variables, the normal I-prior model is highly inappropriate for categorical responses. For one, it violates the normality and homoscedasticity assumptions of the errors. For another, predicted values may be out of the range  $[0, m]$  and thus poorly calibrated. Furthermore, it would be more suitable if the class probabilities—the probability of an observation belonging to a particular class—were also part of the model. In the next section, we propose an improvement to this naïve I-prior classification model by considering a probit-like transformation of the regression functions.



## 5.2 A latent variable motivation: the I-probit model

It is convenient, as we did in the previous subsection, to again think of the responses  $y_i \in \{1, \dots, m\} = \mathcal{M}$  as comprising of a binary vector  $(y_{i1}, \dots, y_{im})$ , with a single ‘1’ at the position corresponding to the value that  $y_i$  takes. In this formulation, each  $y_{ij}$  is distributed as Bernoulli with probability  $p_{ij}$ . Now, assume that, for each  $y_i = (y_{i1}, \dots, y_{im})$ , there exists corresponding *continuous, underlying, latent variables*  $y_{i1}^*, \dots, y_{im}^*$  such that

$$y_i = \begin{cases} 1 & \text{if } y_{i1}^* \geq y_{i2}^*, y_{i3}^*, \dots, y_{im}^* \\ 2 & \text{if } y_{i2}^* \geq y_{i1}^*, y_{i3}^*, \dots, y_{im}^* \\ \vdots & \\ m & \text{if } y_{im}^* \geq y_{i2}^*, y_{i3}^*, \dots, y_{i,m-1}^*. \end{cases} \quad (5.7)$$

{eq:latentmodel}

In other words,  $y_{ij} = \arg \max_{k=1}^m y_{ik}^*$ . Such a formulation is common in economic choice models, and is rationalised by a utility-maximisation argument: an agent faced with a choice from a set of alternatives will choose the one which benefits them most. In this sense, the  $y_{ij}^*$ ’s represent individual  $i$ ’s *latent propensities* for choosing alternative  $j$ .

Instead of modelling the observed  $y_{ij}$ ’s directly, we model instead the  $n$  latent variables in each class  $j = 1, \dots, m$  according to the regression problem

$$\begin{aligned} y_{ij}^* &= \alpha + \alpha_j + f_j(x_i) + \epsilon_{ij} \\ (\epsilon_{i1}, \dots, \epsilon_{im})^\top &\stackrel{\text{iid}}{\sim} \mathbf{N}_m(\mathbf{0}, \mathbf{\Psi}^{-1}). \end{aligned} \quad (5.8)$$

{eq:multinomial-latent}

We can see some semblance of this model with the one in (5.4), and ultimately the aim is to assign I-priors to the regression function of these latent variables, and we will describe this shortly. For now, write  $\boldsymbol{\mu}(x_i) \in \mathbb{R}^m$  whose  $j$ ’th component is  $\alpha + \alpha_j + f_j(x_i)$ , and realise that each  $(y_{i1}^*, \dots, y_{im}^*)^\top$  has the distribution  $\mathbf{N}_m(\boldsymbol{\mu}(x_i), \mathbf{\Psi}^{-1})$ , conditional on the data  $x_i$ , the intercepts  $\alpha, \alpha_1, \dots, \alpha_m$ , the evaluations of the functions at  $x_i$  for each class  $f_1(x_i), \dots, f_m(x_i)$ , and the error covariance matrix  $\mathbf{\Psi}^{-1}$ .

The probability of belonging to class  $j$  for observation  $i$ , i.e.  $p_{ij}$  is calculated as

$$\begin{aligned}
 p_{ij} &= \mathbb{P}(y_i = j) \\
 &= \mathbb{P}(\{y_{ij}^* > y_{ik}^* \mid \forall k \neq j\}) \\
 &= \int \cdots \int_{\{y_{ij}^* > y_{ik}^* \mid \forall k \neq j\}} \phi(y_{i1}^*, \dots, y_{im}^* \mid \boldsymbol{\mu}(x_i), \boldsymbol{\Psi}^{-1}) dy_{i1}^* \cdots dy_{im}^*, \tag{5.9}
 \end{aligned}$$

{eq:p<sub>ij</sub>}

where  $\phi(\cdot \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is the density of the multivariate normal with mean  $\boldsymbol{\mu}$  and variance  $\boldsymbol{\Sigma}$ . This is the probability that the normal random variable  $\mathbf{y}_i^*$  belongs to the set  $\{y_{ij}^* > y_{ik}^* \mid \forall k \neq j\}$ , which are cones in  $\mathbb{R}^m$ . Since the union of these cones is the entire  $m$ -dimensional space of reals, the probabilities add up to one and hence they represent a proper probability mass function of the classes. While this does not have a closed-form expression and highlights one of the difficulties of working with probit models, the integral is by no means impossible to compute—see [Section 5.7.1](#) for a note regarding this matter.

Now we'll see how to specify an I-prior on the regression problem [\(5.8\)](#). In the naïve I-prior model, we wrote  $f(x_i, j) = \alpha_j + f_j(x_i)$ , and specified for  $f$  to belong to an ANOVA RKKS with kernel defined in [\(5.3\)](#). Instead of doing the same, we take a different approach. Treat the  $\alpha_j$ 's in [\(5.8\)](#) as intercept parameters to estimate with the additional requirement that  $\sum_{j=1}^m \alpha_j = 0$ . Further, let  $\mathcal{F}$  be a (centred) RKHS/RKKS of functions over  $\mathcal{X}$  with reproducing kernel  $h_\eta$ . Now, consider putting an I-prior on the regression functions  $f_j \in \mathcal{F}$ ,  $j = 1 \dots, m$ , defined by

$$f_j(x_i) = f_0(x_i) + \sum_{k=1}^n h_\eta(x_i, x_k) w_{ik}$$

with  $\mathbf{w}_i := (w_{i1}, \dots, w_{im})^\top \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \boldsymbol{\Psi})$ . This is similar to the naïve I-prior specification [\(5.4\)](#), except that the intercept have been treated as parameters rather than accounting for them using an RKHS of constant functions. In particular, the overall regression relationship still satisfies the ANOVA functional decomposition. We find that this method bodes well down the line computationally.

We call the multinomial probit regression model of [\(5.7\)](#) subject to [\(5.8\)](#) and I-priors on  $f_j \in \mathcal{F}$ , the *I-probit model*. For completeness, this is stated again: for  $i = 1, \dots, n$ ,

$y_i = \arg \max_{k=1}^m y_{ik}^* \in \{1, \dots, m\}$ , where, for  $j = 1, \dots, m$ ,

$$y_{ij}^* = \alpha + \alpha_j + f_0(x_i) + \overbrace{\sum_{k=1}^n h_\eta(x_i, x_k) w_{ik}}^{f_j(x_i)} + \epsilon_{ij} \quad (5.10)$$

$$(\epsilon_{i1}, \dots, \epsilon_{im})^\top \stackrel{\text{iid}}{\sim} N_m(\mathbf{0}, \Psi^{-1})$$

$$(w_{i1}, \dots, w_{im})^\top \stackrel{\text{iid}}{\sim} N_m(\mathbf{0}, \Psi).$$

{eq:iprobit  
mod}

The parameters of the I-probit model are denoted by  $\theta = \{\alpha_1, \dots, \alpha_m, \eta, \Psi\}$ . To establish notation, let

- $\epsilon \in \mathbb{R}^{n \times m}$  denote the matrix containing  $(i, j)$  entries  $\epsilon_{ij}$ , whose rows are  $\epsilon_i$ . and columns are  $\epsilon_{.j}$ . Its distribution is  $\epsilon \sim \text{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \Psi^{-1})$ ;
- $\mathbf{w} \in \mathbb{R}^{n \times m}$  denote the matrix containing  $(i, j)$  entries  $w_{ij}$ , whose rows are  $\mathbf{w}_i$ . and columns are  $\mathbf{w}_{.j}$ . Its distribution is  $\mathbf{w} \sim \text{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \Psi)$ ;
- $\mathbf{f} \in \mathbb{R}^{n \times m}$  denote the matrix containing  $(i, j)$  entries  $f_j(x_i)$ , and  $\mathbf{f}_0$  a vector equal to  $(f_0(x_1), \dots, f_0(x_n))^\top$ . We then have  $\mathbf{f} = \mathbf{1}_n \mathbf{f}_0^\top + \mathbf{H}_\eta \mathbf{w} \sim \text{MN}_{n,m}(\mathbf{1}_n \mathbf{f}_0^\top, \mathbf{H}_\eta^2, \Psi)$ ;
- $\alpha = (\alpha + \alpha_1, \dots, \alpha + \alpha_m)^\top \in \mathbb{R}^m$  be the vector of intercepts;
- $\mu = \mathbf{1}_n \alpha^\top + \mathbf{f}$ , whose  $(i, j)$  entries are  $\mu_j(x_i) = \alpha + \alpha_j + f_j(x_i)$ ; and
- $\mathbf{y}^* \in \mathbb{R}^{n \times m}$  denote the matrix containing  $(i, j)$  entries  $y_{ij}^*$ . That is,  $\mathbf{y}^* = \mu + \epsilon$ , so  $\mathbf{y}^* | \mathbf{w} \sim \text{MN}_{n,m}(\mu = \mathbf{1}_n \alpha^\top + \mathbf{H}_\eta \mathbf{w}, \mathbf{I}_n, \Psi^{-1})$  and  $\text{vec } \mathbf{y}^* \sim N_{nm}(\text{vec}(\mathbf{1}_n \alpha^\top), (\Psi \otimes \mathbf{H}_\eta^2) + (\Psi^{-1} \otimes \mathbf{I}_n))$ . The marginal distribution of  $\mathbf{y}^*$  cannot be written as a matrix normal, except when  $\Psi = \mathbf{I}_m$ .

Before proceeding with estimating the I-probit model (5.10), we lay out several standing assumptions:

**A4 Centred responses.** Set  $\alpha = 0$ .

**A5 Zero prior mean.** Assume a zero prior mean  $f_0(x) = 0$  for all  $x \in \mathcal{X}$ .

**A6 Fixed error precision.** Assume  $\Psi$  is fixed.

Assumption A4 is a requirement for identifiability. Assumption A5 is motivated by a similar argument to assumption A2 in the normal I-prior model. As for assumption A6, we do not consider estimation of the error precision in this thesis mainly due to time limitations. More on this in Section X.

ass:A4

ass:A5

ass:A6

sec:iaa

### 5.3 Identifiability and IIA

The parameters in a linear multinomial probit model is well known to be unidentified (Michael P. Keane, 1992; Train, 2009), and the reason for this is two-fold. Firstly, an addition of a constant to the latent variables  $y_{ij}^*$ 's in (5.7) will not change which latent variable is maximal, and therefore leaves the model unchanged. Secondly, all latent variables can be scaled by some positive constant without changing which latent variable is largest. Therefore, a *linear parameterisation* for the multinomial probit model is not identified as there can be more than one set of parameters for which the class probabilities are the same. To fix this issue, constraints are imposed on location and scale of the latent variables.

However, for the I-probit model, this is not the case, because the model is not related to the parameters  $\theta = \{\alpha_1, \dots, \alpha_m, \eta, \Psi\}$  linearly. One cannot simply add to or multiply  $\theta$  by a constant and expect the model to be left unchanged. Thus, the I-probit model is identified in the parameter set  $\theta$  without having to impose any restrictions, particularly on the precision matrix  $\Psi$  (if this is to be estimated).

To see how the I-probit model is location identified, suppose assumptions A4 and A5 hold, and consider a constant  $a$  added to the latent propensities. This would then imply the relationship

$$a + y_{ij}^* = \overbrace{a + \alpha_j}^{\alpha_j^*} + f_j(x_i) + \epsilon_{ij},$$

which is similar to adding the constant  $a$  to all of the intercept parameters  $\alpha_j$ —denote these new intercepts by  $\alpha_j^*$ . As a requirement of the functional ANOVA decomposition, the  $\alpha_j^*$ 's need to sum to zero, but we already have that  $\sum_{j=1}^m \alpha_j = 0$ , so it must be that  $a = 0$ . This also highlights the reason behind assumption A4 and A5 for fixing the grand intercept  $\alpha$  to zero.

As for identification in scale, consider multiplying the latent variables by  $c > 0$ . Denote by  $\mathbf{V}_y^*(\omega) \in \mathbb{R}^{nm \times nm}$  the marginal covariance matrix of the latent propensities, which depends on the scale parameters  $\omega = \{\eta, \Psi\}$ . The scaled latent variables  $\{c^{1/2}y_{ij}^* \mid \forall i, j = 1, \dots\}$ , which collectively has (marginal) variance and covariances given by the matrix  $c\mathbf{V}_y^*(\omega)$ , is expected to have been generated from the model with param-

eters  $c\omega$ . However, we have that

$$\begin{aligned} c\mathbf{V}_y^*(\omega) &= c(\Psi \otimes \mathbf{H}_\eta^2) + c(\Psi^{-1} \otimes \mathbf{I}_n) \\ &= (c\Psi \otimes \mathbf{H}_\eta^2) + (c\Psi^{-1} \otimes \mathbf{I}_n) \\ &\neq \mathbf{V}_y^*(c\omega). \end{aligned}$$

Now, we turn to a discussion of the role of  $\Psi$  in the model. In decision theory, the independence axiom states that an agent's choice between a set of alternatives should not be affected by the introduction or elimination of a choice option. The probit model is suitable for modelling multinomial data where the independence axiom, which is also known as the *independence of irrelevant alternatives* (IIA) assumption, is not desired. Such cases arise frequently in economics and social science, and the famous Red-Bus-Blue-Bus example is often used to illustrate IIA: suppose commuters face the decision between taking cars and red busses. The addition of blue busses to commuters' choices should, in theory, be more likely chosen by those who prefer taking the bus over cars. That is, assuming commuters are indifferent about the colour of the bus, commuters who are predisposed to taking the red bus would see the blue bus as an identical alternative. Yet, if IIA is imposed, then the three choices are distinct, and the fact that red and blue busses are substitutable is ignored.

To put it simply, the model is IIA if choice probabilities depend only on the choice in consideration, and not on any other alternatives. In the I-probit model, or rather, in probit models in general, choice dependency is controlled by the error precision matrix  $\Psi$ . Specifically, the off-diagonal elements  $\Psi_{jk}$  capture the correlation between alternatives  $j$  and  $k$ . Allowing all  $m(m+1)/2$  covariance elements of  $\Psi$  to be non-zero leads to the *full I-probit model*, and would not assume an IIA position.

While it is an advantage to be able to model the correlations across choices (unlike in logistic models), there are applications where the IIA assumption would not adversely affect the analysis, such as when all the choices are mutually exclusive and exhaustive. Some analyses might also be indifferent as to whether or not choice dependency exists. In these situations, it would be beneficial, algorithmically speaking, to reduce the I-probit model to a simpler version by assuming  $\Psi = \text{diag}(\psi_1, \dots, \psi_m)$ , which would trigger the IIA assumption in the I-probit model. We refer to this model as the *independent I-probit model*.

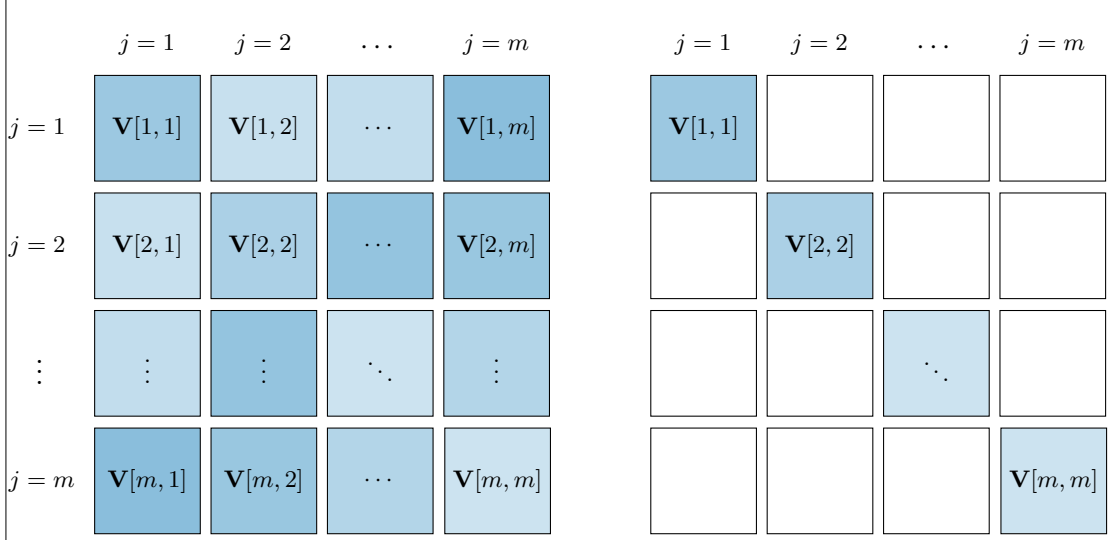


Figure 5.1: Covariance structure of the full I-probit model (left) and the independent I-probit model (right). The full model has  $m^2$  blocks of  $n \times n$  symmetric matrices, and the blocks themselves are arranged symmetrically about the diagonal. The independent model, on the other hand, has a block diagonal structure.

fig:iprobco  
vstr

The independence assumption causes the distribution of the latent variables to be  $y_{ij}^* \sim N(\mu_k(x_i), \sigma_j^2)$  for  $j = 1, \dots, m$ , where  $\sigma_j^2 = \psi_j^{-1}$ . As a continuation of line (5.9), we can show the class probability  $p_{ij}$  to be

$$\begin{aligned}
 p_{ij} &= \int \cdots \int \prod_{\substack{k=1 \\ \{y_{ik}^* > y_{ik}^* | \forall k \neq j\}}}^m \left\{ \phi(y_{ik}^* | \mu_k(x_i), \sigma_k^2) dy_{ik}^* \right\} \\
 &= \int \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( \frac{y_{ij}^* - \mu_k(x_i)}{\sigma_k} \right) \cdot \phi(y_{ij}^* | \mu_j(x_i), \sigma_j^2) dy_{ij}^* \\
 &= E_Z \left[ \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( \frac{\sigma_j}{\sigma_k} Z + \frac{\mu_j(x_i) - \mu_k(x_i)}{\sigma_k} \right) \right] \tag{5.11}
 \end{aligned}$$

{eq:pij2}

where  $Z \sim N(0, 1)$ ,  $\Phi(\cdot)$  its cdf, and  $\phi(\cdot | \mu, \sigma^2)$  is the pdf of  $X \sim N(\mu, \sigma^2)$ . The equation (5.9) is thus simplified to a unidimensional integral involving the Gaussian pdf and cdf, which can be computed fairly efficiently using quadrature methods. The probit link function is evidently seen in the above equation.

## 5.4 Estimation

As with the normal I-prior model, an estimate of the posterior regression function with optimised hyperparameters is sought. The log likelihood function  $L(\cdot)$  for  $\theta$  using all  $n$  observations  $\{(y_1, x_1), \dots, (y_n, x_n)\}$  is obtained by integrating out the I-prior from the categorical likelihood, as follows:

$$\begin{aligned} L(\theta) &= \log \int p(\mathbf{y}|\mathbf{w}, \theta) p(\mathbf{w}|\theta) d\mathbf{w} \\ &= \log \int \prod_{i=1}^n \prod_{j=1}^m \left( g_j^{-1} \left( \alpha_k + \overbrace{f_k(x_i)}^{\sum_{i'=1}^n h_\eta(x_i, x_{i'}) w_{i'j}} \right)_{k=1}^m \right)^{[y_i=j]} \cdot \text{MN}_{n,m}(\mathbf{w}|\mathbf{0}, \mathbf{I}_n, \Psi) d\mathbf{w} \end{aligned} \quad (5.12)$$

{eq:intractablelikelihood}

where we have denoted the probit relationship from (5.9) using the function  $g_j^{-1} : \mathbb{R}^m \rightarrow [0, 1]$ . Unlike in the continuous response models, the integral does not present itself in closed form due to the conditional categorical PMF of the  $y_i$ 's, which they themselves involve integrals of multivariate normal densities. For binary response models,  $g^{-1}$  is simply the probit function, but for multinomial responses, this can be quite challenging to evaluate—more on this in [Section X](#).

Furthermore, the posterior distribution of the regression function, which requires the density of  $\mathbf{w}|\mathbf{y}$ , depends on the marginalisation provided by (5.12). The challenge of estimation is then to first overcome this intractability by means of a suitable approximation of the integral. We present three possible avenues to achieve this aim, namely the Laplace approximation, Markov chain Monte Carlo (MCMC) methods, and variational Bayes.

### 5.4.1 Laplace approximation

To compute the posterior density  $p(\mathbf{w}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{w})p(\mathbf{w}) =: e^{Q(\mathbf{w})}$  with normalising constant equal to the marginal density of  $\mathbf{y}$ ,  $p(\mathbf{y}) = \int e^{Q(\mathbf{w})} d\mathbf{w}$ , we have established that this is intractable. Laplace's method ([Kass and Raftery, 1995, §4.1.1, pp. 777–778](#)) entails expanding a Taylor series for  $Q$  about its posterior mode  $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$ ,

which gives the relationship

$$\begin{aligned}
 Q(\mathbf{w}) &= Q(\hat{\mathbf{w}}) + \underbrace{(\mathbf{w} - \hat{\mathbf{w}})^\top \nabla Q(\hat{\mathbf{w}})}_{\rightarrow 0} - \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{\Omega}(\mathbf{w} - \hat{\mathbf{w}}) + \dots \\
 &\approx Q(\hat{\mathbf{w}}) - \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{\Omega}(\mathbf{w} - \hat{\mathbf{w}}),
 \end{aligned}$$

because, assuming that  $Q$  has a unique maxima,  $\nabla Q$  evaluated at its mode is zero. This is recognised as the logarithm of an unnormalised Gaussian density, implying  $\mathbf{w}|\mathbf{y} \sim N_n(\hat{\mathbf{w}}, \mathbf{\Omega}^{-1})$ . Here,  $\mathbf{\Omega} = -\nabla^2 Q(\mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}}$  is the negative Hessian of  $Q$  evaluated at the posterior mode, and is typically obtained as a byproduct of the maximisation routine of  $Q$  using gradient or quasi-gradient based methods.

The marginal distribution is then approximated by

$$\begin{aligned}
 p(\mathbf{y}) &\approx \int \exp \underbrace{Q(\hat{\mathbf{w}}) - \frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{\Omega}(\mathbf{w} - \hat{\mathbf{w}})}_{\widehat{Q}(\mathbf{w})} d\mathbf{w} \\
 &= (2\pi)^{n/2} |\mathbf{\Omega}|^{-1/2} e^{Q(\hat{\mathbf{w}})} \int (2\pi)^{-n/2} |\mathbf{\Omega}|^{1/2} \exp \left( -\frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{\Omega}(\mathbf{w} - \hat{\mathbf{w}}) \right) d\mathbf{w} \\
 &= (2\pi)^{n/2} |\mathbf{\Omega}|^{-1/2} p(\mathbf{y}|\hat{\mathbf{w}}) p(\hat{\mathbf{w}}).
 \end{aligned}$$

The log marginal density of course depends on the parameters  $\theta$ , which becomes the objective function to maximise in a likelihood maximising approach. Note that, should a fully Bayesian approach be undertaken, i.e. priors prescribed on the model parameters using  $\theta \sim p(\theta)$ , then this approach is viewed as a maximum a posteriori approach.

In any case, each evaluation of the objective function  $L(\theta) = \log p(\mathbf{y}|\theta)$  involves finding the posterior modes  $\hat{\mathbf{w}}$ . This is a slow and difficult undertaking, especially for large sample sizes  $n$ —even assuming computation of the class probabilities  $g^{-1}$  is efficient—because the dimension of this integral is exactly the sample size. Furthermore, obtaining standard errors for the parameters are cumbersome, and it is likely that a computationally burdensome bootstrapping approach is needed. Lastly, as a comment, Laplace’s method only approximates the true marginal likelihood well if the true function is small far away from the mode.



### 5.4.2 Variational inference

We turn to variational inference as a method of estimation. Variational methods are widely discussed in the machine learning literature, but there have been efforts to popularise it in statistics (Blei et al., 2017). In a fully Bayesian setting, one obtains an approximation to the intractable posterior distribution of interest, which is then used for inferential purposes in lieu of the actual posterior distribution.

In addition to the I-probit model, suppose that prior distributions are assigned on the hyperparameters of the model,  $\theta \sim p(\theta)$ . By appending the latent variables  $\{\mathbf{y}^*, \mathbf{w}\}$  to the hyperparameters  $\theta$ , we seek an approximation

$$p(\mathbf{y}^*, \mathbf{w}, \theta | \mathbf{y}) \approx \tilde{q}(\mathbf{y}^*, \mathbf{w}, \theta),$$

where  $\tilde{q}$  satisfies  $\tilde{q} = \arg \min_q \text{KL}(q \| p)$ , subject to certain constraints. The constraint considered by us in this thesis is that  $q$  satisfies a *mean-field* factorisation

$$q(\mathbf{y}^*, \mathbf{w}, \theta) = q(\mathbf{y}^*)q(\mathbf{w})q(\theta).$$

Under this scheme, the posterior for  $\mathbf{y}^*$  is found to be a *conically truncated multivariate normal* distribution, and for  $\mathbf{w}$ , a multivariate normal distribution. The posterior density  $q(\theta)$  is often of a recognisable form, and usually one of the exponential family densities (normal, Wishart or gamma). This is useful, because point estimates of the hyperparameters can be taken to be either the mean or mode of these well-known distributions. In cases where  $q(\theta)$  does not conform to an exponential family type density, then inference can still be done by sampling methods.

It can be shown that, for some variational density  $q$ , the marginal log-likelihood is an upper-bound for the quantity  $\mathcal{L}$

$$\log p(\mathbf{y}) \geq \mathbb{E}_q \log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta) - \mathbb{E}_q \log \tilde{q}(\mathbf{y}^*, \mathbf{w}, \theta) =: \mathcal{L},$$

a quantity often referred to as the *evidence lower bound* (ELBO). It turns out that minimising  $\text{KL}(q \| p)$  is equivalent to maximising the ELBO, a quantity that is more practical to work with than the KL divergence. That is, if  $\tilde{q}$  approximates the true posterior well, then the ELBO is a suitable proxy for the maximised marginal log-likelihood.

The algorithm to obtain  $\tilde{q}$  which maximises the ELBO is known as the *coordinate ascent variational inference* (CAVI) algorithm. The algorithm itself typically condenses to that of a simple, sequential updating scheme, akin to the expectation-maximisation (EM) algorithm for exponential families we saw in Chapter 4, which is very fast to implement compared to the other methods described in the previous subsections. A full derivation of the variational algorithm used by us will be described in [Section 5.5](#).

### 5.4.3 Markov chain Monte Carlo methods

As an alternative to the deterministic Bayesian approach of variational inference, it is possible to use Markov chain Monte Carlo sampling methods as an approach to stochastically approximate the intractable posterior distribution.

[Albert and Chib \(1993\)](#) showed that the latent variable approach to probit models can be analysed using exact Bayesian methods, due to the underlying normality structure. Paired with corresponding conjugate prior choices, sampling from the posterior is very simple using a Gibbs sampling approach. On the other hand, this data augmentation scheme enlarges the variable space to  $n + q$  dimensions, where  $q$  is the number of parameters to estimate, which is inefficient and computationally challenging especially when  $n$  is large. It is no longer possible to marginalise the normal latent variables from the model, as this is intractable, as discussed previously.

Hamiltonian Monte Carlo is another possibility, since it does not require conjugacy. For binary models, this is a feasible approach because the class probabilities are a function of the normal CDF, which means that it is doable using off-the-shelf software such as Stan. However, with multinomial responses, the arduous task of computing class probabilities, which involve integration of an at most  $m$ -dimensional normal density, must be addressed separately.

### 5.4.4 Comparison of estimation methods

**Compare: Laplace, variational and HMC.**

The three estimation methods described aim to overcome the intractable integral by means of either a deterministic approximation (Laplace and variational inference) or a stochastic approximation (MCMC). In the Laplace and variational method, the posterior distribution of  $\mathbf{w}$  ends up being approximated by a Gaussian distribution, although the

mean and variance is different in each method. In essence, once  $\mathbf{w}|\mathbf{y}$  is approximately normal, then estimation of the parameters  $\theta$  using a direct optimisation approach or an EM-type approach is straightforward. On the other hand, MCMC approximates the density  $p(\mathbf{w}|\mathbf{y})$  using samples generated via Gibbs sampling or HMC, and these samples would asymptotically be representative of draws from the true posterior.

Consider the data set... Plot the data. Explain priors for HMC and variational. Compare.

## 5.5 A variational algorithm

We present a variational inference algorithm to estimate the I-probit latent variables  $\mathbf{y}^*$  and  $\mathbf{w}$ , together with the parameters  $\theta = \{\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)^\top, \eta\}$  with *fixed error precision*  $\Psi^2$ . Begin by choosing prior distributions on the parameters,  $p(\theta) = p(\boldsymbol{\alpha})p(\eta)$ . The following flat, uninformative priors are suggested:

- **Kernel parameters**  $\eta$ . This may include parameters such as the Hurst index, lengthscale and offset parameters, in addition to the RKHS scale parameters  $\lambda_1, \dots, \lambda_p$ , and each with their own support. For the scale parameters, assign each  $\lambda_k$  the vague prior

$$\lambda_k \stackrel{\text{iid}}{\sim} \text{N}(0, v_\lambda = 0.001^{-1}), \quad k = 1, \dots, p.$$

As  $v_k^{-1} \rightarrow 0$ , the prior becomes  $p(\lambda_k) \propto \text{const.}$ , an improper prior. The default choice for the rest of the kernel parameters is an improper prior  $p(\eta) \propto \text{const.}$

- **Intercepts**  $\alpha_1, \dots, \alpha_m$ . Assign independent, vague normal priors for each intercept

$$\alpha_j \stackrel{\text{iid}}{\sim} \text{N}(0, v_\alpha = 0.001^{-1}).$$

Although one may devote more attention to the prior specification of these parameters, for our purposes it suffices that they are independent component-wise, and that they are conjugate priors for the complete conditional density  $p(\theta|\mathbf{y}, \mathbf{y}^*, \mathbf{w})$ .

---

<sup>2</sup>It turns out that the variational algorithm as presented is not suited to estimate  $\Psi$ . This issue is discussed further in Section X.

sec:iprobit  
var

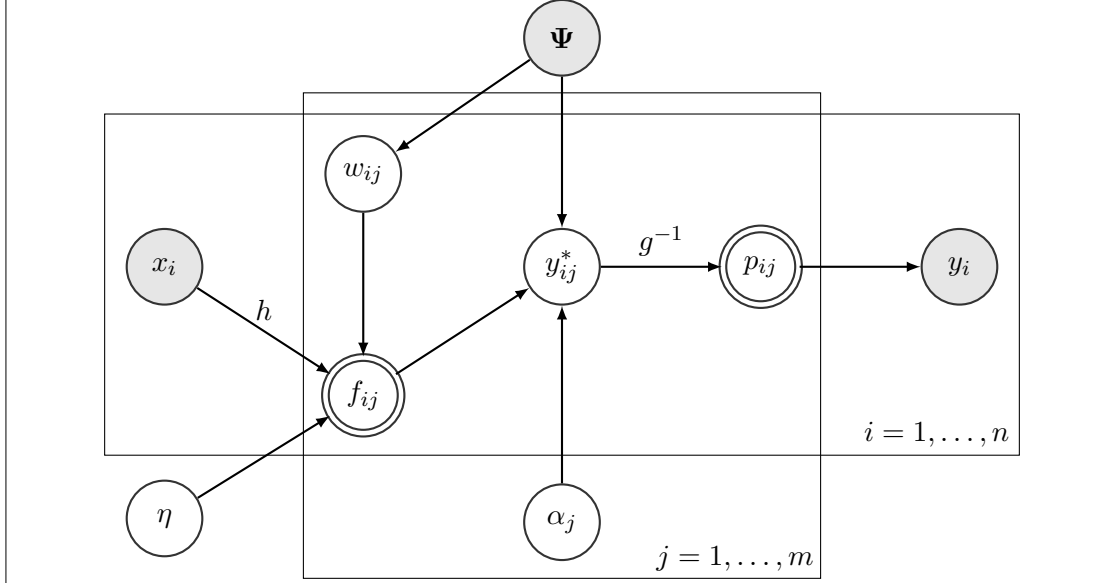


Figure 5.2: A DAG of the I-probit model. Observed/fixed nodes are shaded, while double-lined nodes represents calculable quantities.

The posterior density of  $\mathcal{Z} = \{\mathbf{y}^*, \mathbf{w}, \theta\}$  is approximated by a mean-field variational density  $q$ , i.e.

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

Additionally, we assume independence among the components of  $\theta$  so that  $q(\theta) = \prod_k q(\theta_k)$ . We now present the mean-field variational distributions for each of unknowns in  $\mathcal{Z}$ . On notation: we will typically refer to posterior means of the parameters  $\mathbf{y}^*$ ,  $\mathbf{w}$ ,  $\theta$  and so on by the use of a tilde. For instance, we write  $\tilde{\mathbf{w}}$  to mean  $\mathbb{E}_{\mathbf{w} \sim q}[\mathbf{w}]$ , the expected value of  $\mathbf{w}$  under the pdf  $q(\mathbf{w})$ . The distributions are simply stated, but a full derivation is given in the appendix.

### 5.5.1 Latent propensities $\mathbf{y}^*$

The fact that the rows  $\mathbf{y}_{i\cdot}^* \in \mathbb{R}^m$ ,  $i = 1, \dots, n$  of  $\mathbf{y}^* \in \mathbb{R}^{n \times m}$  are independent can be exploited, which yields an induced factorisation  $q(\mathbf{y}^*) = \prod_{i=1}^n q(\mathbf{y}_{i\cdot}^*)$ . Define the set  $\mathcal{C}_j = \{y_{ij}^* > y_{ik}^* \mid \forall k \neq j\}$ . Then  $q(\mathbf{y}_{i\cdot}^*)$  is the density of a multivariate normal distribution with mean  $\tilde{\boldsymbol{\mu}}_{i\cdot} = \tilde{\boldsymbol{\alpha}} + \tilde{\mathbf{w}}^\top \tilde{\mathbf{h}}_\eta(x_i)$ , and variance  $\Psi^{-1}$  subject to the truncation of its components to the set  $\mathcal{C}_{y_i}$ . That is, for each  $i = 1, \dots, n$  and noting the observed value

$y_i \in \{1, \dots, m\}$ , the  $\mathbf{y}_i^*$ 's are distributed according to

$$\mathbf{y}_i^* \stackrel{\text{iid}}{\sim} \begin{cases} N_m(\tilde{\boldsymbol{\mu}}_{i\cdot}, \boldsymbol{\Psi}^{-1}) & \text{if } y_{iy_i}^* > y_{ik}^*, \forall k \neq y_i \\ 0 & \text{otherwise.} \end{cases} \quad (5.13)$$

{eq:ystardist}

We denote this by  $\mathbf{y}_i^* \stackrel{\text{iid}}{\sim} \text{tN}(\tilde{\boldsymbol{\mu}}_{i\cdot}, \boldsymbol{\Psi}^{-1}, \mathcal{C}_{y_i})$ , and the important properties of this distribution are explored in the appendix.

The required expectations  $E\mathbf{y}_i^* = E(y_{i1}^*, \dots, y_{im}^*)^\top$  are tricky to compute. One strategy might be Monte Carlo integration: using samples from  $N_m(\tilde{\boldsymbol{\mu}}_{i\cdot}, \boldsymbol{\Psi}^{-1})$ , disregard those that do not satisfy the condition  $y_{iy_i}^* > y_{ik}^*, \forall k \neq j$ , and then take the sample average. This works reasonably well so long as the truncation region does not fall into the extreme tails of the multivariate normal. Alternatively, a fast, Gibbs based approach to estimating the mean or any other quantity  $E r(\mathbf{y}_i^*)$  can be implemented, and this is detailed in the appendix.

If the independent I-probit model is considered, where the covariance matrix has the independent structure  $\boldsymbol{\Psi} = \text{diag}(\sigma_1^{-2}, \dots, \sigma_m^{-2})$ , then the expected value can be considered component-wise, and each component of this expectation is given by

$$\tilde{y}_{ik}^* = \begin{cases} \tilde{\mu}_{ik} - \sigma_k C_i^{-1} \int \phi_{ik}(z) \prod_{l \neq k, j} \Phi_{il}(z) \phi(z) dz & \text{if } k \neq y_i \\ \tilde{\mu}_{iy_i} - \sigma_{y_i} \sum_{k \neq y_i} (\tilde{y}_{ik}^* - \tilde{\mu}_{ik}) & \text{if } k = y_i \end{cases} \quad (5.14)$$

{eq:ystarupdate}

with

$$\begin{aligned} \phi_{ik}(Z) &= \phi\left(\frac{\sigma_{y_i}}{\sigma_k} Z + \frac{\tilde{\mu}_{iy_i} - \tilde{\mu}_{ik}}{\sigma_k}\right) \\ \Phi_{ik}(Z) &= \Phi\left(\frac{\sigma_{y_i}}{\sigma_k} Z + \frac{\tilde{\mu}_{iy_i} - \tilde{\mu}_{ik}}{\sigma_k}\right) \\ C_i &= \int \prod_{l \neq j} \Phi_{il}(z) \phi(z) dz \end{aligned}$$

and  $Z \sim N(0, 1)$  with pdf and cdf  $\phi(\cdot)$  and  $\Phi(\cdot)$  respectively. The integrals that appear above are functions of a unidimensional Gaussian pdf, and these can be computed rather efficiently using quadrature methods.

### 5.5.2 I-prior random effects $\mathbf{w}$

Given that both  $\text{vec } \mathbf{y}^* | \text{vec } \mathbf{w}$  and  $\text{vec } \mathbf{w}$  are normally distributed, we find that the conditional posterior distribution  $p(\mathbf{w} | \mathcal{Z}_{-\mathbf{w}}, \mathbf{y})$  is also normal, and therefore the approximate posterior density  $q$  for  $\text{vec } \mathbf{w} \in \mathbb{R}^{nm}$  is also normal with mean and precision given by

$$\text{vec } \tilde{\mathbf{w}} = \tilde{\mathbf{V}}_w (\boldsymbol{\Psi} \otimes \tilde{\mathbf{H}}_\eta) \text{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top) \quad \text{and} \quad \tilde{\mathbf{V}}_w^{-1} = (\boldsymbol{\Psi} \otimes \tilde{\mathbf{H}}_\eta^2) + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n). \quad (5.15)$$

{eq:varipostw}

We note the similarity between (5.15) above and the posterior distribution for the I-prior random effects in a normal model (4.7) seen in the previous chapter. Naïvely computing the inverse  $\tilde{\mathbf{V}}_w^{-1}$  presents a computational challenge, as this takes  $O(n^3 m^3)$  time. By exploiting the Kronecker product structure in  $\tilde{\mathbf{V}}_w$ , we are able to efficiently compute the required inverse in roughly  $O(n^3 m)$  time—see the appendix for details.

If the independent I-probit model is assumed, i.e.  $\tilde{\boldsymbol{\Psi}} = \text{diag}(\tilde{\psi}_1, \dots, \tilde{\psi}_m)$ , then the posterior covariance matrix  $\tilde{\mathbf{V}}_w$  has a simpler structure: random matrix  $\mathbf{w}$  will have columns which are independent of each other. By writing  $\mathbf{w}_{\cdot j} = (w_{1j}, \dots, w_{nj})^\top \in \mathbb{R}^n$ ,  $j = 1, \dots, m$ , to denote the column vectors of  $\mathbf{w}$  and with a slight abuse of notation, we have that

$$N_{nm}(\text{vec } \mathbf{w} | \text{vec } \tilde{\mathbf{w}}, \tilde{\mathbf{V}}_w) = \prod_{j=1}^m N_n(\mathbf{w}_{\cdot j} | \tilde{\mathbf{w}}_{\cdot j}, \tilde{\mathbf{V}}_{w_j}),$$

where

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

The consequence of this is that the posterior regression functions are class independent, the exact intended effect by specifying a diagonal precision matrix  $\boldsymbol{\Psi}$ .

### 5.5.3 Kernel parameters $\eta$

sec:varupdatea

The posterior density  $q$  involving the kernel parameters is of the form

$$\log q(\eta) = -\frac{1}{2} \text{tr } E_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ (\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top - \mathbf{H}_\eta \mathbf{w}) \boldsymbol{\Psi} (\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top - \mathbf{H}_\eta \mathbf{w})^\top \right] + \log p(\eta) + \text{const.}$$

where  $p(\eta)$  is an appropriate prior density for  $\eta$ . Generally, samples  $\eta^{(1)}, \dots, \eta^{(T)}$  from  $\tilde{q}(\eta)$  may be obtained using a Metropolis algorithm, so that quantities such as  $\tilde{\mathbf{H}}_\eta = \mathbf{E}_{\eta \sim q} \mathbf{H}_\eta$  and the like may be approximated using  $\frac{1}{T} \sum_{t=1}^T \mathbf{H}_{\eta^{(t)}}$ . Details of the Metropolis sampler is available in the appendix.

When only RKHS scale parameters are involved, then the distribution  $q$  can be found in closed-form, much like in the exponential family EM algorithm described in [Section 4.3.3](#). Under the same setting as in that subsection, assume that only  $\eta = \{\lambda_1, \dots, \lambda_p\}$  need be estimated, and for each  $k = 1, \dots, p$ , we can decompose the kernel matrix as  $\mathbf{H}_\eta = \lambda_k \mathbf{R}_k + \mathbf{S}_k$  and its square as  $\mathbf{H}_\eta^2 = \lambda_k^2 \mathbf{R}_k^2 + \lambda_k \mathbf{U}_k + \mathbf{S}_k^2$ . Additionally, we impose a further mean-field restriction on  $q(\eta)$ , i.e.,  $q(\eta) = \prod_{k=1}^p p(\lambda_k)$ . Then, by using independent and identical normal priors on the  $\lambda_k$ 's, such as the one listed at the beginning of this section, we find that  $q(\lambda_k)$  is the density of a normal distribution with mean  $d_k c_k^{-1}$  and variance  $c_k^{-1}$ , where

$$c_k = \text{tr}(\Psi \mathbf{E}[\mathbf{w}^\top \mathbf{R}_k^2 \mathbf{w}]) + v_\lambda^{-2}$$

and

$$d_k = \text{tr}\left(\Psi(\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top)^\top \mathbf{R}_k \tilde{\mathbf{w}} - \frac{1}{2} \Psi \mathbf{E}[\mathbf{w}^\top \mathbf{U}_k \mathbf{w}]\right).$$

For a method of evaluating quantities such as  $\text{tr}(\mathbf{C} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}])$  for suitably sized matrices  $\mathbf{C}$  and  $\mathbf{D}$ , refer to the appendix.

#### 5.5.4 Intercepts $\boldsymbol{\alpha}$

Finally, the posterior distribution for the intercepts follow a normal distribution with the normal priors specified earlier. The posterior mean and variance for the intercepts are given by  $\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{A}}^{-1} \tilde{\mathbf{a}}$  and  $\tilde{\mathbf{A}}^{-1}$  respectively, where

$$\tilde{\mathbf{a}} = \sum_{i=1}^n \Psi(\tilde{\mathbf{y}}_{i\cdot}^* - \tilde{\mathbf{w}}^\top \tilde{\mathbf{h}}_\eta(x_i)) \quad \text{and} \quad \tilde{\mathbf{A}} = n\Psi + v_\alpha \mathbf{I}_m.$$

If  $\Psi$  is diagonal, the components of  $\boldsymbol{\alpha}$  would be independent, and each would be distributed according to

$$\mathcal{N}\left(\frac{\psi_j \sum_{i=1}^n (\tilde{y}_{ij}^* - \tilde{f}_{ij})^2}{n\psi_j + v_\alpha^{-1}}, \frac{1}{n\psi_j + v_\alpha^{-1}}\right).$$

Here, we used the notation  $\tilde{f}_{ij}$  to mean the  $(i, j)$ 'th element of  $\mathbb{E}[\mathbf{H}_\eta \mathbf{w}] = \tilde{\mathbf{H}}_\eta \tilde{\mathbf{w}} \in \mathbb{R}^{n \times m}$ . Note that it is necessary, as discussed earlier, that  $\sum_{j=1}^m \alpha_j = 0$  for identifiability.

### 5.5.5 The CAVI algorithm

One will have noticed that the evaluation of each component of the posterior depends on knowing the posterior distribution of the rest of the components. This circular dependence is dealt with by way of an iterative updating scheme of the components. Using an arbitrary starting value, each component is updated in turn according to the above derivations, until a maximum number of iterations is reached, or ideally, until a convergence criterion is met. In variational inference, the ELBO is used to assess convergence. The expression for the ELBO for the I-probit model is derived in the appendix. The CAVI algorithm for the I-probit model is summarised in [Algorithm 1](#).

Similar to the EM algorithm, each iteration of the algorithm increases the ELBO to a stationary point ([Blei et al., 2017](#))—hence the name coordinate ascent variational inference (CAVI). Unlike the EM algorithm though, the CAVI algorithm does *not* guarantee an increase in the marginal log-likelihood at each step, nor does it guarantee convergence to the global maxima of the log-likelihood.

Further, the ELBO expression to be maximised is often not convex, which means the CAVI algorithm may terminate at local modes, for which there may be many. Note that the variational distribution with the higher ELBO value is the distribution that is closer, in terms of the KL divergence, to the true posterior distribution. In our experience, multiple random starts alleviates this issue for the I-probit model.

## 5.6 Post-estimation

Working within a variational Bayesian framework means that we are able to perform inferences on any quantity of interest using the (approximate) posterior distributions obtained. Any of the post estimation procedures explained in the previous chapter when dealing with normal I-prior models can be extended here.

Prediction of a new data point  $x_{\text{new}}$  is described. Step one is to determine the distribution of the posterior regression functions in each class,  $\mathbf{f}(x_{\text{new}}) = \mathbf{w}^\top \mathbf{h}_\eta(x_{\text{new}})$ , given values for the parameters  $\theta$  of the I-probit model. To this end, we use the posterior mean estimate for  $\theta$ , and denote them with tildes, as we have done so far in this chapter.



alg:caviipr  
obit

#### Algorithm 1 CAVI for the I-probit model

```

1: procedure INITIALISATION
2:   Initialise  $\tilde{\mathbf{y}}^{*(0)}, \tilde{\mathbf{w}}^{(0)}, \tilde{\boldsymbol{\alpha}}^{(0)}, \tilde{\mathbf{H}}_{\eta^{(0)}}, \Psi$ 
3:    $t \leftarrow 0$ 
4: end procedure

5: while not converged do
6:   for  $i = 1, \dots, n$  do ▷ Update  $\mathbf{y}^*$ 
7:      $q^{(t+1)}(\mathbf{y}_{i.}^*) \leftarrow {}^t\text{N}_m(\tilde{\boldsymbol{\alpha}}^{(t)} + \tilde{\mathbf{w}}^{(t)\top} \tilde{\mathbf{h}}_{\eta^{(t)}}(x_i), \Psi, \mathcal{C}_{y_i})$ 
8:      $\tilde{\mathbf{y}}_{i.}^{*(t+1)} \leftarrow \text{E}_{q^{(t+1)}}[\mathbf{y}_{i.}^*]$ 
9:   end for

10:   $\mathbf{V}_w^{(t+1)} \leftarrow ((\Psi \otimes \tilde{\mathbf{H}}_{\eta^{(t)}}^2) + (\Psi^{-1} \otimes \mathbf{I}_n))^{-1}$  ▷ Update  $\mathbf{w}$ 
11:   $\tilde{\mathbf{w}}^{(t+1)} \leftarrow \tilde{\mathbf{V}}_w^{(t+1)}(\Psi \otimes \tilde{\mathbf{H}}_{\eta^{(t)}}) \text{vec}(\tilde{\mathbf{y}}^{*(t+1)} - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^{(t)\top})$ 
12:   $q^{(t+1)}(\mathbf{w}) \leftarrow \text{N}_{nm}(\tilde{\mathbf{w}}^{(t+1)}, \mathbf{V}_w^{(t+1)})$ 

13:  Update  $q^{(t+1)}(\eta)$  as per Section 5.5.3 ▷ Update  $\eta$ 
14:  Sample  $\eta^{[1]}, \dots, \eta^{[T]} \sim q^{(t+1)}(\eta)$ 
15:   $\tilde{\mathbf{H}}_{\eta^{(t+1)}} \leftarrow \frac{1}{T} \sum_{i=1}^T \mathbf{H}_{\eta^{[i]}}$  and  $\tilde{\mathbf{H}}_{\eta^{(t+1)}}^2 \leftarrow \frac{1}{T} \sum_{i=1}^T \mathbf{H}_{\eta^{[i]}}^2$ 

16:   $\tilde{\boldsymbol{\alpha}}^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n (\tilde{\mathbf{y}}_{i.}^{*(t+1)} - \tilde{\mathbf{w}}^{(t+1)\top} \tilde{\mathbf{h}}_{\eta^{(t+1)}}(x_i))$  ▷ Update  $\boldsymbol{\alpha}$ 
17:   $q^{(t+1)}(\boldsymbol{\alpha}) \leftarrow \text{N}_m(\tilde{\boldsymbol{\alpha}}^{(t+1)}, \frac{1}{n} \Psi^{-1})$ 

18:  Calculate ELBO  $\mathcal{L}^{(t+1)}$ 
19:   $t \leftarrow t + 1$ 
20: end while

```

As we know,  $\text{vec } \mathbf{w}$  is normally distributed with mean and variance according to (5.15). By writing  $\text{vec } \tilde{\mathbf{w}} = (\tilde{\mathbf{w}}_{\cdot 1}, \dots, \tilde{\mathbf{w}}_{\cdot m})^\top$  to separate out the I-prior random effects per class, we have that  $\mathbf{w}_{\cdot j} | \tilde{\theta} \sim \text{N}_n(\tilde{\mathbf{w}}_{\cdot 1}, \tilde{\mathbf{V}}_w[j, j])$ , and  $\text{Cov}(\mathbf{w}_{\cdot j}, \mathbf{w}_{\cdot k}) = \tilde{\mathbf{V}}_w[j, k]$ , where the  $[\cdot, \cdot]$  indexes the  $n \times n$  sub-block of the block matrix structured matrix  $\mathbf{V}_w$ . Thus, for each class  $j = 1, \dots, m$  and any  $x \in \mathcal{X}$ ,

$$f_j(x) | \mathbf{y}, \tilde{\theta} \sim \text{N}(\tilde{\mathbf{h}}_{\eta}(x)^\top \mathbf{w}_{\cdot j}, \tilde{\mathbf{h}}_{\eta}(x)^\top \tilde{\mathbf{V}}_w[j, j] \tilde{\mathbf{h}}_{\eta}(x)),$$

and the covariance between the regression functions in two different classes is

$$\text{Cov}[f_j(x), f_k(x) | \mathbf{y}, \tilde{\theta}] = \tilde{\mathbf{h}}_{\eta}(x)^\top \tilde{\mathbf{V}}_w[j, k] \tilde{\mathbf{h}}_{\eta}(x).$$

Then, in step two, using the results obtained in the previous chapter in [Section 4.4](#),

we have that the latent propensities  $y_{\text{new},j}^*$  for each class are normally distributed with mean, variance, and covariances

$$\begin{aligned} \mathbb{E}[y_{\text{new},j}^* | \mathbf{y}, \tilde{\theta}] &= \tilde{\alpha}_j + \mathbb{E}[f_j(x_{\text{new}}) | \mathbf{y}, \tilde{\theta}] && =: \hat{\mu}_j(x_{\text{new}}) \\ \text{Var}[y_{\text{new},j}^* | \mathbf{y}, \tilde{\theta}] &= \text{Var}[f_j(x_{\text{new}}) | \mathbf{y}, \tilde{\theta}] + \boldsymbol{\Psi}_{jj}^{-1} && =: \hat{\sigma}_j^2(x_{\text{new}}) \\ \text{Cov}[y_{\text{new},j}^*, y_{\text{new},k}^* | \mathbf{y}, \tilde{\theta}] &= \text{Cov}[f_j(x), f_k(x) | \mathbf{y}, \tilde{\theta}] + \boldsymbol{\Psi}_{jk}^{-1} && =: \hat{\sigma}_{jk}(x_{\text{new}}). \end{aligned}$$

From here, step three would be to extract class information of data point  $x_{\text{new}}$ , which are contained in the normal distribution  $N_m(\hat{\boldsymbol{\mu}}_{\text{new}}, \hat{\mathbf{V}}_{\text{new}})$ , where

$$\hat{\boldsymbol{\mu}}_{\text{new}} = (\mu_1(x_{\text{new}}), \dots, \mu_m(x_{\text{new}}))^{\top} \quad \text{and} \quad \hat{\mathbf{V}}_{\text{new},jk} = \begin{cases} \hat{\sigma}_j^2(x_{\text{new}}) & \text{if } i = j \\ \hat{\sigma}_{jk}(x_{\text{new}}) & \text{if } i \neq j. \end{cases}$$

The predicted class is inferred from the latent variables via

$$\hat{y}_{\text{new}} = \arg \max_k \hat{\mu}_k(x_{\text{new}}),$$

while the probabilities for each class are obtained via integration of a multivariate normal density, as per (5.9), and restated here for convenience:

$$\hat{p}_{\text{new},j} = \int \cdots \int_{\{y_{ij}^* > y_{ik}^* \mid \forall k \neq j\}} \phi(y_{i1}^*, \dots, y_{im}^* | \hat{\boldsymbol{\mu}}_{\text{new}}, \hat{\mathbf{V}}_{\text{new}}) dy_{i1}^* \cdots dy_{im}^*.$$

For the independent I-probit model, class probabilities are obtained in a more compact manner via

$$\hat{p}_{\text{new},j} = \mathbb{E}_Z \left[ \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( \frac{\hat{\sigma}_j(x_{\text{new}})}{\hat{\sigma}_k(x_{\text{new}})} Z + \frac{\hat{\mu}_j(x_{\text{new}}) - \hat{\mu}_k(x_{\text{new}})}{\hat{\sigma}_k^2(x_{\text{new}})} \right) \right],$$

as per (5.11), since the  $m$  components of  $\mathbf{f}(x_{\text{new}})$ , and hence the  $\mathbf{y}_{\text{new},j}^*$ 's, are independent of each other ( $\boldsymbol{\Psi}$  and  $\hat{\mathbf{V}}_{\text{new}}$  are diagonal).

In this Bayesian setting, the analogue of standard errors for the parameters are their posterior standard deviations, which explain the uncertainty surrounding parameters. For the most part, these are easy to come by, and their posterior densities are easy to sample from. This allows us to conduct inference on transformed parameters, such as log odds ratios, quite easily. The procedure would be like this: first obtain samples

of  $\theta^{(1)}, \dots, \theta^{(T)}$  from their respective distributions, then sample  $\mathbf{w}^{(i)} \sim p(\mathbf{w}|\theta^{(i)})$  for  $i = 1, \dots, T$ , and finally obtain samples of class probabilities  $\hat{p}_{xj}^{(1)}, \dots, \hat{p}_{xj}^{(T)}$ ,  $j = 1, \dots, m$ , for a given data point  $x \in \mathcal{X}$ . To obtain a statistic of interest, say, a 95% credibility interval of a function  $r(p_{xj})$  of the probabilities, simply take the empirical lower 2.5th and upper 97.5th percentile of this transformed sample. In this manner, all aspects of uncertainty, from the parameters to the latent variables of the generative model, are accounted for.

It is possible to perform model comparison by comparing the maximised ELBO quantity of several candidate models (Beal and Ghahramani, 2003), and the justification for this is that it supposedly gives a tight lower bound to the marginal likelihood (model evidence), especially if the variational density is close in the KL divergence sense to the true posterior density. This would allow model selection using Bayes factor as a model selection criterion. Kass and Raftery (1995) suggest the following interpretation of observed Bayes factor values for comparing model  $M_1$  against model  $M_0$ .

Table 5.1: Guidelines for interpreting Bayes factors.

$2 \log \text{BF}(M_1, M_0)$	$\text{BF}(M_1, M_0)$	Evidence against $M_0$
0–2	1–3	Not worth more than a bare mention
2–6	3–20	Positive
6–10	20–150	Strong
>10	>150	Very strong

It should be noted that while this works in practice, there is no theoretical basis for model comparison using the ELBO (Blei et al., 2017).

## 5.7 Computational consideration

Computational challenges for the I-probit model stems from two sources: 1) calculation of the class probabilities (5.9); and 2) storage and time requirements for the CAVI. We also discuss issues faced with the estimation of the error precision  $\Psi$ , and suggest ways to overcome this for future work.

tab:bf

misc:mnint

### 5.7.1 Efficient computation of class probabilities

As an opening remark, note that the dimension of the integral (5.9) is  $m - 1$ , since the  $j$ 'th coordinates is fixed relative to the others. An alternative specification of the I-probit model can be made in terms of *relative differences* of the latent propensities. Choosing the first category as the reference category, define new random variables  $z_{ij} = y_{ij}^* - y_{i1}^*$ , for  $j = 2, \dots, m$ . The model (5.7) is equivalently represented by

$$y_i = \begin{cases} 1 & \text{if } \max(z_{i2}, \dots, z_{im}) < 0 \\ j & \text{if } \max(z_{i2}, \dots, z_{im}) = z_{ij} \geq 0. \end{cases} \quad (5.16)$$

Write  $\mathbf{z}_i = (z_{i2}, \dots, z_{im})^\top \in \mathbb{R}^{m-1}$ . Then  $\mathbf{z}_i = \mathbf{Q}\mathbf{y}_i^*$ , where  $\mathbf{Q} \in \mathbb{R}^{(m-1) \times m}$  is the  $(m-1)$  identity matrix pre-augmented with a column vector of minus ones. We have that  $\mathbf{z}_i \stackrel{\text{iid}}{\sim} N_{m-1}(\boldsymbol{\nu}, \boldsymbol{\Omega})$ , where  $\boldsymbol{\nu}_i = \mathbf{Q}\boldsymbol{\mu}_i$  and  $\boldsymbol{\Omega} = \mathbf{Q}\boldsymbol{\Psi}^{-1}\mathbf{Q}^\top$ . Note that if  $\boldsymbol{\Psi}$  is diagonal, then the transformation to  $\boldsymbol{\Omega}$  will not retain diagonality—indeed, each component will undoubtedly be correlated with one another as they are all anchored on the same latent variable.

Now, the class probabilities for  $j = 2, \dots, m$  are

$$p_{ij} = \int_{\{z_{ik} < 0 \mid \forall k \neq 1, j\}} \mathbb{1}[z_{ij} \geq 0] \phi(\mathbf{z}_i | \boldsymbol{\nu}_i, \boldsymbol{\Omega}) d\mathbf{z}_i. \quad (5.17)$$

{eq:pij3}

The class probability  $p_{i1}$  is simply  $p_{i1} = 1 - \sum_{k \neq 1} p_{ik}$ . From this representation of the model, with  $m = 2$  (binary outcomes) we see that

$$p_{i1} = \Phi\left(\frac{z_{i2} - \nu}{\omega^{1/2}}\right) \quad \text{and} \quad p_{i2} = 1 - \Phi\left(\frac{z_{i2} - \nu}{\omega^{1/2}}\right),$$

where  $\Phi(\cdot)$  is the CDF of the standard normal univariate distribution, and  $\nu$  and  $\omega$  are the mean and variance of the univariate random variable  $\mathbf{z}_i = z_{i2}$ . The probit link function involving the cdf of a standard normal is clearly seen here, especially if the error precision is treated as fixed such that  $\omega = 1$ .

The issue at hand here is that for  $m > 4$ , the evaluation of the class probabilities in (5.9) is computationally burdensome using classical methods such as quadrature methods Geweke et al. (1994).

The simplest strategy to overcome this is a frequency simulator (otherwise known as Monte Carlo integration): obtain random samples from  $N_{m-1}(\boldsymbol{\nu}_i, \boldsymbol{\Omega})$ , and calculate how many of these samples fall within the required region. This method is fast and yields unbiased estimates of the class probabilities. However, accuracy of this method is questionable when the mean  $\boldsymbol{\nu}_i$  of the multivariate normal is many standard deviations away from zero (the cutoff region as per (5.17)).

A more reliable method is the probability simulator of Geweke-Hajivassiliou-Keane (GHK) (Geweke, 1991; Hajivassiliou et al., 1996; Michael P Keane, 1994), which we describe now. For clarity, we drop the subscript  $i$  denoting individuals, and write  $\mathbf{z} = (z_1, \dots, z_m)$ , remembering that  $z_1 = 0$ . Suppose that an observation  $y = j$  has been made. Rewrite the model by anchoring on the  $j$ 'th latent variable  $z_j$  as follows:

$$\tilde{\mathbf{z}} := (\overbrace{z_1 - z_j}^{\tilde{z}_1}, \dots, \overbrace{z_{j-1} - z_j}^{\tilde{z}_{j-1}}, \overbrace{z_{j+1} - z_j}^{\tilde{z}_{j+1}}, \dots, \overbrace{z_m - z_j}^{\tilde{z}_m})^\top \in \mathbb{R}^{m-1}.$$

Let  $\boldsymbol{\nu}_{(j)}$  and  $\boldsymbol{\Omega}_{(j)}$  be the appropriately transformed mean vector and covariance matrix for  $\tilde{\mathbf{z}}$ . These are indexed by ' $(j)$ ' because the transformation is dependent on which latent variable the  $\mathbf{z}$ 's are anchored on. Since this transformation is linear,  $\tilde{\mathbf{z}} \sim N_{m-1}(\boldsymbol{\nu}_{(j)}, \boldsymbol{\Omega}_{(j)})$ . For the symmetric and positive definite matrix  $\boldsymbol{\Psi}^{-1}$ , obtain its Cholesky decomposition as  $\boldsymbol{\Omega}_{(j)} = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L}$  is a lower triangular matrix. Then,  $\tilde{\mathbf{z}} = \boldsymbol{\nu}_{(j)} + \mathbf{L}\boldsymbol{\zeta}$ , where  $\boldsymbol{\zeta} \sim N_{m-1}(\mathbf{0}, \mathbf{I}_{m-1})$ . That is,

$$\begin{pmatrix} \tilde{z}_1 \\ \tilde{z}_2 \\ \vdots \\ \tilde{z}_m \end{pmatrix} = \begin{pmatrix} \nu_{(j)1} \\ \nu_{(j)2} \\ \vdots \\ \nu_{(j)m} \end{pmatrix} + \begin{pmatrix} L_{11} & & & \\ L_{21} & L_{22} & & \\ \vdots & \vdots & \ddots & \\ L_{m1} & L_{m2} & \cdots & L_{mm} \end{pmatrix} \begin{pmatrix} \zeta_1 \\ \zeta_2 \\ \vdots \\ \zeta_m \end{pmatrix} = \begin{pmatrix} \nu_{(j)1} + L_{11}\zeta_1 \\ \nu_{(j)2} + \sum_{k=1}^2 L_{k2}\zeta_k \\ \vdots \\ \nu_{(j)m} + \sum_{k=1}^m L_{km}\zeta_k \end{pmatrix}.$$

With this setup, we can calculate  $p_j$ , the probability of class  $j$ , which is equivalent to the probability that each  $\tilde{z}_k = z_k - z_j < 0$ , as follows

$$\begin{aligned} p_j &= P(\tilde{z}_1 < 0, \dots, \tilde{z}_{j-1} < 0, \tilde{z}_{j+1} < 0, \dots, \tilde{z}_m < 0) \\ &= P(\zeta_1 < u_1, \dots, \zeta_{j-1} < u_{j-1}, \zeta_{j+1} < u_{j+1}, \dots, \zeta_m < u_m) \\ &= P(\zeta_1 < u_1) P(\zeta_2 < u_2 | \zeta_1 < u_1) \cdots \\ &\quad \cdots P(\zeta_m < u_m | \zeta_1 < u_1, \dots, \zeta_{j-1} < u_{j-1}, \zeta_{j+1} < u_{j+1}, \dots, \zeta_{m-1} < u_{m-1}), \end{aligned}$$

where  $u_i = u_i(\zeta_1, \dots, \zeta_{i-1}) = -(\nu_{(j)i} + \sum_{k=1}^{i-1} L_{ki}\zeta_k)/L_{ii}$ . Thus, the integral involving a  $(m-1)$ -variate normal density (5.17) is turned into a product of  $m-1$  univariate normal cdfs, which can be computed fairly efficiently in modern computer systems.

As an aside, the GHK probability simulator, can be used to sample from a truncated multivariate normal distribution:

- Draw  $\tilde{\zeta}_1 \sim {}^t\text{N}(0, 1, -\infty, u_1)$ .
- Draw  $\tilde{\zeta}_2 \sim {}^t\text{N}(0, 1, -\infty, \tilde{u}_2)$ , where  $\tilde{u}_2 = u_2(\tilde{\zeta}_1)$ .
- ...
- Draw  $\tilde{\zeta}_{j-1} \sim {}^t\text{N}(0, 1, -\infty, \tilde{u}_{j-1})$ , where  $\tilde{u}_{j-1} = u_{j-1}(\tilde{\zeta}_1, \dots, \tilde{\zeta}_{j-2})$ .
- Draw  $\tilde{\zeta}_{j+1} \sim {}^t\text{N}(0, 1, -\infty, \tilde{u}_{j+1})$ , where  $\tilde{u}_{j+1} = u_{j+1}(\tilde{\zeta}_1, \dots, \tilde{\zeta}_{j-1})$ .
- ...
- Draw  $\tilde{\zeta}_m \sim {}^t\text{N}(0, 1, -\infty, \tilde{u}_m)$ , where  $\tilde{u}_m = u_m(\tilde{\zeta}_1, \dots, \tilde{\zeta}_{j-1}, \tilde{\zeta}_{j+1}, \dots, \tilde{\zeta}_{m-1})$ .

Then,  $\tilde{z} = \boldsymbol{\nu}_{(j)}\mathbf{L}\tilde{\boldsymbol{\zeta}}$  will be distributed according to  $\text{N}_{m-1}(\boldsymbol{\nu}_{(j)}, \boldsymbol{\Omega}_{(j)})$ . Any quantity of interest, e.g.  $\text{Er}(\tilde{z})$ , can then be estimated by the sample mean. In the variational algorithm, we require quantities such as first and second moments and also the entropy of a truncated multivariate normal distribution. Alternative methods are also discussed in the appendix.

Finally, a point on independent probit models. As we alluded to earlier in the chapter, the class probabilities condense to a unidimensional integral involving products of normal cdfs (see (5.11)) if  $\boldsymbol{\Psi}$  is diagonal. While this represents a massive simplification, care should be taken when dealing with the formula in (5.11). When at least one of the normal cdfs in the product is extremely small, this can cause loss of significance due to floating-point errors. In the **iprobit** package, the product of normal cdfs is handled as a sum on the log scale to avoid this issue.

### 5.7.2 Computational complexity of the CAVI algorithm

This is where talk about computational complexity. Of course,  $O(n^3)$  (at least) for binary, otherwise  $O(mn^3)$  in general, although can be  $O((m-1)n^3)$ . Worst case is  $O(m^3n^3)$ , but manage to reduce this. Storage is  $O(n^2)$ . Prediction is  $O(mn^2)$ .

### 5.7.3 Difficulties faced with estimating $\Psi$

also require moments involving this truncated normal distribution.

## 5.8 Examples

## 5.9 Discussion

I-prior extended to non-normal data. Naive works good, but can be better. Simply transform the normal model through a squashing function. All the nice things about I-prior can be applied here too. Probit model variety of binary and multinomial regression models.

Laplace slow, unreliable modes. MCMC also slow. Variational has similarity to EM, but advantageous: easier to calculate posterior s.d., ability to do inference on transformed parameters.

As with the normal model, storage and time requirements slow. again, look to machine learning. improvements in variational algorithm.

Extend to include class-specific covariates.

improvement in calculating the normal integral? Need to see timing where takes longest

In terms of similarity to other works, the generalised additive models (GAMs) of [Hastie and Tibshirani, 1986](#) comes close. The setup of GAMs is near identical to the I-probit model, although estimation is done differently. GAMs do not assume smooth functions from any RKHS, but instead estimates the  $f$ 's using a local scoring method or a local likelihood method. Kernel methods for classification are extremely popular in computer science and machine learning; examples include support vector machines ([Schölkopf and Smola, 2002](#)) and Gaussian process classification ([Rasmussen and Williams, 2006](#)), with the latter being more closely related to the I-probit method. I-priors differ from Gaussian process priors in the specification of the covariance kernel. Gaussian process classification typically uses the logistic link function (or squashing function, to use machine learning nomenclature), and estimation is done most commonly using the Laplace approximation, but other methods such as expectation propagation ([Minka, 2001](#)) and MCMC ([Neal, 1999](#)) have been explored as well. Variational inference

for Gaussian process probit models have been studied by [Girolami and Rogers, 2006](#), with their work providing a close reference to the variational algorithm employed by us.

## 5.10 Miscellanea

### 5.10.1 A brief introduction to variational inference

Suppose that, in a fully Bayesian setting, we append the unknown model parameters to the latent variables to form  $\mathbf{z} = \{\mathbf{y}^*, \mathbf{w}, \theta\}$ . The crux of variational inference is this: find a suitably close distribution function  $q(\mathbf{z})$  that approximates the true posterior  $p(\mathbf{z}|\mathbf{y})$ , where closeness here is defined in the Kullback-Leibler divergence sense,

$$\text{KL}(q||p) = \int \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{y})} q(\mathbf{z}) \, d\mathbf{z}.$$

One may then show that log marginal density (the log of the intractable integral) holds the following bound:

$$\begin{aligned} \log p(\mathbf{y}) &= \log p(\mathbf{y}, \mathbf{z}) - \log p(\mathbf{z}|\mathbf{y}) \\ &= \int \left\{ \log \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z})} - \log \frac{p(\mathbf{z}|\mathbf{y})}{q(\mathbf{z})} \right\} q(\mathbf{z}) \, d\mathbf{z} \\ &= \mathcal{L}(q) + \text{KL}(q||p) \\ &\geq \mathcal{L}(q) \end{aligned} \tag{5.18}$$

{eq:varbound}

since the KL divergence is a non-negative quantity. The functional  $\mathcal{L}(q)$  given by

$$\begin{aligned} \mathcal{L}(q) &= \int \log \frac{p(\mathbf{y}, \mathbf{z})}{q(\mathbf{z})} q(\mathbf{z}) \, d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{y}, \mathbf{z})] + H(q), \end{aligned} \tag{5.19}$$

{eq:elbo1}

where  $H$  is the entropy functional, is known as the *evidence lower bound* (ELBO), which serves as the proxy objective function in the likelihood maximisation problem. Evidently, the closer  $q$  is to the true  $p$ , the better, and this is achieved by maximising  $\mathcal{L}$ , or equivalently, minimising the KL divergence<sup>3</sup> from  $p$  to  $q$ . Note that the bound (5.18) achieves equality if and only if  $q \equiv p$ , but of course the true form of the posterior

---

<sup>3</sup>The astute reader will realise that  $\text{KL}(q||p)$  is impossible to compute, since one does not know the true distribution  $p(\mathbf{z}|\mathbf{y})$ . Efforts are concentrated on maximising the ELBO instead.



is unknown to us. Maximising  $\mathcal{L}(q)$  or minimising  $\text{KL}(q\|p)$  with respect to the density  $q$  is a problem of calculus of variations, which incidentally, is where variational inference takes its name.

Maximising  $\mathcal{L}$  over all possible density functions  $q$  is not possible without considering certain constraints. Two such constraints are described. The first, is to make a distributional assumption regarding  $q$ , for which it is parameterised by  $\nu$ . For instance, we might choose the closest normal distribution to the posterior  $p(\mathbf{z}|\mathbf{y})$  in terms of KL divergence. In this case, the task is to find optimal mean and variance parameters of a normal distribution.

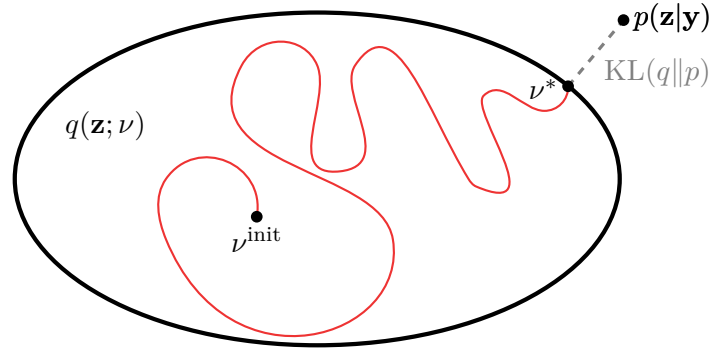


Figure 5.3: Schematic view of variational inference. The aim is to find the closest distribution  $q$  (parameterised by a variational parameter  $\nu$ ) to  $p$  in terms of KL divergence within the set of variational distributions, represented by the ellipse.

The second type of constraint, and the one considered in this thesis, is simply an assumption that the approximate posterior  $q$  factorises into  $M$  disjoint factors. Supposing that the elements of  $\mathbf{z}$  may indeed be partitioned into  $M$  disjoint groups  $\mathbf{z} = (z^{(1)}, \dots, z^{(M)})$ , then the structure

$$q(\mathbf{z}) = \prod_{k=1}^M q_k(z^{(k)})$$

for  $q$  is considered. This factorised form of variational inference is known in the statistical physics literature as the *mean-field theory* (Itzykson and Drouffe, 1991).

Denote by  $\tilde{q}$  the distributions which minimise the Kullbeck-Leibler divergence (maximise the variational lower bound). By appealing to Bishop (2006, equation 10.9, p.

466), we find that for each  $\xi \in \{\mathbf{y}^*, \mathbf{w}, \theta\} =: \mathcal{Z}$ ,  $\tilde{q}$  satisfies

$$\log \tilde{q}(\xi) = \mathbb{E}_{-\xi}[\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta)] + \text{const.} \quad (5.20)$$

{eq:qtilde}

where expectation of the log joint density of  $(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta)$  is taken with respect to all of the unknowns  $\mathcal{Z}$  except the one currently in consideration, under their respective  $q$  densities. Estimates of the latent variables and parameters are then obtained by taking the mean of their respective approximate posterior distribution.

In practice, rather than an explicit calculation of the normalising constant, one simply needs to inspect (5.20) to recognise it as a known log-density function, which is the case when exponential family distributions are considered. That is, suppose that each complete conditional  $p(\xi | \mathcal{Z}_{-\xi}, \mathbf{y})$  follows an exponential family distribution,

$$p(\xi | \mathcal{Z}_{-\xi}, \mathbf{y}) = B(\xi) \exp(\langle \zeta_{\xi}(\mathcal{Z}_{-\xi}, \mathbf{y}), \xi \rangle - A(\zeta_{\xi})).$$

Then, from (5.20),

$$\begin{aligned} \tilde{q}(\xi) &\propto \exp(\mathbb{E}_{-\xi}[\log p(\xi | \mathcal{Z}_{-\xi}, \mathbf{y})]) \\ &= \exp\left(\log B(\xi) + \mathbb{E}\langle \zeta_{\xi}(\mathcal{Z}_{-\xi}, \mathbf{y}), \xi \rangle - \mathbb{E}[A(\zeta_{\xi})]\right) \\ &\propto B(\xi) \exp \mathbb{E}\langle \zeta_{\xi}(\mathcal{Z}_{-\xi}, \mathbf{y}), \xi \rangle \end{aligned}$$

is also in the same exponential family. In situations where there is no closed form expression for  $\tilde{q}$ , then one resorts to sampling methods such as a Metropolis random walk to estimate quantities of interest. This stochastic step within a deterministic algorithm has been explored before in the context of a Monte Carlo EM algorithm—see [Meng and Van Dyk \(1997, §4, pp. 537–538\)](#) and references therein.

### 5.10.2 The EM algorithm is intractable—variational Bayes EM

## Appendix

### 5.11 Some distributions and their properties

This is a reference relating to the multivariate normal, matrix normal, truncated univariate and multivariate normal, Wishart, and gamma distributions which are collated

from various sources for convenience. Of interest are probability density functions, first and second moments, and entropy (as defined in [Chapter 3](#)).

### 5.11.1 Multivariate normal distribution

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.**  $H(p) = \frac{1}{2} \log|2\pi e \Sigma| = \frac{d}{2}(1 + \log 2\pi) + \frac{1}{2} \log|\Sigma|$ .

**Lemma 5.1** (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 5.2.** 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 \(2008, §8.1.6\)](#). □

### 5.11.2 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.**  $H(p) = \frac{1}{2} \log |2\pi e(\Psi \otimes \Sigma)| = \frac{nm}{2}(1 + \log 2\pi) + \frac{1}{2} \log |\Sigma|^m |\Psi|^n$ .

In the above, ‘ $\otimes$ ’ denotes the Kronecker matrix product defined by

$$\Psi \otimes \Sigma = \begin{pmatrix} \Psi_{11}\Sigma & \Psi_{12}\Sigma & \cdots & \Psi_{1m}\Sigma \\ \Psi_{21}\Sigma & \Psi_{22}\Sigma & \cdots & \Psi_{2m}\Sigma \\ \vdots & \vdots & \ddots & \vdots \\ \Psi_{m1}\Sigma & \Psi_{m2}\Sigma & \cdots & \Psi_{mm}\Sigma \end{pmatrix} \in \mathbb{R}^{nm \times nm}.$$

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  $\lambda$ ,

$$\begin{aligned} A \otimes (B + C) &= A \otimes B + A \otimes C \\ (A + B) \otimes C &= A \otimes C + B \otimes C \\ \lambda A \otimes B &= A \otimes \lambda B = \lambda(A \otimes B) \\ (A \otimes B) \otimes C &= A \otimes (B \otimes C) \end{aligned}$$

- **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 vectorisation operation ‘vec’ stacks the columns of the matrices into one long vector, for instance,

$$\text{vec } \Psi = (\Psi_{11}, \dots, \Psi_{m1}, \Psi_{12}, \dots, \Psi_{m2}, \dots, \Psi_{1m}, \dots, \Psi_{mm})^\top \in \mathbb{R}^{m \times m}.$$

**Lemma 5.3** (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.  $\square$

Some useful properties of the matrix normal distribution are listed:

- **Expected values.**

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

- **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  $A X B \sim \text{MN}_{a,b}(\mu^\top, A \Sigma A^\top, B^\top \Psi B)$ .
- **Iid.** If  $X_i \stackrel{\text{iid}}{\sim} \text{N}_m(\mu, \Psi)$  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  $X \sim \text{MN}(1_n \mu^\top, I_n, \Psi)$ .

### 5.11.3 Truncated univariate normal distribution

Let  $X \sim N(\mu, \sigma^2)$  with  $X$  lying in the interval  $(a, b)$ . 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.**  $p(X|\mu, \sigma, a, b) = C^{-1}(2\pi\sigma^2)^{-1/2}e^{-\frac{1}{2\sigma^2}(X-\mu)^2} = \sigma C^{-1}\phi\left(\frac{x-\mu}{\sigma}\right)$ .

- **Moments.**

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

- **Entropy.**

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

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

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  ${}^tN(0, \sigma^2, 0, +\infty)$  is called the *half-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,  $C = \Phi(-\mu/\sigma) - \Phi(-\infty) = 1 - \Phi(\mu/\sigma)$ .

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.

#### 5.11.4 Truncated multivariate normal distribution

Consider the restriction of  $X \sim N_d(\mu, \Sigma)$  to a convex subset<sup>4</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  $EX$  (mean),  $EX^2$  (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) and Damien and Walker (2001). In the latter, the authors explore a simple Gibbs-based approach that is easy to implement in practice. 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 distribution, the full conditionals of  $X_j$  given  $X_{-j}$  is

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

<sup>4</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.



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  $\Sigma^{-1}$ . 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 integral of products of normal cdfs. We present some results that we have not previously seen before elsewhere.

**Lemma 5.4.** *Let  $X \sim \text{N}_d(\mu, \Sigma, \mathcal{C}_j)$ , 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:*

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

thm:contrun  
cn

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  $Z \sim N(0, 1)$ .

(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

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

(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 ?? for the proof. □

## 5.12 Derivation of the CAVI algorithm

Let  $\mathcal{Z} = \{\mathbf{y}^*, \mathbf{w}, \boldsymbol{\alpha}, \eta, \boldsymbol{\Psi}\}$ . Approximate the posterior for  $\mathcal{Z}$  by a mean-field variational distribution

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

The first line is by assumption, while the second line follows from an induced factorisation on the latent propensities, as we will see later. If needed, we also assume that  $q(\eta)$  factorises into its constituents components. Recall that, for each  $\xi \in \mathcal{Z}$ , the optimal mean-field variational density  $\tilde{q}$  for  $\xi$  satisfies

$$\log \tilde{q}(\xi) = \mathbb{E}_{-\xi}[\log p(\mathbf{y}, \mathcal{Z})] + \text{const.} \quad (5.20)$$

Write  $\mathbf{f} = \mathbf{H}_\eta \mathbf{w} \in \mathbb{R}^{n \times m}$ . The joint likelihood  $p(\mathbf{y}, \mathcal{Z})$  is given by

$$\begin{aligned} p(\mathbf{y}, \mathcal{Z}) &= p(\mathbf{y} | \mathcal{Z})p(\mathcal{Z}) \\ &= p(\mathbf{y} | \mathbf{y}^*)p(\mathbf{y}^* | \boldsymbol{\alpha}, \mathbf{w}, \eta, \boldsymbol{\Psi})p(\mathbf{w} | \boldsymbol{\Psi})p(\eta)p(\boldsymbol{\Psi})p(\boldsymbol{\alpha}). \end{aligned}$$

For reference, the 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.

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

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

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

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

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

where  $\mathbf{y}_i^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}_i \in \mathbb{R}^m$  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  $\mathbf{y}_i^*$  are independent multivariate normal with mean  $\boldsymbol{\mu}_i$  and variance  $\boldsymbol{\Psi}^{-1}$ .

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

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

- $p(\eta)$ . The most common scenario would be  $\eta = \{\lambda_1, \dots, \lambda_p\}$  only. In this case, choose independent normal priors for each  $\lambda_k \sim \text{N}(m_k, v_k)$ ,  $k = 1, \dots, p$ , whose pdf is

$$p(\eta) = \prod_{k=1}^p \exp \left[ -\frac{1}{2} \log 2\pi - \frac{1}{2} \log v_k - \frac{1}{2v_k} (\lambda_k - m_k)^2 \right].$$

An improper prior  $p(\eta) \propto \text{const.}$  can be used as well, and this is the same as letting  $m_k \rightarrow 0$  and  $v_k \rightarrow 0$ . The resulting posterior will be proper. If  $\eta$  contains other parameters as well, such as the Hurst coefficient  $\gamma \in (0, 1)$ , SE lengthscale  $l > 0$  or polynomial offset  $c > 0$ , then appropriate priors should be used to match the support of the parameter. Choices include  $p(\gamma) = \mathbf{1}(\gamma \in (0, 1))$  and  $l, c \sim \Gamma(a, b)$ .

- $p(\boldsymbol{\Psi})$ . Our analysis shows that regardless of prior choice of  $\boldsymbol{\Psi}$ , be it in the full or independent I-probit model, the posterior for  $\boldsymbol{\Psi}$  will not be of a recognisable form. Without giving too much thought, assume an improper prior on  $\boldsymbol{\Psi}$ , i.e.  $p(\boldsymbol{\Psi}) \propto \text{const.}$

- $\mathbf{p}(\boldsymbol{\alpha})$ . Choose independent normal priors for the intercept,  $\alpha_j \sim \mathcal{N}(a_j, A_j)$  for  $j = 1, \dots, m$ . The pdf is

$$p(\boldsymbol{\alpha}) = \prod_{j=1}^m \exp \left[ -\frac{1}{2} \log 2\pi - \frac{1}{2} \log A_j - \frac{1}{2A_j} (\alpha_j - a_j)^2 \right].$$

*Remark 5.1.* The priors on the parameters  $\{\boldsymbol{\alpha}, \eta\}$  can be set to very vague or even improper priors, and the resulting posterior will still yield a proper distribution. Using improper priors eases the algebra slightly. For the precision matrix  $\boldsymbol{\Psi}$ , it is best to stick with the Wishart prior to avoid positive-definite issues, unless the independent I-probit model is used, in which case Jeffreys' prior for the precisions  $p(\sigma_j^{-2}) \propto \sigma_j^2$  is a convenient choice.

### 5.12.1 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  $j \in \{1, \dots, m\}$ . The mean-field density  $q(\mathbf{y}_i^*)$  for each  $i = 1, \dots, n$  is found to be

$$\begin{aligned} \log \tilde{q}(\mathbf{y}_i^*) &= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \mathbb{E}_{\mathcal{Z} \setminus \{\mathbf{y}^*\} \sim q} \left[ -\frac{1}{2} (\mathbf{y}_i^* - \boldsymbol{\mu}_i)^\top \boldsymbol{\Psi} (\mathbf{y}_i^* - \boldsymbol{\mu}_i) \right] + \text{const.} \\ &= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \left[ -\frac{1}{2} (\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i)^\top \tilde{\boldsymbol{\Psi}} (\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i) \right] + \text{const.} \quad (\star) \\ &\equiv \begin{cases} \phi(\mathbf{y}_i^* | \tilde{\boldsymbol{\mu}}_i, \tilde{\boldsymbol{\Psi}}) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

where  $\tilde{\boldsymbol{\mu}}_i = \mathbb{E} \boldsymbol{\alpha} + (\mathbb{E} \mathbf{H}_\eta \mathbb{E} \mathbf{w})_i$ , and expectations are taken under the optimal mean-field distribution  $\tilde{q}$ . The distribution  $q(\mathbf{y}_i^*)$  is a truncated  $m$ -variate normal distribution such that the  $j$ 'th component is always largest. Unfortunately, the expectation of this distribution cannot be found in closed-form, and must be approximated by techniques such as Monte Carlo integration. If, however, the independent I-probit model is used and  $\tilde{\boldsymbol{\Psi}}$  is diagonal, then [Lemma X](#) provides a simplification.

*Remark 5.2.* In  $(\star)$  above, we needn't consider the second order terms in the expectations because they do not involve  $\mathbf{y}^*$  and can be absorbed into the constant. To see this,

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

We will see this occurring a lot later on and we shall take note of this fact.

### 5.12.2 Derivation of $\tilde{q}(\mathbf{w})$

The terms involving  $\mathbf{w}$  in (5.20) are the  $p(\mathbf{y}^*|\boldsymbol{\alpha}, \mathbf{w}, \eta, \boldsymbol{\Psi})$  and  $p(\mathbf{w}|\boldsymbol{\Psi})$  terms, and 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,

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

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

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

and making a first-order approximation  $(\mathbb{E} \mathbf{A})^{-1} \approx \mathbb{E}[\mathbf{A}^{-1}]^5$ ,

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

Ideally, we do not want to work with the  $nm \times nm$  matrix  $\mathbf{V}_w$ , since its inverse is expensive to compute. We can exploit the Kronecker product structure to compute the inverse efficiently. Perform an orthogonal eigendecomposition of  $\mathbf{H}_\eta$  to obtain  $\mathbf{H}_\eta = \mathbf{V} \mathbf{U} \mathbf{V}^\top$  and of  $\boldsymbol{\Psi}$  to obtain  $\boldsymbol{\Psi} = \mathbf{Q} \mathbf{P} \mathbf{Q}^\top$ . This process takes  $O(n^3 + m^3) \approx O(n^3)$  time if  $m \ll n$ . Then, manipulate  $\mathbf{V}_w^{-1}$  as follows (for clarity, we drop the tildes from the notations):

$$\begin{aligned}
 \mathbf{V}_w^{-1} &= (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta^2) + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \\
 &= (\mathbf{Q} \mathbf{P} \mathbf{Q}^\top \otimes \mathbf{V} \mathbf{U}^2 \mathbf{V}^\top) + (\mathbf{Q} \mathbf{P}^{-1} \mathbf{Q}^\top \otimes \mathbf{V} \mathbf{V}^\top) \\
 &= (\mathbf{Q} \otimes \mathbf{V}) (\mathbf{P} \otimes \mathbf{U}^2) (\mathbf{Q}^\top \otimes \mathbf{V}^\top) + (\mathbf{Q} \otimes \mathbf{V}) (\mathbf{P}^{-1} \otimes \mathbf{I}_n) (\mathbf{Q}^\top \otimes \mathbf{V}^\top) \\
 &= (\mathbf{Q} \otimes \mathbf{V}) (\mathbf{P} \otimes \mathbf{U}^2 + \mathbf{P}^{-1} \otimes \mathbf{I}_n) (\mathbf{Q}^\top \otimes \mathbf{V}^\top)
 \end{aligned}$$

Its inverse is

$$\begin{aligned}
 \mathbf{V}_w &= (\mathbf{Q}^\top \otimes \mathbf{V}^\top)^{-1} (\mathbf{P} \otimes \mathbf{U}^2 + \mathbf{P}^{-1} \otimes \mathbf{I}_n)^{-1} (\mathbf{Q} \otimes \mathbf{V})^{-1} \\
 &= (\mathbf{Q} \otimes \mathbf{V}) (\mathbf{P} \otimes \mathbf{U}^2 + \mathbf{P}^{-1} \otimes \mathbf{I}_n)^{-1} (\mathbf{Q}^\top \otimes \mathbf{V}^\top)
 \end{aligned}$$

which is easy to compute since the middle term is an inverse of diagonal matrices.

<sup>5</sup>Groves and Rothenberg (1969) show that  $\mathbb{E}[\mathbf{A}^{-1}] = (\mathbb{E} \mathbf{A})^{-1} + \mathbf{B}$ , where  $\mathbf{B}$  is a positive-definite matrix.

In the case of the I-probit model, where  $\Psi = \text{diag}(\sigma_1^{-2}, \dots, \sigma_m^{-2})$ , then the covariance  $\mathbf{V}_w$  takes a simpler form. Specifically, it has the block diagonal structure:

$$\begin{aligned}\mathbf{V}_w &= (\text{diag}(\tilde{\sigma}_1^{-2}, \dots, \tilde{\sigma}_m^{-2}) \otimes \tilde{\mathbf{H}}_\eta^2 + (\text{diag}(\tilde{\sigma}_1^2, \dots, \tilde{\sigma}_m^2) \otimes \mathbf{I}_n)^{-1} \\ &= \text{diag}\left((\tilde{\sigma}_1^{-2} \tilde{\mathbf{H}}_\eta^2 + \tilde{\sigma}_1^2 \mathbf{I}_n)^{-1}, \dots, (\tilde{\sigma}_m^{-2} \tilde{\mathbf{H}}_\eta^2 + \tilde{\sigma}_m^2 \mathbf{I}_n)^{-1}\right) \\ &=: \text{diag}(\tilde{\mathbf{V}}_{w_1}, \dots, \tilde{\mathbf{V}}_{w_m}).\end{aligned}$$

The mean  $\tilde{\mathbf{w}}$  in matrix form is

$$\begin{aligned}\tilde{\mathbf{w}} &= \tilde{\mathbf{V}}_w (\text{diag}(\tilde{\sigma}_1^{-2}, \dots, \tilde{\sigma}_m^{-2}) \otimes \tilde{\mathbf{H}}_\eta) (\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top) \\ &= \text{diag}(\tilde{\mathbf{V}}_{w_1}, \dots, \tilde{\mathbf{V}}_{w_m}) \text{diag}(\tilde{\sigma}_1^{-2} \tilde{\mathbf{H}}_\eta, \dots, \tilde{\sigma}_m^{-2} \tilde{\mathbf{H}}_\eta) (\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top) \\ &= \text{diag}(\tilde{\sigma}_1^{-2} \tilde{\mathbf{V}}_{w_1} \tilde{\mathbf{H}}_\eta, \dots, \tilde{\sigma}_m^{-2} \tilde{\mathbf{V}}_{w_m} \tilde{\mathbf{H}}_\eta) (\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top) \\ &= \begin{pmatrix} \tilde{\mathbf{w}}_{\cdot 1} & \dots & \tilde{\mathbf{w}}_{\cdot m} \end{pmatrix} \\ &= \begin{pmatrix} \tilde{\sigma}_1^{-2} \tilde{\mathbf{V}}_{w_1} \tilde{\mathbf{H}}_\eta (\tilde{\mathbf{y}}_{\cdot 1}^* - \tilde{\alpha}_1 \mathbf{1}_n) & \dots & \tilde{\sigma}_m^{-2} \tilde{\mathbf{V}}_{w_m} \tilde{\mathbf{H}}_\eta (\tilde{\mathbf{y}}_{\cdot m}^* - \tilde{\alpha}_m \mathbf{1}_n) \end{pmatrix}.\end{aligned}$$

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

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

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

$$\begin{aligned}\text{tr}(\mathbf{C} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]) &= \sum_{i,j=1}^m \mathbf{C}_{ij} \mathbf{E}[\mathbf{w}^\top \mathbf{D} \mathbf{w}]_{ij} \\ &= \sum_{i,j=1}^m \mathbf{C}_{ij} \mathbf{E}[\mathbf{w}_{\cdot i}^\top \mathbf{D} \mathbf{w}_{\cdot j}].\end{aligned} \tag{5.21}$$

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 (5.21) leads us to

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

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

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

### 5.12.3 Derivation of $\tilde{q}(\eta)$

By looking at only the terms involving  $\eta$  in (5.20), we deduce that  $\tilde{q}$  for  $\eta$  satisfies

$$\begin{aligned} \log \tilde{q}(\eta) &= -\frac{1}{2} \text{tr} \mathbf{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ (\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top - \mathbf{H}_\eta \mathbf{w}) \boldsymbol{\Psi} (\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top - \mathbf{H}_\eta \mathbf{w})^\top \right] + \log p(\eta) \\ &\quad + \text{const.} \\ &= -\frac{1}{2} \text{tr} \mathbf{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left( \boldsymbol{\Psi} \mathbf{w}^\top \mathbf{H}_\eta^2 \mathbf{w} - 2 \boldsymbol{\Psi} \mathbf{w}^\top \mathbf{H}_\eta (\mathbf{y}^* - \boldsymbol{\alpha}) \right) + \log p(\eta) + \text{const.} \\ &= -\frac{1}{2} \text{tr} \left( \tilde{\boldsymbol{\Psi}} \mathbf{E}[\mathbf{w}^\top \mathbf{H}_\eta^2 \mathbf{w}] - 2 \tilde{\boldsymbol{\Psi}} \tilde{\mathbf{w}}^\top \mathbf{H}_\eta (\tilde{\mathbf{y}}^* - \tilde{\boldsymbol{\alpha}}) \right) + \log p(\eta) + \text{const.} \end{aligned}$$

with some appropriate prior  $p(\eta)$ . In general, this does not have a recognisable form in  $\eta$ , especially when it is not linearly dependent on the kernel matrix. This happens when considering parameters other than the scales of the RKHSs. Our interest would be to obtain  $\tilde{\mathbf{H}}_\eta := \mathbf{E}_{\eta \sim q} \mathbf{H}_\eta$  and  $\tilde{\mathbf{H}}_\eta^2 := \mathbf{E}_{\eta \sim q} \mathbf{H}_\eta^2$ . We use a Metropolis random-walk algorithm to obtain these quantities, as detailed in the algorithm below.

Now consider the case where  $\eta = \{\lambda_1, \dots, \lambda_p\}$  (RKHS scale parameters only), and the scenario described in the exponential family EM algorithm of [Section 4.3.3](#) applies. In particular, for  $k = 1, \dots, p$ , we can decompose the kernel matrix as  $\mathbf{H}_\eta = \lambda_k \mathbf{R}_k + \mathbf{S}_k$  and its square as  $\mathbf{H}_\eta^2 = \lambda_k^2 \mathbf{R}_k^2 + \lambda_k \mathbf{U}_k + \mathbf{S}_k^2$ . Then, for  $j = 1, \dots, m$ , assuming each of

**Algorithm 2** Metropolis random-walk to sample  $\eta$

- 1: **inputs**  $\tilde{\alpha}$ ,  $\tilde{\mathbf{w}}$ ,  $\tilde{\Psi}$ , and  $s$  Metropolis sampling s.d.
- 2: **initialise**  $\eta^{(0)} \in \mathbb{R}^q$  and  $t \leftarrow 0$
- 3: **for**  $t = 1, \dots, T$  **do**
- 4:     Draw  $\eta^* \sim N_q(\eta^{(t)}, s^2)$
- 5:     Accept/reject proposal state, i.e.

$$\eta^{(t+1)} \leftarrow \begin{cases} \eta^* & \text{if } u \sim \text{Unif}(0, 1) < \pi_{\text{acc}} \\ \eta^{(t)} & \text{otherwise} \end{cases}$$

where

$$\pi_{\text{acc}} = \min \left( 1, \exp \left( \log \tilde{q}(\eta^*) - \log \tilde{q}(\eta^{(t)}) \right) \right).$$

6: **end for**

7:  $\tilde{\mathbf{H}}_\eta \leftarrow \frac{1}{T} \sum_{i=1}^T \mathbf{H}_{\eta^{(i)}}$  and  $\tilde{\mathbf{H}}_\eta^2 \leftarrow \frac{1}{T} \sum_{i=1}^T \mathbf{H}_{\eta^{(i)}}^2$

the  $q(\lambda_k)$  densities are independent of each other, we find that

$$\begin{aligned} \log \tilde{q}(\lambda_k) &= \mathbb{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ -\frac{1}{2} \text{tr} \left( (\mathbf{y}^* - \boldsymbol{\mu}) \Psi (\mathbf{y}^* - \boldsymbol{\mu})^\top \right) \right] - \frac{1}{2v_k^2} (\lambda_k - m_k)^2 + \text{const.} \\ &= -\frac{1}{2} \text{tr} \mathbb{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ \Psi \mathbf{w}^\top \mathbf{H}_\eta^2 \mathbf{w} - 2\Psi (\mathbf{y}^* - \mathbf{1}\boldsymbol{\alpha}^\top)^\top \mathbf{H}_\eta \mathbf{w} \right] \\ &\quad - \frac{1}{2v_k^2} (\lambda_k - m_k)^2 + \text{const.} \\ &= -\frac{1}{2} \text{tr} \mathbb{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ \Psi \mathbf{w}^\top (\lambda_k^2 \mathbf{R}_k^2 + \lambda_k \mathbf{U}_k) \mathbf{w} - 2\Psi (\mathbf{y}^* - \mathbf{1}\boldsymbol{\alpha}^\top)^\top (\lambda_k \mathbf{R}_k) \mathbf{w} \right] \\ &\quad - \frac{1}{2v_k^2} (\lambda_k^2 - 2m_k \lambda_k) + \text{const.} \\ &= -\frac{1}{2} \text{tr} \mathbb{E}_{\mathcal{Z} \setminus \{\eta\} \sim q} \left[ \lambda_k^2 \Psi \mathbf{w}^\top \mathbf{R}_k^2 \mathbf{w} - 2\lambda_k \left( \Psi (\mathbf{y}^* - \mathbf{1}\boldsymbol{\alpha}^\top)^\top \mathbf{R}_k \mathbf{w} - \frac{1}{2} \Psi \mathbf{w}^\top \mathbf{U}_k \mathbf{w} \right) \right] \\ &\quad - \frac{1}{2} \left( \frac{1}{v_k^2} \lambda_k^2 - 2\frac{m_k}{v_k^2} \lambda_k \right) + \text{const.} \\ &= -\frac{1}{2} \left[ \lambda_k^2 \overbrace{(\text{tr}(\tilde{\Psi} \mathbb{E}[\mathbf{w}^\top \mathbf{R}_k^2 \mathbf{w}]) + v_k^{-2})}^{c_k} \right. \\ &\quad \left. - 2\lambda_k \overbrace{\left( \text{tr} \left( \tilde{\Psi} (\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top)^\top \mathbf{R}_k \tilde{\mathbf{w}} - \frac{1}{2} \tilde{\Psi} \mathbb{E}[\mathbf{w}^\top \mathbf{U}_k \mathbf{w}] \right) + m_k v_k^{-2} \right)}^{d_k} \right] \end{aligned}$$

By completing the squares, we recognise this is as the kernel of a univariate normal density. Specifically,  $\lambda_k \sim N(d_k/c_k, 1/c_k)$ . The quantity  $\tilde{\mathbf{H}}_\eta$  can be obtained by substituting

tuting  $\lambda_k \mapsto \mathbb{E}_{\lambda_k \sim q}[\lambda_k]$  in the **expression XXX**. However, in the calculation of  $\tilde{\mathbf{H}}_\eta^2$ , we must replace  $\lambda_k^2 \mapsto \mathbb{E}_{\lambda_k \sim q}[\lambda_k]^2 + \text{Var}_{\lambda_k \sim q}[\lambda_k]$  in all occurrences of square terms. This can be cumbersome, so if felt necessary, use the approximation  $\lambda_k^2 \mapsto \mathbb{E}_{\lambda_k \sim q}[\lambda_k]^2$  instead.

**Example 5.1.** Suppose  $k = 1$ , and we only have  $\lambda$  to estimate. Then,  $\mathbf{H}_\eta = \lambda \mathbf{H}$ ,  $\mathbf{R}_k = \mathbf{H}$ ,  $\mathbf{R}_k^2 = \mathbf{H}^2$ , and  $\mathbf{U}_k = \mathbf{0}$ . Suppose also we use an improper prior  $\lambda_k \propto \text{const.}$ , which is the same as having  $v_k^2 \rightarrow 0$  and  $m_k v_k^{-2} \rightarrow 0$ . The mean field distribution for  $\lambda$  is then

$$\lambda \sim \text{N} \left( \frac{\text{tr}(\tilde{\Psi}(\tilde{\mathbf{y}}^* - \mathbf{1}\tilde{\alpha}^\top)^\top \mathbf{H}\tilde{\mathbf{w}})}{\text{tr}(\tilde{\Psi} \mathbb{E}[\mathbf{w}^\top \mathbf{H}^2 \mathbf{w}])}, \frac{1}{\text{tr}(\tilde{\Psi} \mathbb{E}[\mathbf{w}^\top \mathbf{H}^2 \mathbf{w}])} \right)$$

Further, if  $\tilde{\Psi} = \tilde{\psi} \mathbf{I}_m$ , then

$$\lambda \sim \text{N} \left( \frac{\sum_{j=1}^m (\tilde{\mathbf{y}}_{\cdot j}^* - \tilde{\alpha}_j \mathbf{1})^\top \mathbf{H} \tilde{\mathbf{w}}_{\cdot j}}{\sum_{j=1}^m \text{tr}(\mathbf{H}^2 \mathbb{E}[\mathbf{w}_{\cdot j} \mathbf{w}_{\cdot j}^\top])}, \frac{1}{\sum_{j=1}^m \text{tr}(\mathbf{H}^2 \mathbb{E}[\mathbf{w}_{\cdot j} \mathbf{w}_{\cdot j}^\top])} \right)$$

which bears a resemblance to the exponential family EM algorithm solutions described in Chapter 4. Now,  $\tilde{\mathbf{H}}_\eta = \mathbb{E}[\lambda \mathbf{H}] = \tilde{\lambda} \mathbf{H}$ , and  $\tilde{\mathbf{H}}_\eta^2 = \mathbb{E}[\lambda^2 \mathbf{H}^2] = (\text{Var } \lambda + \tilde{\lambda}^2) \mathbf{H}^2$ .

#### 5.12.4 Derivation of $\tilde{q}(\Psi)$

We find that  $q(\Psi)$  satisfies

$$\begin{aligned} \log q(\Psi) &= \mathbb{E}_{\mathcal{Z} \setminus \{\Psi\} \sim q} \left[ -\frac{1}{2} \text{tr}((\mathbf{y}^* - \boldsymbol{\mu})^\top (\mathbf{y}^* - \boldsymbol{\mu}) \Psi) - \frac{1}{2} \text{tr}(\mathbf{w}^\top \mathbf{w} \Psi^{-1}) \right] \\ &\quad + \log p(\Psi) + \text{const.} \\ &= -\frac{1}{2} \text{tr} \left( \overbrace{(\mathbb{E}[(\mathbf{y}^* - \boldsymbol{\mu})^\top (\mathbf{y}^* - \boldsymbol{\mu})])}^{\mathbf{G}_1} \Psi + \overbrace{\mathbb{E}[\mathbf{w}^\top \mathbf{w}]}^{\mathbf{G}_2} \Psi^{-1} \right) \\ &\quad + \log p(\Psi) + \text{const.} \end{aligned}$$

This seems to be the pdf of  $\text{Wis}(\mathbf{G} + \mathbf{G}_1, g)$  plus the pdf of a distribution which almost resembles an inverse Wishart pdf. Unfortunately, the properties such as its moments and entropy are unknown.

The matrix  $\mathbf{G}_1$  is

$$\begin{aligned}\mathbf{G}_1 &= \mathbb{E}[(\mathbf{y}^* - \boldsymbol{\mu})^\top (\mathbf{y}^* - \boldsymbol{\mu})] \\ &= \mathbb{E}[\mathbf{y}^{*\top} \mathbf{y}^* + \boldsymbol{\alpha} \mathbf{1}_n^\top \mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{w}^\top \mathbf{H}_\eta^2 \mathbf{w} - 2\mathbf{y}^{*\top} \mathbf{1}_n \boldsymbol{\alpha}^\top - 2\mathbf{y}^{*\top} \mathbf{H}_\eta \mathbf{w} - 2\boldsymbol{\alpha} \mathbf{1}_n^\top \mathbf{H}_\eta \mathbf{w}] \\ &= \mathbb{E}[\mathbf{y}^{*\top} \mathbf{y}^*] + n \mathbb{E}[\boldsymbol{\alpha} \boldsymbol{\alpha}^\top] + \mathbb{E}[\mathbf{w}^\top \mathbf{H}_\eta \mathbf{w}] - 2(\tilde{\mathbf{y}}^{*\top} \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top + \tilde{\mathbf{y}}^{*\top} \tilde{\mathbf{H}}_\eta \tilde{\mathbf{w}} + \tilde{\boldsymbol{\alpha}} \mathbf{1}_n^\top \tilde{\mathbf{H}}_\eta \tilde{\mathbf{w}}),\end{aligned}$$

and this involves second order moments of a conically truncated multivariate normal distribution, which needs to be obtained via simulation. Meanwhile,

$$\begin{aligned}\mathbf{G}_{2,ij} &= \mathbb{E}[\mathbf{w}^\top \mathbf{w}]_{ij} \\ &= \mathbb{E}[\mathbf{w}_{\cdot i}^\top \mathbf{w}_{\cdot j}] \\ &= \tilde{\mathbf{V}}_w[i, j] + \tilde{\mathbf{w}}_{\cdot i}^\top \tilde{\mathbf{w}}_{\cdot j}.\end{aligned}$$

In the case of the independent I-probit model, we use a gamma prior on each of the precisions in the diagonal entries of  $\boldsymbol{\Psi} = \text{diag}(\psi_1, \dots, \psi_m)$ . Then, the variational density for each  $\psi_j$  is found to be

$$\begin{aligned}\log q(\psi_j) &= \mathbb{E}_{\mathcal{Z} \setminus \{\boldsymbol{\Psi}\} \sim q} \left[ \frac{n}{2} \log(\psi_1 \cdots \psi_m) - \frac{1}{2} \sum_{j=1}^m \sum_{i=1}^n \psi_j (\mathbf{y}_{ij}^* - \boldsymbol{\mu}_{ij})^2 \right] \\ &\quad + \mathbb{E}_{\mathcal{Z} \setminus \{\boldsymbol{\Psi}\} \sim q} \left[ -\frac{n}{2} \log(\psi_1 \cdots \psi_m) - \frac{1}{2} \sum_{j=1}^m \sum_{i=1}^n \psi_j^{-1} \mathbf{w}_{ij}^2 \right] \\ &\quad + \sum_{j=1}^m ((s_j - 1) \log \psi_j - r_j \psi_j) + \text{const.} \\ &= (s_j - 1) \log \psi_j - \psi_j \left( \frac{1}{2} \mathbb{E} \|\mathbf{y}_{\cdot j}^* - \boldsymbol{\mu}_{\cdot j}\|^2 + r_j \right) \\ &\quad - \psi_j^{-1} \left( \frac{1}{2} \mathbb{E} \|\mathbf{w}_{\cdot j}\|^2 \right) + \text{const.}\end{aligned}$$

which again, is a pdf of an unknown distribution. However, its posterior mode can be computed. Write  $a = -\left(\frac{1}{2} \mathbb{E} \|\mathbf{y}_{\cdot j}^* - \boldsymbol{\mu}_{\cdot j}\|^2 + r_j\right)$ ,  $b = s_j - 1$ , and  $c = \left(\frac{1}{2} \mathbb{E} \|\mathbf{w}_{\cdot j}\|^2\right)$ . Then,

$$\frac{\partial}{\partial \psi_j} \log q(\psi_j) = \frac{\partial}{\partial \psi_j} (a\psi_j + b \log \psi_j - c\psi_j^{-1}) = a + b\psi_j^{-1} + c\psi_j^{-2}$$

equated to zero means solving a quadratic equation in  $\psi_j$ . Suppose that  $p(\psi_j) \propto \text{const.}$ , then  $s_j = 1$  and  $r_j = 0$  so  $\tilde{\psi}_j$  can be solved directly to be

$$\hat{\psi}_j = \sqrt{\frac{\mathbb{E}\|\mathbf{y}_{\cdot,j}^* - \boldsymbol{\mu}_{\cdot,j}\|^2}{\mathbb{E}\|\mathbf{w}_{\cdot,j}\|^2}}.$$

If the posterior mean is close to its mode, then  $\hat{\psi}_j$  is a good approximation for  $\tilde{\psi}_j$ .

To calculate  $\mathbb{E}\|y_{\cdot,j}^* - \boldsymbol{\mu}_{\cdot,j}\|^2 = \mathbb{E}\sum_{i=1}^n (\mathbf{y}_{ij}^* - \mu_{ij})^2$ , one first needs  $\mathbb{E}(y_{ij}^* - \alpha_j - \mathbf{w}_{\cdot,j}^\top \mathbf{h}_\eta(x_i))^2$ . This, in itself, presents a challenge to compute analytically, because it requires, among other things, the second moments  $\mathbb{E}y_{ij}^{*2}$  and  $\mathbb{E}[\mathbf{w}_{\cdot,j}^\top \mathbf{h}_\eta(x_i) \mathbf{h}_\eta(x_i)^\top \mathbf{w}_{\cdot,j}]$ . Although not entirely accurate, it is simpler to use the approximation

$$\mathbb{E}\|y_{\cdot,j}^* - \boldsymbol{\mu}_{\cdot,j}\|^2 \approx \|\tilde{y}_{\cdot,j}^* - \tilde{\boldsymbol{\mu}}_{\cdot,j}\|^2.$$

(see note on [page](#)). Also, we have  $\mathbf{w}_{\cdot,j} \sim N_n(\tilde{\mathbf{w}}_{\cdot,j}, \tilde{\mathbf{V}}_{w_j})$ , and so  $\mathbb{E}\|\mathbf{w}_{\cdot,j}\|^2 = \text{tr}(\tilde{\mathbf{V}}_{w_j} + \tilde{\mathbf{w}}_{\cdot,j} \tilde{\mathbf{w}}_{\cdot,j}^\top)$ .

### 5.12.5 Derivation of $\tilde{q}(\boldsymbol{\alpha})$

Let  $\mathbf{A} = \text{diag}(A_1, \dots, A_m)$  and  $\mathbf{a} = (a_1, \dots, a_m)^\top$ . The terms involving  $\alpha_j$  in (5.20) are

$$\begin{aligned} \log q(\boldsymbol{\alpha}) &= \mathbb{E}_{\mathcal{Z} \setminus \{\boldsymbol{\alpha}\} \sim q} \left[ -\frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i^* - \boldsymbol{\alpha} - \mathbf{w}^\top \mathbf{h}_\eta(x_i))^\top \boldsymbol{\Psi} (\mathbf{y}_i^* - \boldsymbol{\alpha} - \mathbf{w}^\top \mathbf{h}_\eta(x_i)) \right] \\ &\quad - \frac{1}{2} (\boldsymbol{\alpha} - \mathbf{a})^\top \mathbf{A}^{-1} (\boldsymbol{\alpha} - \mathbf{a}) + \text{const.} \\ &= -\frac{1}{2} \left[ \boldsymbol{\alpha}^\top \overbrace{(n\boldsymbol{\Psi} + \mathbf{A}^{-1})}^{\tilde{\mathbf{A}}} \boldsymbol{\alpha} - 2 \overbrace{\left( \sum_{i=1}^n \boldsymbol{\Psi} (\tilde{\mathbf{y}}_i^* - \tilde{\mathbf{w}}^\top \tilde{\mathbf{h}}_\eta(x_i)) + \mathbf{A}^{-1} \mathbf{a} \right)}^{\tilde{\mathbf{a}}} \boldsymbol{\alpha} \right] \end{aligned}$$

which implies a normal mean-field distribution for  $\boldsymbol{\alpha}$  whose mean and variance are  $\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{A}}^{-1} \tilde{\mathbf{a}}$  and  $\tilde{\mathbf{A}}^{-1}$  respectively. If  $\boldsymbol{\Psi}$  is diagonal, the components of  $\boldsymbol{\alpha}$  would be independent.

As a remark, due to identifiability, only  $m - 1$  of these intercept are estimable. We can either put a constraint that one of the intercepts is fixed at zero, or the sum of the intercepts equals zero. The latter constraint is implemented in this thesis, and this is realised by estimating all the intercepts and then centring them.

## 5.13 Deriving the ELBO expression

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

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

As we saw earlier, the distribution of  $q(\boldsymbol{\Psi})$  is not of recognisable form. This makes computation of  $\mathbb{E} \log |\boldsymbol{\Psi}|$ ,  $\mathbb{E} \log p(\boldsymbol{\Psi})$ , and  $H[q(\boldsymbol{\Psi})]$ , which are required in the expression of the ELBO, problematic. For simplicity, we present the ELBO calculations for when  $\boldsymbol{\Psi}$  is treated to be fixed.

*Remark 5.3.* As discussed, given the latent propensities  $\mathbf{y}^*$ , the pdf of  $\mathbf{y}$  is degenerate and hence can be disregarded.

*Remark 5.4.* When using improper priors for the hyperparameters, i.e.  $p(\eta, \boldsymbol{\alpha}) \propto \text{const.}$ , then these terms can be disregarded.

### 5.13.1 Terms involving distributions of $\mathbf{y}^*$

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

where  $C_i$  is the normalising constant for the distribution of multivariate truncated normal  $\mathbf{y}_{i\cdot}$ .

Notes:

1.  $p(\mathbf{y}_{i\cdot}^*)$  is the pdf of  $N(\boldsymbol{\mu}_{i\cdot}, \boldsymbol{\Psi}^{-1})$ , and  $q(\mathbf{y}_{i\cdot}^*)$  is the pdf of  ${}^tN(\tilde{\boldsymbol{\mu}}_{i\cdot}, \boldsymbol{\Psi}^{-1}, \mathcal{C}_{y_i})$ , where  $\boldsymbol{\mu}_{i\cdot} = \boldsymbol{\alpha} + \mathbf{w}^\top \mathbf{h}_\eta(x_i) \in \mathbb{R}^m$ .
2. It is simpler to use the approximation

$$\mathbb{E}(\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot}) \approx \mathbb{E}(\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot}). \quad (5.22)$$

rather than work out the actual quantity, which is

$$\mathbb{E}(\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \boldsymbol{\mu}_{i\cdot}) = \mathbb{E}(\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot})^\top \boldsymbol{\Psi} (\mathbf{y}_{i\cdot}^* - \tilde{\boldsymbol{\mu}}_{i\cdot}) + \text{tr}(\boldsymbol{\Psi} \text{Var } \boldsymbol{\mu}_{i\cdot}) \quad (5.23)$$

where  $\text{Var } \boldsymbol{\mu}_{i\cdot} = \text{Var } \boldsymbol{\alpha} + \text{Var } \mathbf{w}^\top \mathbf{h}_\eta(x_i)$ , obtained by taking expectations with respect to everything except  $\mathbf{y}_{i\cdot}^*$ . The first term is a diagonal matrix of the posterior variances of the intercepts. The second term is where things get complicated. Let  $\boldsymbol{\Omega}_i = \text{Var } \mathbf{w}^\top \mathbf{h}_\eta(x_i)$ . Then  $\boldsymbol{\Omega}_{i,kj} \approx \text{Cov}(\mathbf{w}_{\cdot k}^\top \mathbf{h}_\eta(x_i), \mathbf{w}_{\cdot j}^\top \mathbf{h}_\eta(x_i)) = \mathbf{h}_\eta(x_i)^\top \tilde{\mathbf{V}}_w[k, j] \mathbf{h}_\eta(x_i)$ . So

$$\text{tr}(\boldsymbol{\Psi} \boldsymbol{\Omega}_i) \approx \sum_{k,j=1}^m \boldsymbol{\Psi}_{kj} \mathbf{h}_\eta(x_i)^\top \tilde{\mathbf{V}}_w[k, j] \mathbf{h}_\eta(x_i)$$

However, we know that  $\text{Var } XY = \mathbb{E} X^2 Y^2 - (\mathbb{E} XY)^2 = \text{Var } X \text{Var } Y + \text{Var } X (\mathbb{E} Y)^2 + \text{Var } Y (\mathbb{E} X)^2$ , so there is actually some covariance terms which need to be considered, and these are not so easily computed. In practice, we find that using (5.22) gives satisfactory results as far as determining convergence for the variational algorithm goes.

### 5.13.2 Terms involving distributions of $\mathbf{w}$

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

Notes:

1.  $p(\mathbf{w})$  is the pdf of  $\text{MN}(\mathbf{0}, \mathbf{I}_n, \Psi)$ , and  $q(\mathbf{w})$  is the pdf of  $\text{N}(\text{vec } \tilde{\mathbf{w}}, \tilde{\mathbf{V}}_w)$ .
2.  $\tilde{\mathbf{V}}_w[j, j]$  are the  $n \times n$  sub matrices along the diagonal of  $\tilde{\mathbf{V}}_w$ .

### 5.13.3 Terms involving distributions of $\eta$

If no closed-form expression for  $q(\eta)$  is found, then the expression  $\mathbb{E}[\log p(\eta) - q(\eta)]$  must be obtained by sampling methods. Otherwise, consider the case where  $\eta = \{\lambda_1, \dots, \lambda_p\}$ . Then, the contribution to the ELBO is

$$\begin{aligned} \mathbb{E} \log p(\lambda_1, \dots, \lambda_p) + H[q(\lambda_1, \dots, \lambda_p)] \\ &= -\frac{p}{2} \log 2\pi - \frac{1}{2} \log v_1 \cdots v_p - \frac{1}{2} \sum_{k=1}^p \frac{\mathbb{E}(\lambda_k - m_k)^2}{v_k} \\ &\quad + \frac{p}{2} (1 + \log 2\pi) + \frac{1}{2} \log \tilde{v}_1 \cdots \tilde{v}_p \\ &= \text{const.} + \frac{1}{2} \sum_{k=1}^p \log \tilde{v}_k - \frac{1}{2} \sum_{k=1}^p \frac{\tilde{v}_k + \tilde{\lambda}_k^2 - 2\tilde{\lambda}_k m_k}{v_k} \end{aligned}$$

Notes:

1. The priors on the  $\lambda_k$ 's are  $\text{N}(m_k, v_k)$ , and  $q(\lambda_k)$  is the density of  $\text{N}(\tilde{\lambda}_k, v_{\lambda_k})$ .
2. When using improper priors  $\lambda_k \propto \text{const.}$ , then we need only consider the middle term involving the sums of  $\log \tilde{v}_{\lambda_k}$ .



### 5.13.4 Terms involving distribution of $\alpha$

For the intercepts, consider only

$$\begin{aligned} \mathbb{E} \log p(\alpha) + H[q(\alpha)] &= \text{const.} - \frac{1}{2} \mathbb{E} \sum_{j=1}^m \frac{(\alpha_j - a_j)^2}{A_j} + \frac{1}{2} \log \tilde{v}_{\alpha_1} \cdots \tilde{v}_{\alpha_m} \\ &= \text{const.} + \frac{1}{2} \sum_{j=1}^m \log \tilde{v}_{\alpha_j} - \frac{1}{2} \sum_{j=1}^m \frac{v_{\alpha_j} + \tilde{\alpha}_j^2 - 2a_j \tilde{\alpha}_j}{A_j} \end{aligned}$$

Notes:

1.  $p(\alpha)$  is  $\prod_{j=1}^m \phi(\alpha_j | a_j, A_j)$ , and  $q(\alpha) \prod_{j=1}^m \phi(\alpha_j | \tilde{\alpha}_j, \tilde{v}_{\alpha_j})$ .

### 5.13.5 ELBO summarised

In the example section of Chapter 5, we considered only 1) the independent I-probit model; 2) fixed  $\Sigma = \mathbf{I}_m$ ; 3) only RKHS scale parameters to estimate; and 4) and improper priors on the hyperparameters. In such situations, the ELBO expression is simply

$$\mathcal{L} = \text{const.} + \sum_{i=1}^n \log C_i - \frac{1}{2} \sum_{j=1}^m \text{tr}(\tilde{\mathbf{V}}_{w_j} + \tilde{\mathbf{w}}_{\cdot j} \tilde{\mathbf{w}}_{\cdot j}^\top) + \frac{1}{2} \sum_{k=1}^p \log \tilde{v}_k.$$

As a final remark, often times the ELBO is treated as a proxy for the (penalised) marginal likelihood of the model, in which case it must be noted that the ELBO as we had derived is correct up to a constant. We find that keeping track of the constants is slightly tedious, and hence decided not to do so. When comparing ELBOs of two or more models, the comparison is still valid as only differences between the ELBOs matter, in which case the constants would cancel out.

# Bibliography

- |  |  |
|--|--|
| albert1993bayesian                       | Albert, James H and Siddhartha Chib (1993). “Bayesian analysis of binary and polychotomous response data”. In: <i>Journal of the American statistical Association</i> 88.422, pp. 669–679.   |
| beal2003                                 | Beal, M. J. and Z. Ghahramani (2003). “The variational Bayesian EM algorithm for incomplete data: With application to scoring graphical model structures”. In: <i>Bayesian Statistics 7</i> . Proceedings of the Seventh Valencia International Meeting. Ed. by José M. Bernardo, A. Philip Dawid, James O. Berger, Mike West, David Heckerman, M.J. Bayarri, and Adrian F.M. Smith. Oxford: Oxford University Press, pp. 453–464. |
| bishop2006pattern<br>blei2017variational | Bishop, Christopher (2006). <i>Pattern Recognition and Machine Learning</i> . Springer-Verlag.   |
|  | Blei, David M, Alp Kucukelbir, and Jon D McAuliffe (2017). “Variational inference: A review for statisticians”. In: <i>Journal of the American Statistical Association</i> just-accepted.  |
| chopin2011fast                           | Chopin, Nicolas (2011). “Fast simulation of truncated Gaussian distributions”. In: <i>Statistics and Computing</i> 21.2, pp. 275–288.  |
| damien2001sampling                       | Damien, Paul and Stephen G Walker (2001). “Sampling truncated normal, beta, and gamma densities”. In: <i>Journal of Computational and Graphical Statistics</i> 10.2, pp. 206–215.  |
| geweke1991efficient                      | Geweke, John (1991). <i>Efficient simulation from the multivariate normal and student-t distributions subject to linear constraints and the evaluation of constraint probabilities</i> .   |
| geweke1994alternative                    | Geweke, John, Michael Keane, and David Runkle (1994). “Alternative computational approaches to inference in the multinomial probit model”. In: <i>The review of economics and statistics</i> , pp. 609–632.  |

girolami2006variational	Girolami, Mark and Simon Rogers (2006). “Variational Bayesian Multinomial Probit Regression with Gaussian Process Priors”. In: <i>Neural Computation</i> 18.8, pp. 1790–1817.
groves1969note	Groves, Theodore and Thomas Rothenberg (1969). “A note on the expected value of an inverse matrix”. In: <i>Biometrika</i> 56.3, pp. 690–691.
hajivassiliou1996simulation	Hajivassiliou, Vassilis, Daniel McFadden, and Paul Ruud (1996). “Simulation of multivariate normal rectangle probabilities and their derivatives theoretical and computational results”. In: <i>Journal of econometrics</i> 72.1-2, pp. 85–134.
hastie1986	Hastie, Trevor and Robert Tibshirani (Aug. 1986). “Generalized Additive Models”. In: <i>Statist. Sci.</i> 1.3, pp. 297–310. DOI: <a href="https://doi.org/10.1214/ss/1177013604">10.1214/ss/1177013604</a> . URL: <a href="https://doi.org/10.1214/ss/1177013604">https://doi.org/10.1214/ss/1177013604</a> .
itzykson1991statistica	Itzykson, Claude and Jean Michel Drouffe (1991). <i>Statistical Field Theory: Volume 2, Strong Coupling, Monte Carlo Methods, Conformal Field Theory and Random Systems</i> . Cambridge University Press.
jamil2017	Jamil, Haziq and Wicher Bergsma (2017). “iprior: An R Package for Regression Modelling using I-priors”. In: <i>Manuscript in submission</i> .
kass1995bayes	Kass, Robert E and Adrian E Raftery (1995). “Bayes factors”. In: <i>Journal of the american statistical association</i> 90.430, pp. 773–795.
keane1994computationally	Keane, Michael P (1994). “A computationally practical simulation estimator for panel data”. In: <i>Econometrica: Journal of the Econometric Society</i> , pp. 95–116.
Keane1992	Keane, Michael P. (1992). “A Note on Identification in the Multinomial Probit Model”. In: <i>Journal of Business &amp; Economic Statistics</i> 10.2, pp. 193–200. ISSN: 0735-0015. DOI: <a href="https://doi.org/10.2307/1391677">10.2307/1391677</a> . URL: <a href="http://www.jstor.org/stable/1391677">http://www.jstor.org/stable/1391677</a> <a href="http://www.jstor.org/stable/pdfplus/1391677.pdf?acceptTC=true">http://www.jstor.org/stable/pdfplus/1391677.pdf?acceptTC=true</a> .
marsaglia2000ziggurat	Marsaglia, George and Wai Wan Tsang (2000). “The ziggurat method for generating random variables”. In: <i>Journal of statistical software</i> 5.8, pp. 1–7.
mccullagh1989	McCullagh, P. and John A. Nelder (1989). <i>Generalized Linear Models</i> . 2nd. Chapman & Hall/CRC Press.
meng1997algorithm	Meng, Xiao-Li and David Van Dyk (1997). “The EM Algorithm—an Old Folk-song Sung to a Fast New Tune”. In: <i>Journal of the Royal Statistical Society: Series B (Statistical Methodology)</i> 59.3, pp. 511–567.

minka2001expectation	Minka, Thomas P (2001). “Expectation propagation for approximate Bayesian inference”. In: <i>Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence</i> . Morgan Kaufmann Publishers Inc., pp. 362–369.
neal1999	Neal, Radford M. (1999). “Regression and Classification using Gaussian Process Priors”. In: <i>Bayesian Statistics</i> . Ed. by J M Bernardo, J O Berger, A P Dawid, and A F M Smith. Vol. 6. Oxford University Press. (with discussion), pp. 475–501.
petersen2008matrix	Petersen, Kaare Brandt and Michael Syskind Pedersen (2008). “The matrix cookbook”. In: <i>Technical University of Denmark</i> 7.15, p. 510.
rasmussen2006gaussian	Rasmussen, Carl Edward and Christopher K I Williams (2006). <i>Gaussian Processes for Machine Learning</i> . The MIT Press.
robert1995simulation	Robert, Christian P (1995). “Simulation of truncated normal variables”. In: <i>Statistics and computing</i> 5.2, pp. 121–125.
scholkopf2002learning	Schölkopf, Bernhard and Alexander J Smola (2002). <i>Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond</i> . MIT Press.
train2009discrete	Train, Kenneth E (2009). <i>Discrete choice methods with simulation</i> . Cambridge university press.
zhang2013kronecker	Zhang, Huamin and Feng Ding (2013). “On the Kronecker products and their applications”. In: <i>Journal of Applied Mathematics</i> 2013.