# To-do list

# Contents

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 $y_1, \ldots, y_n$, where each $y_i$ takes on exactly one of the values $\{1, \ldots, m\}$ 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

$$y_i \sim \mathrm{Cat}(p_{i1}, \ldots, p_{im}),$$

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

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

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

{eq:catdist}

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

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

where $g : [0,1] \to \mathbb{R}^m$ 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 $j \in \{1, \ldots, m\}$ 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 $\sigma_{jk} := \mathrm{Cov}[\epsilon_{ij}, \epsilon_{ik}]$, for each $j, k \in \{1, \ldots, m\}$ 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

---

[1] $[A]$ 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

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

## 5.3 Identifiability and IIA

## 5.4 Estimation

## 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 $\theta = \{\boldsymbol{\alpha} := (\alpha_1, \dots, \alpha_m), \eta, \boldsymbol{\Psi}\}$. Begin by assuming some prior distribution on the parameters $p(\theta) = p(\boldsymbol{\alpha})p(\eta)p(\boldsymbol{\Psi})$. 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 the PDFs belong to the exponential family of distributions with known hyperparameters. The exponential family requirement greatly eases the complexity of deriving the variational algorithm later on[2].

Recall that $\mathbf{y}^*|\mathbf{w} \sim \mathrm{MN}_{n,m}(\mathbf{1}_n\boldsymbol{\alpha}^\top + \mathbf{H}_\eta\mathbf{w}, \mathbf{I}_n, \boldsymbol{\Psi}^{-1})$ and $\mathbf{w} \sim \mathrm{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \boldsymbol{\Psi})$. The required posterior distribution is then $p(\mathbf{y}^*, \mathbf{w}, \theta|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{y}^*)p(\mathbf{y}^*|\mathbf{w}, \theta)p(\mathbf{w}|\theta)p(\theta)$. This is approximated by a mean-field distribution of the form $q(\mathbf{y}^*, \mathbf{w}, \theta) \equiv q(\mathbf{y}^*)q(\mathbf{w})q(\theta)$, and also $q(\theta) = q(\boldsymbol{\alpha})q(\eta)q(\boldsymbol{\Psi})$. 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

$$\log \tilde{q}(\xi) = \mathrm{E}_{-\xi}[\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta)] + \text{const.} \tag{5.2}$$

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.

---

[2] Of interest, one may even opt to assign improper priors on $\theta$ and the algorithm would still work. This is akin to obtaining empirical Bayes estimate of the $\theta$ if seen from an EM algorithm standpoint.
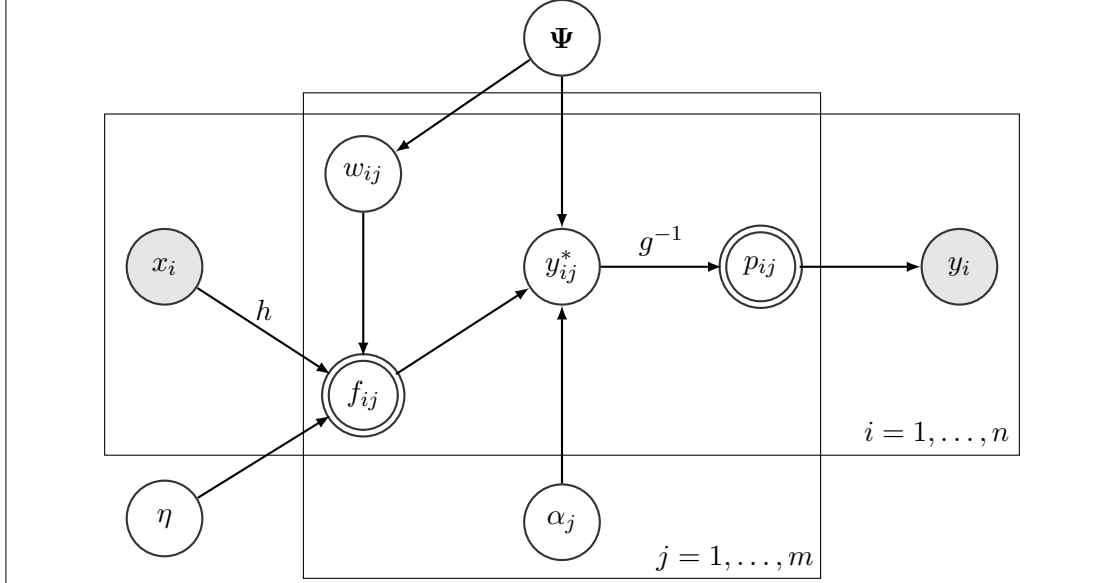
Figure 5.1: A DAG of the I-probit model. Observed nodes are shaded, while double-lined nodes represents calculable quantities.

In practice, rather than an explicit calculation of the normalising constant, one simply needs to inspect (5.2) 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,

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

Then, from (5.2),

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

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.

We now present the mean-field variational distributions $\tilde{q}$. 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 $\mathrm{E}_{\mathbf{w} \sim \tilde{q}}[\mathbf{w}]$, the expected value of $\mathbf{w}$ under the pdf $\tilde{q}(\mathbf{w})$. 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 of $\mathbf{y}^*$ are independent can be exploited. Write $\mathbf{y}_i^* = (y_{i1}^*, \ldots, y_{im}^*)^\top$. Then $\mathbf{y}_i^* | \theta, x_i \sim \mathrm{N}_m(\boldsymbol{\alpha} + \mathbf{f}(x_i), \boldsymbol{\Psi}^{-1})$, and we have the induced factorisation of the distribution $q(\mathbf{y}^*) = \prod_{i=1}^n q(\mathbf{y}_i^*)$, where each $q(\mathbf{y}_i^*)$ is the density of a *conically truncated multivariate normal disribution*. That is, for each $i = 1, \ldots, n$ and noting the observed values $y_i = j \in \{1, \ldots, m\}$, the $\mathbf{y}_i^*$'s are distributed according to

$$\mathbf{y}_i^* \overset{\text{iid}}{\sim} \begin{cases} \mathrm{N}_m(\tilde{\boldsymbol{\alpha}} + \tilde{\mathbf{f}}(x_i), \tilde{\boldsymbol{\Sigma}}) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise.} \end{cases} \tag{5.3}$$

The required expectations $\mathrm{E}\,\mathbf{y}_i^* = \mathrm{E}(y_{i1}^*, \ldots, y_{im}^*)^\top$ are tricky to compute. One strategy might be Monte Carlo integration: using samples from $\mathrm{N}_m(\tilde{\alpha} + \tilde{\mathbf{f}}(x_i), \tilde{\boldsymbol{\Psi}}^{-1})$, zero out those that do not satisfy the condition $y_{ij}^* > y_{ik}^*, \forall k \neq j$, then take the sample average. If the independent I-probit model is considered, where $\boldsymbol{\Psi} = \mathrm{diag}(\sigma_1^{-2}, \ldots, \sigma_m^{-2})$, the expected value can be considered component-wise, and each component of this expectation is given by

$$\tilde{y}_{ik}^* = \begin{cases} \tilde{\alpha}_k + \tilde{f}_{ik} - \tilde{\sigma}_k C_i^{-1} \int \phi_{ik}(z) \prod_{l \neq k, j} \Phi_{il}(z) \phi(z) \, \mathrm{d}z & \text{if } k \neq y_i \\ \tilde{\alpha}_{y_i} + \tilde{f}_{iy_i} - \tilde{\sigma}_{y_i} \sum_{k \neq y_i} \left( \tilde{y}_{ik}^* - \tilde{f}_{ik} \right) & \text{if } k = y_i \end{cases} \tag{5.4}$$

with

$$\phi_{ik}(Z) = \phi \left( \frac{\tilde{\sigma}_{y_i}}{\tilde{\sigma}_k} Z + \frac{\tilde{\alpha}_{y_i} + \tilde{f}_{iy_i} - \tilde{\alpha}_k - \tilde{f}_{ik}}{\tilde{\sigma}_k} \right)$$

$$\Phi_{ik}(Z) = \Phi \left( \frac{\tilde{\sigma}_{y_i}}{\tilde{\sigma}_k} Z + \frac{\tilde{\alpha}_{y_i} + \tilde{f}_{iy_i} - \tilde{\alpha}_k - \tilde{f}_{ik}}{\tilde{\sigma}_k} \right)$$

$$C_i = \int \prod_{l \neq j} \Phi_{il}(z) \phi(z) \, \mathrm{d}z$$

and $Z \sim \mathrm{N}(0,1)$ 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 $\mathrm{vec}\,\mathbf{y}^*|\,\mathrm{vec}\,\mathbf{w}$ and $\mathrm{vec}\,\mathbf{w}$ 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 $\tilde{q}$ for $\mathrm{vec}\,\mathbf{w} \in \mathbb{R}^{nm}$ is also normal with mean and precision given by

$$\mathrm{vec}\,\tilde{\mathbf{w}} = \tilde{\mathbf{V}}_w(\tilde{\mathbf{\Psi}} \otimes \tilde{\mathbf{H}}_\eta)\,\mathrm{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n\tilde{\boldsymbol{\alpha}}^\top) \quad \text{and} \quad \tilde{\mathbf{V}}_w^{-1} = (\tilde{\mathbf{\Psi}} \otimes \tilde{\mathbf{H}}_\eta^2) + (\tilde{\mathbf{\Psi}}^{-1} \otimes \mathbf{I}_n). \quad (5.5)$$

<div align="right">{eq:varipostw}</div>

We note the similarity between (5.5) 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 $O(n^3m^3)$ time. By exploiting the Kronecker product structure in $\tilde{\mathbf{V}}_w^{-1}$, we are able to efficiently compute the required inverse in roughly $O(n^3m)$ time—see the appendix for details. Equivalently, we can express the distribution for $\mathbf{w} \sim \tilde{q}$ as a matrix normal distribution

$$\mathrm{MN}_{nm}\left( \overbrace{\tilde{\mathbf{H}}_\eta^{-1}(\tilde{\mathbf{y}}^* - \mathbf{1}_n\tilde{\boldsymbol{\alpha}}^\top)\tilde{\mathbf{\Psi}}^2}^{\tilde{\mathbf{w}}}, \tilde{\mathbf{H}}_\eta^{-2}, \tilde{\mathbf{\Psi}} \right). \quad (5.6)$$

<div align="right">{eq:varipostw2}</div>

If the independent I-probit model is assumed, i.e. $\tilde{\mathbf{\Psi}} = \mathrm{diag}(\tilde{\sigma}_1^{-2},\ldots,\tilde{\sigma}_m^{-2})$, then the posterior covariance matrix $\tilde{\mathbf{V}}_w$ has a simpler structure. This means that the random matrix $\mathbf{w}$ will have columns which are independent of each other. By writing $\mathbf{w}_j = (w_{1j},\ldots,w_{nj})^\top \in \mathbb{R}^n$, $j = 1,\ldots,m$, to denote the column vectors of $\mathbf{w}$ and with a slight abuse of notation, we have that

$$\mathrm{N}_{nm}(\mathrm{vec}\,\mathbf{w}|\,\mathrm{vec}\,\tilde{\mathbf{w}},\tilde{\mathbf{V}}_w) = \prod_{j=1}^{m} \mathrm{N}_n(\mathbf{w}_j|\tilde{\mathbf{w}}_j,\tilde{\mathbf{V}}_{w_j}),$$

where

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

### 5.5.3   RKHS parameters $\eta$

The posterior density $\tilde{q}$ involving the RKHS parameters is of the form

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

where $p(\eta)$ is an appropriate prior distribution for $\eta$. Generally, samples $\eta^{(1)}, \ldots, \eta^{(T)}$ from $\tilde{q}(\eta)$ may be obtained using a Metropolis algorithm, and quantities such as $\tilde{\mathbf{H}}_\eta = \operatorname{E}_q[\mathbf{H}_\eta]$ may be approximated using $\frac{1}{T} \sum_{t=1}^{T} \mathbf{H}_{\eta^{(t)}}$.

However, when only RKHS scale parameters are involved, then the distribution $\tilde{q}$ 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 $\eta = \{\lambda_1, \ldots, \lambda_p\}$ need be estimated, and for each $k = 1, \ldots, p$, 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 $q(\eta)$, and that is $q(\eta) = \prod_{k=1}^{p} p(\lambda_k)$. Then, by using independent and identical normal priors for $\lambda_k$, say $\lambda_k \sim \mathrm{N}(0, v_\lambda)$, each $\tilde{q}(\lambda_k)$ density is normal with mean and variance

Write down the mean and variance for lambda

### 5.5.4   Error precision $\boldsymbol{\Psi}$

A small reparameterisation of the I-prior random effects is necessary to achieve conjugacy for the $\boldsymbol{\Psi}$ parameter. Let $\mathbf{u} \in \mathbb{R}^{n \times m}$ be a matrix defined by $\boldsymbol{\Psi}^{-1}\mathbf{w}$. Then $\mathbf{u} \sim \mathrm{MN}_{n,m}(\mathbf{0}, \mathbf{I}_n, \boldsymbol{\Psi}^{-1})$ a priori. From (5.6), the optimal variational distribution for $\mathbf{u}$ would be $\mathrm{MN}_{n,m}(\tilde{\mathbf{w}}\tilde{\boldsymbol{\Psi}}^{-1}, \tilde{\mathbf{H}}_\eta^2, \tilde{\boldsymbol{\Psi}}^{-1})$. With a Wishart prior on the precision matrix $\boldsymbol{\Psi} \sim \mathrm{Wis}_m(\mathbf{G}, g)$, where $g \geq m$, the optimal variational density for $\boldsymbol{\Psi}$ is found to satisfy

$$\log \tilde{q}(\boldsymbol{\Psi}) = \text{const.} - \frac{1}{2} \sum_{i=1}^{n} \operatorname{tr} \left( (\mathbf{G}_1 + \mathbf{G}_2 + \mathbf{G}) \boldsymbol{\Psi} \right) + \frac{g - m - 1}{2} \log|\boldsymbol{\Psi}|$$

which is recognised as the log density of a Wishart distribution with scale matrix $\mathbf{G}_1 + \mathbf{G}_2 + \mathbf{G}$ and $g$ degrees of freedom, where

$$\mathbf{G}_1 = \mathrm{E}_{\mathcal{Z}\setminus\{\boldsymbol{\Psi}\}\sim q}\left[\sum_{i=1}^{n}(\mathbf{y}_i^* - \boldsymbol{\alpha} - \mathbf{f}(x_i))(\mathbf{y}_i^* - \boldsymbol{\alpha} - \mathbf{f}(x_i))^\top\right]$$

$$\mathbf{G}_2 = \sum_{i=1}^{n}\mathrm{E}_{\mathbf{u}\sim q}\left[\mathbf{u}_i\mathbf{u}_i^\top\right].$$

(5.7)

The challenge here is that it involves the second posterior moment of the conically truncated multivariate normal distribution for $\mathbf{y}^*$, which may be obtained by sampling or numerical integration as described earlier.

If the independent I-probit model is considered, then $\boldsymbol{\Sigma} = \mathrm{diag}(\sigma_1^2,\ldots,\sigma_m^2)$, class independence holds so we can use independent inverse gamma distributions as a prior for $\boldsymbol{\Sigma}$, i.e. $p(\boldsymbol{\Sigma}) = \prod_{j=1}^{m}p(\sigma_j^2)$, where each $p(\sigma_j) \equiv \Gamma^{-1}(r,s)$. The posterior for $\boldsymbol{\Sigma}$ will also be of a similar factorised form , namely $\tilde{q}(\boldsymbol{\Sigma}) = \prod_{j=1}^{m}\tilde{q}(\sigma_j^2)$, where $\tilde{q}(\sigma_j^2)$ is the PDF of an inverse gamma distribution with shape and scale parameters $\tilde{r} = 2n + r - 1$ and $\tilde{s} = \frac{1}{2}\|\tilde{\mathbf{y}}_j^* - \tilde{\alpha}_j - \tilde{\mathbf{f}}_j\|^2 + \frac{1}{2}\|\tilde{\mathbf{u}}_j\|^2 + s$ respectively.

Finally, the posterior distribution for the intercepts follow a normal distribution should the prior specified on the intercepts also be a normal distribution, e.g. $\boldsymbol{\alpha} \sim \mathrm{N}_m(\mathbf{0}, \mathbf{A})$. The posterior mean and variance for the intercepts are given by

$$\tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{V}}_\alpha\tilde{\boldsymbol{\Sigma}}^{-1}\big(\tilde{\mathbf{y}}_i^* - \tilde{\mathbf{f}}(x_i)\big) \quad\text{and}\quad \tilde{\mathbf{V}}_\alpha = \big(n\tilde{\boldsymbol{\Sigma}}^{-1} + \mathbf{A}^{-1}\big)^{-1}.$$

Note that the evaluation of each of the component of the posterior depends on some of the components itself, and so this circular dependence is dealt with by using some arbitrary starting values and after which an iterative updating scheme of the components ensues. The updating scheme is performed until a maximum number of iterations is reached, or ideally until some of convergence criterion is met. In variational inference, the *variational lower bound* is typically used to asses convergence. The lower bound is given by

$$\mathcal{L} = \int q(\mathbf{y}^*, \mathbf{w}, \theta)\log\left[\frac{p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta)}{q(\mathbf{y}^*, \mathbf{w}, \theta)}\right]\mathrm{d}\mathbf{y}^*\mathrm{d}\mathbf{w}\mathrm{d}\theta$$

$$= \mathrm{E}[\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \theta)] - \mathrm{E}[\log q(\mathbf{y}^*, \mathbf{w}, \theta)].$$

These are calculable once the posterior distributions $\tilde{q}$ are known—the first term is the expectation of the logarithm of the joint density, whereas the second term factorises into the entropy of its individual components. Similar to the EM algorithm, this quantity is expected to increase with every iteration.

> 4. Proof?

## 5.6    Post-estimation

## 5.7    Examples

## 5.8    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.

## 5.9 Miscellanea

### 5.9.1 A note on computing the multivariate normal integral

`misc:mnint`

How is this calculated? Simulation usually, but also quadrature methods not too bad if $m$ not too large. Stata sheet useful? Talk about if iid errors.

Much research has been devoted into developing efficient computational methods for computing these integral, and MCMC methods seem to be the tool of choice in Bayesian analysis (R. McCulloch and Rossi, 1994; Nobile, 1998; R. E. McCulloch et al., 2000). Things get more tractable if $\boldsymbol{\Sigma}$ is assumed to be diagonal (which corresponds to abandoning the independence of irrelevant alternatives assumption) and much more so if we assume that $\boldsymbol{\Sigma} = \mathbf{I}_m$. The latter yields the *normalised I-probit model*, and a discussion of the merits of this model is given later.

6. can use Hamiltonian Monte Carlo?

### 5.9.2 Similarity of EM algorithm and variational Bayes

### 5.9.3 Conically truncated multivariate normal distributions

Crucial to the probit model, the properties of conically truncated multivariate normal distributions are worth investigating.

`definition: conically-truncated-normal`

**Definition 5.1** (Conically-truncated multivariate normal distribution)**.** Let $\mathbf{X} = (X_1, \ldots, X_d)$ be a $d$-dimensional random variable with pdf defined as

$$p(\mathbf{x}) = \begin{cases} \prod_{i=1}^{d} \mathrm{N}(\mu_i, \sigma_i) & \text{if } X_j > X_i, \forall i \neq j \\ 0 & \text{otherwise} \end{cases}$$

for some $j \in \{1, \ldots, d\}$. We denote the distribution of $\mathbf{X}$ by $\mathrm{N}^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_d)$ and $\boldsymbol{\Sigma} = \mathrm{diag}(\sigma_1^2, \ldots, \sigma_d^2)$. The pdf of $\mathbf{X}$ has support on the set $\{\mathbb{R}^d \mid x_j > x_i, \forall i \neq j\}$ and the following functional form:

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

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

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

where $Z \sim \mathrm{N}(0, 1)$. In the case where all variances are unity, the pdf of $\mathbf{X} \sim \mathrm{N}^{(j)}(\boldsymbol{\mu}, \mathbf{I}_d)$ is

$$p(\mathbf{x}) = \left\{(2\pi)^{d/2} \mathrm{E}_Z\left[\prod_{\substack{i=1 \\ i \neq j}}^{d} \Phi\left(Z + \mu_j - \mu_i\right)\right]\right\}^{-1} \exp\left[-\frac{1}{2} \sum_{i=1}^{d} (x_i - \mu_i)^2\right].$$

*Proof.* A derivation of the functional form for the pdf of $X \sim \mathrm{N}^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is given. Using the fact that $\int p(x)\mathrm{d}x = 1$, and that

$$
\int \mathbb{1}[x_i < x_j, \forall i \neq j] \prod_{i=1}^{d} \mathrm{N}(\mu_i, \sigma_i^2)\mathrm{d}x_1 \cdots \mathrm{d}x_d
$$

$$
= \int \mathbb{1}[x_i < x_j, \forall i \neq j] \prod_{i=1}^{d} \left[ \frac{1}{\sigma_i} \phi \left( \frac{x_i - \mu_i}{\sigma} \right) \right] \mathrm{d}x_1 \cdots \mathrm{d}x_d
$$

$$
= \int \mathbb{1}[x_i < x_j, \forall i \neq j] \frac{1}{\sigma_j} \phi \left( \frac{x_j - \mu_j}{\sigma_j} \right) \prod_{\substack{i=1 \\ i \neq j}}^{d} \left[ \frac{1}{\sigma_i} \phi \left( \frac{x_i - \mu_i}{\sigma_i} \right) \right] \mathrm{d}x_1 \cdots \mathrm{d}x_d
$$

$$
= \int \prod_{\substack{i=1 \\ i \neq j}}^{d} \Phi \left( \frac{x_j - \mu_i}{\sigma_i} \right) \frac{1}{\sigma_j} \phi \left( \frac{x_j - \mu_j}{\sigma_j} \right) \mathrm{d}x_j
$$

$$
= \int \prod_{\substack{i=1 \\ i \neq j}}^{d} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_i}{\sigma_i} \right) \phi(z_j)\mathrm{d}z_j
$$

$$
\text{(by using the standardisation } z_j = (x_j - \mu_j)/\sigma_j)
$$

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

the proof follows directly. □

**Lemma 5.1.** *Let* $X \sim \mathrm{N}^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ *with pdf* $p(\mathbf{x})$ *as defined in Definition 5.1. Then*

*(i) The expectation* $\mathrm{E}[\mathbf{X}] = \big( \mathrm{E}[X_1], \ldots, \mathrm{E}[X_d] \big)$ *is given by*

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

*(ii) The differential entropy* $\mathcal{H}(p)$ *is given by*

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

where $C = \mathrm{E}\left[\prod_{i \neq j} \Phi_i\right]$, and we had defined

$$\phi_i = \phi_i(Z) = \phi\left(\frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i}\right)$$

$$\Phi_i = \Phi_i(Z) = \Phi\left(\frac{\sigma_j Z + \mu_j - \mu_i}{\sigma_i}\right)$$

with $Z \sim \mathrm{N}(0,1)$, and $\phi(\cdot)$ and $\Phi(\cdot)$ the pdf and cdf of $Z$ respectively.

# Appendix

Fact: $X \sim \mathrm{N}(a, A)$ and $Y \sim \mathrm{N}(b, B)$, then

$$p(x)p(y) \propto \mathrm{N}(c, C)$$

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

## 5.9.4  Proof of Lemma

*Proof.*    (i) Due to the independence structure in the pdf of $\mathbf{X}$, it is easy to consider the expectations of each of the components separately and marginalising out the

rest of the components. For $i \neq j$, we have

$$
\mathrm{E}[x_i] = C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_i \prod_{k=1}^{d} \frac{1}{\sigma_k} \phi \left( \frac{x_k - \mu_k}{\sigma_k} \right) \mathrm{d}x_1 \cdots \mathrm{d}x_d
$$

$$
= C^{-1} \iint \mathbb{1}[x_i < x_j] \frac{x_i}{\sigma_i} \phi \left( \frac{x_i - \mu_i}{\sigma_i} \right) \prod_{k \neq i,j} \Phi \left( \frac{x_j - \mu_k}{\sigma_k} \right) \frac{1}{\sigma_j} \phi \left( \frac{x_j - \mu_j}{\sigma_j} \right) \mathrm{d}x_i \mathrm{d}x_j
$$

$$
= C^{-1} \iint \mathbb{1}[\sigma_i z_i + \mu_i < \sigma_j z_j + \mu_j](\sigma_i z_i + \mu_i)\phi(z_i) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_i\mathrm{d}z_j
$$

$$
= \mu_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i]\phi(z_i) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_i\mathrm{d}z_j
$$

$$
+ \sigma_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i]z_i\phi(z_i) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_i\mathrm{d}z_j
$$

$$
= \mu_i C^{-1} \overbrace{\int \prod_{k \neq j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_j}^{C}
$$

$$
+ \sigma_i C^{-1} \int \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i]z_i\phi(z_i) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_i\mathrm{d}z_j
$$

The integral involving $z_i$ in the second part of the sum is recognised as the (unnormalised) expectation of the lower-tail of a univariate standard normal distribution truncated at $\tau_{ij} = (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i$. That is,

$$
\mathrm{E}[Z_i | Z_i < \tau_{ij}] = \left[ \Phi(\tau_{ij}) \right]^{-1} \int \mathbb{1}[z_i < \tau_{ij}]z_i\phi(z_i)\mathrm{d}z_i = -\frac{\phi(\tau_{ij})}{\Phi(\tau_{ij})}
$$

Plugging this expression back into the derivation of this expectation, we get

$$
\mathrm{E}[X_i] = \mu_i - \sigma_i C^{-1} \int \phi \left( \frac{\sigma_j z_j + \mu_j - \mu_i}{\sigma_i} \right) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k} \right) \phi(z_j)\mathrm{d}z_j
$$

$$
= \mu_i - \sigma_i C^{-1} \mathrm{E} \left[ \phi \left( \frac{\sigma_j Z_j + \mu_j - \mu_i}{\sigma_i} \right) \prod_{k \neq i,j} \Phi \left( \frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k} \right) \right].
$$

The expectation for the $j$th component is

$$
\mathrm{E}[X_j] = C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_j \prod_{k=1}^{d} \frac{1}{\sigma_k} \phi\left(\frac{x_k - \mu_k}{\sigma_k}\right) \mathrm{d}x_1 \cdots \mathrm{d}x_d
$$

$$
= C^{-1} \int x_j \prod_{k \neq j} \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) \cdot \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) \mathrm{d}x_j
$$

$$
= C^{-1} \int (\sigma_j z_j + \mu_j) \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot \phi(z_j) \mathrm{d}z_j
$$

$$
= \mu_j C^{-1} \overbrace{\int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot \phi(z_j) \mathrm{d}z_j}^{C}
$$

$$
+ \sigma_j C^{-1} \int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot z_j \phi(z_j) \mathrm{d}z_j
$$

$$
= \mu_j + \sigma_j C^{-1} \mathrm{E}\left[ Z_j \prod_{k \neq j} \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right) \right]
$$

$$
= \mu_j + \sigma_j \sum_{\substack{i=1 \\ i \neq j}}^{d} \sigma_i C^{-1} \mathrm{E}\left[ \phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i,j}}^{d} \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right) \right]
$$

$$
= \mu_j - \sigma_j \sum_{i \neq j} \left( \mathrm{E}[X_i] - \mu_i \right)
$$

where we have made use of Lemma 5.2 in the second last step of the above.

(ii) The differential entropy is given by

$$
\mathcal{H}(p) = - \int p(\mathbf{x}) \log p(\mathbf{x}) \mathrm{d}\mathbf{x} = - \mathrm{E}\left[\log p(\mathbf{x})\right]
$$

$$
= - \mathrm{E}\left[ -\log C - \frac{d}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^{d} \log \sigma_i^2 - \frac{1}{2} \sum_{i=1}^{d} \left(\frac{x_i - \mu_i}{\sigma_i}\right)^2 \right]
$$

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

□

lem:EZgZ

**Lemma 5.2.** *Let $Z \sim \mathrm{N}(0,1)$. Then for all $m \in \{\mathbb{N} \,|\, m > 1\}$ and $(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+$,*

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

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

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

We adopt the following notation for convenience:

$$\phi_i = \phi(\sigma_i z + \mu_i)$$
$$\Phi_i = \Phi(\sigma_i z + \mu_i)$$

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

$$g_2(z) = \Phi_2$$

$$\Rightarrow g_2'(z) = \sigma_2 \phi_2 = \sum_{\substack{i=1 \\ i \neq 1}}^{2} \left[\sigma_i \phi_i \sum_{\substack{k=1 \\ k \neq 1,2}}^{2} \Phi_k\right].$$

Now assume that the inductive hypothesis holds for some $m \in \{\mathbb{N} \,|\, m > 1\}$. That is, the derivative of

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

which is

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

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

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

is found to be

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

Thus, by induction and linearity of expectations, the proof is complete. $\qquad\square$

## 5.10  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

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

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

$$
\log \tilde{q}(\xi) = \mathrm{E}_{-\xi}[\log p(\mathbf{y}, \mathcal{Z})] + \text{const.} \tag{5.2}
$$

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 $i \in \{1, \ldots, n\}$, given the corresponding latent propensities $\mathbf{y}_i^* = (y_{i1}^*, \ldots, y_{im}^*)$, 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.

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

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

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

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

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

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 $\mathbf{w} \sim \mathrm{N}_{n,m}(\mathbf{0}, \mathbf{I}_n, \boldsymbol{\Psi})$ with pdf

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

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

pdf is

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

An improper prior $p(\eta) \propto$ const. can be used as well, and this is the same as letting $m_k \to 0$ and $v_k \to 0$. 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 $p(\gamma) = \mathbb{1}\left(\gamma \in (0,1)\right)$ and $l, c \sim \Gamma(a, b)$.

- $p(\boldsymbol{\Psi})$. For the precision matrix, a Wishart prior with scale matrix $\mathbf{G}^{-1}$ and $g$ degrees of freedom, denoted $\boldsymbol{\Psi} \sim \mathrm{Wis}_m(\mathbf{G}^{-1}, g)$, is convenient. It has pdf

$$p(\boldsymbol{\Psi}) = \exp\left[\mathrm{const.} + \frac{g - m - 1}{2}\log|\boldsymbol{\Psi}| - \frac{1}{2}\mathrm{tr}(\mathbf{G}\boldsymbol{\Psi})\right].$$

For the independent I-probit model, $\boldsymbol{\Psi} = \mathrm{diag}(\sigma_1^{-2}, \ldots, \sigma_m^{-2})$, and we choose independent Gamma distributions for each precision $\sigma_j^{-2} \sim \Gamma(r_j, s_j)$, where $r_j$ and $s_j$ are the shape and scale parameters. Then,

$$p(\boldsymbol{\Psi}) = \prod_{j=1}^{m} \exp\left[\mathrm{const.} + (r_j - 1)\log \sigma_j^{-2} - \frac{\sigma_j^{-2}}{s_j}\right].$$

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

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

*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.

### 5.10.1 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

$$\log \tilde{q}(\mathbf{y}_i^*) = \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*]\, \mathrm{E}_{\mathcal{Z} \setminus \{\mathbf{y}^*\} \sim q} \left[ -\frac{1}{2}(\mathbf{y}_i^* - \boldsymbol{\mu}_i)^\top \boldsymbol{\Psi}(\mathbf{y}_i^* - \boldsymbol{\mu}_i) \right] + \text{const.}$$

$$= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \left[ -\frac{1}{2}(\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i)^\top \tilde{\boldsymbol{\Psi}}(\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i) \right] + \text{const.} \qquad (\star)$$

$$\equiv \begin{cases} \phi(\mathbf{y}_i^* | \tilde{\boldsymbol{\mu}}_i, \tilde{\boldsymbol{\Psi}}) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise} \end{cases}$$

where $\tilde{\boldsymbol{\mu}}_i = \mathrm{E}\,\boldsymbol{\alpha} + (\mathrm{E}\,\mathbf{H}_\eta\, \mathrm{E}\,\mathbf{w})_i$, 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 <mark>Lemma X</mark> 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,

$$\mathrm{E}[(\mathbf{y}_i^* - \boldsymbol{\mu}_i)^\top \boldsymbol{\Psi}(\mathbf{y}_i^* - \boldsymbol{\mu}_i)] = \mathrm{E}[\mathbf{y}_i^{*\top} \boldsymbol{\Psi} \mathbf{y}_i^* + \boldsymbol{\mu}_i^\top \boldsymbol{\Psi} \boldsymbol{\mu}_i - 2\boldsymbol{\mu}_i^\top \boldsymbol{\Psi} \mathbf{y}_i^*]$$

$$= \mathbf{y}_i^{*\top} \boldsymbol{\Psi} \mathbf{y}_i^* - 2\,\mathrm{E}[\boldsymbol{\mu}_i^\top]\,\mathrm{E}[\boldsymbol{\Psi}]\mathbf{y}_i^* + \text{const.}$$

$$= \mathbf{y}_i^{*\top} \boldsymbol{\Psi} \mathbf{y}_i^* - 2\tilde{\boldsymbol{\mu}}_i^\top \tilde{\boldsymbol{\Psi}} \mathbf{y}_i^* + \text{const.}$$

$$= (\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i)^\top \tilde{\boldsymbol{\Psi}}(\mathbf{y}_i^* - \tilde{\boldsymbol{\mu}}_i) + \text{const.}$$

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

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

The terms involving $\mathbf{w}$ in (5.2) 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

$$\operatorname{vec} \mathbf{y}^* | \boldsymbol{\alpha}, \mathbf{w}, \eta, \boldsymbol{\Psi} \sim \mathrm{N}_{nm} \left( \operatorname{vec}(\mathbf{1}_n \boldsymbol{\alpha}^\top + \mathbf{H}_\eta \mathbf{w}), \boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n \right)$$

and

$$\operatorname{vec} \mathbf{w} | \boldsymbol{\Psi} \sim \mathrm{N}_{nm}(\mathbf{0}, \boldsymbol{\Psi} \otimes \mathbf{I}_n)$$

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

$$\log \tilde{q}(\mathbf{w}) = \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ -\frac{1}{2} (\bar{\mathbf{y}}^* - \mathbf{M} \operatorname{vec} \mathbf{w})^\top (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n)^{-1} (\bar{\mathbf{y}}^* - \mathbf{M} \operatorname{vec} \mathbf{w}) \right]$$

$$+ \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ -\frac{1}{2} (\operatorname{vec} \mathbf{w})^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n)^{-1} \operatorname{vec}(\mathbf{w}) \right] + \text{const.}$$

$$= -\frac{1}{2} \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\operatorname{vec} \mathbf{w})^\top \overbrace{\left( \mathbf{M}^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n) \mathbf{M} + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \right)}^{\mathbf{A}} \operatorname{vec}(\mathbf{w}) \right]$$

$$+ \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ \overbrace{\bar{\mathbf{y}}^{*\top} (\boldsymbol{\Psi} \otimes \mathbf{I}_n) \mathbf{M}}^{\mathbf{a}^\top} \operatorname{vec}(\mathbf{w}) \right] + \text{const.}$$

$$= -\frac{1}{2} \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\operatorname{vec} \mathbf{w} - \mathbf{A}^{-1} \mathbf{a})^\top \mathbf{A} (\operatorname{vec} \mathbf{w} - \mathbf{A}^{-1} \mathbf{a}) \right] + \text{const.}$$

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

$$\mathbf{V}_w^{-1} = \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q}[\mathbf{A}]$$

$$= \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\mathbf{I}_m \otimes \mathbf{H}_\eta)^\top (\boldsymbol{\Psi} \otimes \mathbf{I}_n)(\mathbf{I}_m \otimes \mathbf{H}_\eta) + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \right]$$

$$= \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta^2) + (\boldsymbol{\Psi}^{-1} \otimes \mathbf{I}_n) \right]$$

$$= (\tilde{\boldsymbol{\Psi}} \otimes \tilde{\mathbf{H}}_\eta^2) + (\tilde{\boldsymbol{\Psi}}^{-1} \otimes \mathbf{I}_n)$$

and making a first-order approximation $(\mathrm{E}\,\mathbf{A})^{-1} \approx \mathrm{E}[\mathbf{A}^{-1}]$[3],

$$\operatorname{vec} \tilde{\mathbf{w}} = \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q}[\mathbf{A}^{-1} \mathbf{a}]$$

$$= \tilde{\mathbf{V}}_w \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\mathbf{I}_m \otimes \mathbf{H}_\eta)(\boldsymbol{\Psi} \otimes \mathbf{I}_n) \operatorname{vec}(\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top) \right]$$

$$= \tilde{\mathbf{V}}_w \mathrm{E}_{\mathcal{Z} \backslash \{\mathbf{w}\} \sim q} \left[ (\boldsymbol{\Psi} \otimes \mathbf{H}_\eta) \operatorname{vec}(\mathbf{y}^* - \mathbf{1}_n \boldsymbol{\alpha}^\top) \right]$$

$$= \tilde{\mathbf{V}}_w (\tilde{\boldsymbol{\Psi}} \otimes \tilde{\mathbf{H}}_\eta) \operatorname{vec}(\tilde{\mathbf{y}}^* - \mathbf{1}_n \tilde{\boldsymbol{\alpha}}^\top).$$

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 Kronekcer product structure to compute the inverse efficiently. Perform an orthogonal eigendecomposition of $\mathbf{H}_\eta$ to obtain $\mathbf{H}_\eta = \mathbf{VUV}^\top$ and of $\boldsymbol{\Psi}$ to obtain $\boldsymbol{\Psi} = \mathbf{QPQ}^\top$. 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):

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

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.

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

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

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

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

---

[3]Groves and Rothenberg (1969) show that $\mathrm{E}[\mathbf{A}^{-1}] = (\mathrm{E}\,\mathbf{A})^{-1} + \mathbf{B}$, where $\mathbf{B}$ is a positive-definite matrix.

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

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

A quantity that we will be requiring time and again will be $\mathrm{tr}(\mathbf{C}\,\mathrm{E}[\mathbf{w}^\top\mathbf{D}\mathbf{w}])$, 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

$$
\begin{aligned}
\mathrm{tr}(\mathbf{C}\,\mathrm{E}[\mathbf{w}^\top\mathbf{D}\mathbf{w}]) &= \sum_{i,j=1}^{m} \mathbf{C}_{ij}\,\mathrm{E}[\mathbf{w}^\top\mathbf{D}\mathbf{w}]_{ij} \\
&= \sum_{i,j=1}^{m} \mathbf{C}_{ij}\,\mathrm{E}[\mathbf{w}_{\cdot i}^\top\mathbf{D}\mathbf{w}_{\cdot j}]. \quad (5.8)
\end{aligned}
$$

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

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

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 (5.8) leads us to

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

If $\mathbf{C} = \mathrm{diag}(c_1, \ldots, c_m)$, then

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

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

By looking at only the terms involving $\eta$ in (5.2), we deduce that $\tilde{q}$ for $\eta$ satisfies

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

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 $\tilde{\mathbf{H}}_\eta := \mathrm{E}_{\eta\sim q} \mathbf{H}_\eta$ and $\tilde{\mathbf{H}}_\eta^2 := \mathrm{E}_{\eta\sim q} \mathbf{H}_\eta^2$. We use a Metropolis random-walk algorithm to obtain these quantities, as detailed in the algorithm below.

---

**Algorithm 1** Metropolis random-walk to sample $\eta$

---

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

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

   where

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

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

---

Now consider the case where $\eta = \{\lambda_1, \ldots, \lambda_p\}$ (RKHS scale parameters only), and the scenario described in the exponential family EM algorithm of <mark>Section 4.3.3</mark> applies. In particular, for $k = 1, \ldots, p$, 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 $j = 1, \ldots, m$, assuming each of

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

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

By completing the squares, we recognise this is as the kernel of a univariate normal density. Specifically, $\lambda_k \sim \mathrm{N}(d_k/c_k, 1/c_k)$. The quantity $\tilde{\mathbf{H}}_\eta$ can be obtained by substituting $\lambda_k \mapsto \mathrm{E}_{\lambda_k\sim q}[\lambda_k]$ in the <mark>expression XXX</mark>. However, in the calculation of $\tilde{\mathbf{H}}_\eta^2$, we must replace $\lambda_k^2 \mapsto \mathrm{E}_{\lambda_k\sim q}[\lambda_k]^2 + \mathrm{Var}_{\lambda_k\sim q}[\lambda_k]$ in all occurrences of square terms. This can be cumbersome, so if felt necessary, use the approximation $\lambda_k^2 \mapsto \mathrm{E}_{\lambda_k\sim q}[\lambda_k]^2$ 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 $v_k^2 \to 0$ and $m_k v_k^{-2} \to 0$. The mean field distribution for $\lambda$ is then

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

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

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

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

**Derivation of $\tilde{q}(\boldsymbol{\Psi})$**

Introduce the transformed random matrix $\mathbf{u} = \mathbf{w}\boldsymbol{\Psi}^{-1} \in \mathbb{R}^{n\times m}$. Since we have that $\mathrm{vec}\,\mathbf{u} = (\mathrm{vec}\,\mathbf{w})^\top(\boldsymbol{\Psi}^{-1}\otimes\mathbf{I}_n)$, the optimal mean-field distribution for $\mathbf{u}$ is normal with mean $\mathrm{vec}\,\tilde{\mathbf{u}} = \mathrm{vec}(\tilde{\mathbf{w}}\tilde{\boldsymbol{\Psi}}^{-1})$ and variance

$$\tilde{\mathbf{V}}_u = (\tilde{\boldsymbol{\Psi}}^{-1}\otimes\mathbf{I}_n)\tilde{\mathbf{V}}_w(\tilde{\boldsymbol{\Psi}}^{-1}\otimes\mathbf{I}_n).$$

In the case of the independent model, its mean is $\tilde{\mathbf{u}}_{\cdot j} = \tilde{\psi}_j^{-1}\tilde{\mathbf{u}}_{\cdot j}$ for $j = 1, \ldots, m$ and its variance is

$$\tilde{\mathbf{V}}_u = \mathrm{diag}(\tilde{\psi}_1^{-2}\tilde{\mathbf{V}}_{w_1}, \ldots, \tilde{\psi}_m^{-2}\tilde{\mathbf{V}}_{w_m}).$$

Now, to derive $\tilde{q}(\boldsymbol{\Psi})$ for the full I-probit model, we inspect the equation

$$\log\tilde{q}(\boldsymbol{\Psi}) = \mathrm{E}_{\mathcal{Z}\backslash\{\boldsymbol{\Psi}\}\sim q}\left[\frac{n}{2}\log|\boldsymbol{\Psi}| - \frac{1}{2}\mathrm{tr}\left((\mathbf{y}^*-\boldsymbol{\mu})^\top(\mathbf{y}^*-\boldsymbol{\mu})\boldsymbol{\Psi}\right) + \frac{n}{2}\log|\boldsymbol{\Psi}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{u}^\top\mathbf{u}\boldsymbol{\Psi}\right)\right]$$
$$+ \frac{g-m-1}{2}\log|\boldsymbol{\Psi}| - \frac{1}{2}\mathrm{tr}(\mathbf{G}\boldsymbol{\Psi}) + \mathrm{const.}$$
$$= -\frac{1}{2}\mathrm{tr}\left(\left(\mathbf{G} + \overbrace{\mathrm{E}[(\mathbf{y}^*-\boldsymbol{\mu})^\top(\mathbf{y}^*-\boldsymbol{\mu})]}^{\mathbf{G}_1} + \overbrace{\mathrm{E}[\mathbf{u}^\top\mathbf{u}]}^{\mathbf{G}_2}\right)\boldsymbol{\Psi}\right)$$
$$+ \frac{2n+g-m-1}{2}\log|\boldsymbol{\Psi}| + \mathrm{const.}$$

which we recognise to be a Wishart distribution with scale matrix $(\mathbf{G}+\mathbf{G}_1+\mathbf{G}_2)^{-1}$ and $2n+g-m$ degrees of freedom. The mean of this distribution is $\tilde{\boldsymbol{\Psi}} = (2n+g-m)(\mathbf{G}+\mathbf{G}_1+\mathbf{G}_2)^{-1}$.

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

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

$$\mathbf{G}_2 = \mathrm{E}[\mathbf{u}^\top \mathbf{u}]$$

$$=$$

**Derivation $\tilde{q}(\boldsymbol{\alpha})$**

For $j = 1, \ldots, m$, denote $\mathbf{H}_i$ as the row vector of the kernel matrix $\mathbf{H}$. Then,

$$\log \tilde{q}(\alpha) = \mathrm{E}_{\mathbf{y}^*,\mathbf{w},\lambda}\left[ -\frac{1}{2} \sum_{j=1}^m \sum_{i=1}^n \left( y_{ij}^* - \alpha_j - \lambda_j \sum_{k=1}^n h(x_i, x_k) w_{kj} \right)^2 \right] + \mathrm{const.}$$

$$= -\frac{1}{2} \sum_{j=1}^m \mathrm{E}_{\mathbf{y}^*,\mathbf{w},\lambda}\left[ n\alpha_j^2 - 2\alpha_j \sum_{i=1}^n (y_{ij}^* - \lambda_j \mathbf{H}_i \mathbf{w}_j) \right] + \mathrm{const.}$$

$$= -\frac{n}{2} \sum_{j=1}^m \left[ \left( \alpha_j - \frac{1}{n} \sum_{i=1}^n (\mathrm{E}[y_{ij}^*] - \mathrm{E}[\lambda_j] \mathbf{H}_i \mathbf{w}_j) \right)^2 \right] + \mathrm{const.}$$

which is of course the kernel of the product of $m$ univariate normal densities, each with mean and variance

$$\tilde{\alpha}_j = \frac{1}{n} \sum_{i=1}^n \left( \mathrm{E}[y_{ij}^*] - \mathrm{E}[\lambda_j] \mathbf{H}_i \, \mathrm{E}[\mathbf{w}_j] \right) \ \text{ and } \ v_{\alpha_j} = \frac{1}{n}.$$

Suppose that we use a single intercept parameter $\alpha$. In this case, $\alpha$ is is also normally distributed with mean and variance

$$\tilde{\alpha} = \frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n \left( \mathrm{E}[y_{ij}^*] - \mathrm{E}[\lambda_j] \mathbf{H}_i \, \mathrm{E}[\mathbf{w}_j] \right) \ \text{ and } \ v_\alpha = \frac{1}{nm}.$$

### 5.10.4 Monitoring the lower bound

A convergence criterion would be when there is no more significant increase in the lower bound $\mathcal{L}$, as defined by

$$
\begin{aligned}
\mathcal{L} &= \int q(\mathbf{y}^*, \mathbf{w}, \lambda, \alpha) \log\left[\frac{p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \lambda, \alpha)}{q(\mathbf{y}^*, \mathbf{w}, \lambda, \alpha)}\right] \mathrm{d}\mathbf{y}^* \mathrm{d}\mathbf{w}\mathrm{d}\lambda\mathrm{d}\alpha \\
&= \mathrm{E}[\log p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \lambda, \alpha)] - \mathrm{E}[\log q(\mathbf{y}^*, \mathbf{w}, \lambda, \alpha)] \\
&= \mathrm{E}\left[\log \prod_{i=1}^{n}\prod_{j=1}^{m} \cancel{p(y_i|y_{ij}^*)}\right] + \mathrm{E}\left[\log p(\mathbf{y}^*|\mathbf{f})\right] + \mathrm{E}\left[\log p(\mathbf{w})\right] + \cancel{\underline{\mathrm{E}\left[\log p(\lambda)\right]}} + \cancel{\underline{\mathrm{E}\left[\log p(\alpha)\right]}} \\
&\quad - \mathrm{E}\left[\log q(\mathbf{y}^*)\right] - \mathrm{E}\left[\log q(\mathbf{w})\right] - \mathrm{E}\left[\log q(\lambda)\right] - \mathrm{E}\left[\log q(\alpha)\right]
\end{aligned}
$$

Note that the categorical pmf $p(y_i|y_{ij}^*)$ becomes degenerate once the latent variables are known, so this term is cancelled out. With the exception of $q(\mathbf{y}^*)$, all of the distributions are Gaussian. The following results will be helpful.

**Definition 5.2** (Differential entropy)**.** The differential entropy $\mathcal{H}$ of a pdf $p(x)$ is given by

$$
\mathcal{H}(p) = -\int p(x) \log p(x) \mathrm{d}x = -\mathrm{E}_p[\log p(x)].
$$

thm:normentropy

**Lemma 5.3.** *Let $p(x)$ be the pdf of a random variable $x$. Then if*

*(i) $p$ is a univariate normal distribution with mean $\mu$ and variance $\sigma^2$,*

$$
\mathcal{H}(p) = \frac{1}{2}(1 + \log 2\pi) + \frac{1}{2}\log \sigma^2
$$

*(ii) $p$ is a d-dimensional normal distribution with mean $\mu$ and variance $\Sigma$,*

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

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

$$\mathrm{E}\left[\log p(\mathbf{y}^*|\mathbf{f})\right] - \mathrm{E}\left[\log q(\mathbf{y}^*)\right] = \sum_{i=1}^{n}\sum_{j=1}^{m} \mathrm{E}\left[\log p(y_{ij}^*|f_{ij})\right] + \sum_{i=1}^{n} \mathcal{H}\big(q(y_i^*)\big)$$

$$= \sum_{i=1}^{n}\sum_{j=1}^{m}\left(-\frac{1}{2}\log 2\pi - \frac{1}{2}\mathrm{E}[y_{ij}^* - f_{ij}]^2\right)$$

$$+ \sum_{i=1}^{n}\sum_{j=1}^{m}\left(\frac{1}{2}\log 2\pi + \frac{1}{2}\mathrm{E}[y_{ij}^* - f_{ij}]^2\right) + \sum_{i=1}^{n}\log C_i$$

**Terms involving distributions of $\mathbf{w}$**

$$\mathrm{E}\left[\log p(\mathbf{w})\right] - \mathrm{E}\left[\log q(\mathbf{w})\right] = \sum_{j=1}^{m}\Big(\mathrm{E}\left[\log p(\mathbf{w}_j)\right] - \mathrm{E}\left[\log q(\mathbf{w}_j)\right]\Big)$$

$$= \sum_{j=1}^{m}\left(-\frac{n}{2}\log 2\pi - \frac{1}{2}\mathrm{E}[\mathbf{w}_j^\top\mathbf{w}_j] + \mathcal{H}\big(q(\mathbf{w}_j)\big)\right)$$

$$= \sum_{j=1}^{m}\left(-\frac{n}{2}\log 2\pi - \frac{1}{2}\mathrm{tr}\left(\mathrm{E}[\mathbf{w}_j\mathbf{w}_j^\top]\right) + \frac{n}{2}(1 + \log 2\pi) - \frac{1}{2}\log|\mathbf{A}_j|\right)$$

$$= \frac{nm}{2} - \frac{1}{2}\sum_{j=1}^{m}\left(\mathrm{tr}\,\widetilde{\mathbf{W}}_j + \log|\mathbf{A}_j|\right)$$

**Terms involving distribution of $q(\lambda)$**

$$-\mathrm{E}\left[\log q(\lambda)\right] = \sum_{j=1}^{m}\mathcal{H}\big(q(\lambda_j)\big)$$

$$= \sum_{j=1}^{m}\left(\frac{1}{2}(1 + \log 2\pi) - \frac{1}{2}\log c_j\right)$$

$$= \frac{m}{2}(1 + \log 2\pi) - \frac{1}{2}\sum_{j=1}^{m}\log c_j$$

31

or if using single $\lambda$

$$- \mathrm{E}\left[\log q(\lambda)\right] = \frac{1}{2}(1 + \log 2\pi) - \frac{1}{2}\log\sum_{j=1}^{m}c_j.$$

**Terms involving distribution of $q(\alpha)$**

$$-\mathrm{E}\left[\log q(\alpha)\right] = \sum_{j=1}^{m}\mathcal{H}\big(q(\alpha_j)\big)$$
$$= \frac{m}{2}(1 + \log 2\pi - \log n)$$

or if using single $\alpha$

$$-\mathrm{E}\left[\log q(\alpha)\right] = \frac{1}{2}(1 + \log 2\pi - \log nm).$$

**The lower bound**

$$\mathcal{L} = \sum_{i=1}^{n}\log C_i + \frac{nm}{2} - \frac{1}{2}\sum_{j=1}^{m}\left(\operatorname{tr}\widetilde{\mathbf{W}}_j + \log|\mathbf{A}_j|\right)$$
$$+ \frac{m}{2}(1 + \log 2\pi) - \frac{1}{2}\sum_{j=1}^{m}\log c_j + \frac{m}{2}(1 + \log 2\pi - \log n)$$
$$= \frac{m}{2}\big(n + 2(1 + \log 2\pi) - \log n\big) - \frac{1}{2}\sum_{j=1}^{m}\left(\operatorname{tr}\widetilde{\mathbf{W}}_j + \log|\mathbf{A}_j| + \log c_j\right) + \sum_{i=1}^{n}\log C_i$$

Of course, if using either single $\alpha$ or single $\lambda$, then the formula needs to be adjusted accordingly.

# Bibliography

**bishop2006pattern** Bishop, Christopher (2006). *Pattern Recognition and Machine Learning.* Springer-Verlag.

**girolami2006variational** Girolami, Mark and Simon Rogers (2006). "Variational Bayesian Multinomial Probit Regression with Gaussian Process Priors". In: *Neural Computation* 18.8, pp. 1790–1817.

**groves1969note** Groves, Theodore and Thomas Rothenberg (1969). "A note on the expected value of an inverse matrix". In: *Biometrika* 56.3, pp. 690–691.

**hastie1986** Hastie, Trevor and Robert Tibshirani (Aug. 1986). "Generalized Additive Models". In: *Statist. Sci.* 1.3, pp. 297–310. DOI: 10.1214/ss/1177013604. URL: https://doi.org/10.1214/ss/1177013604.

**jamil2017** Jamil, Haziq and Wicher Bergsma (2017). "iprior: An R Package for Regression Modelling using I-priors". In: *Manuscript in submission.*

**mccullagh1989** McCullagh, P. and John A. Nelder (1989). *Generalized Linear Models.* 2nd. Chapman & Hall/CRC Press.

**mcculloch2000bayesian** McCulloch, Robert E, Nicholas G Polson, and Peter E Rossi (2000). "A Bayesian analysis of the multinomial probit model with fully identified parameters". In: *Journal of econometrics* 99.1, pp. 173–193.

**mcculloch1994exact** McCulloch, Robert and Peter E Rossi (1994). "An exact likelihood analysis of the multinomial probit model". In: *Journal of Econometrics* 64.1, pp. 207–240.

**meng1997algorithm** Meng, Xiao-Li and David Van Dyk (1997). "The EM Algorithm—an Old Folk-song Sung to a Fast New Tune". In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 59.3, pp. 511–567.

| | |
|---|---|
| minka2001ex pectation | Minka, Thomas P (2001). "Expectation propagation for approximate Bayesian inference". In: *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence.* Morgan Kaufmann Publishers Inc., pp. 362–369. |
| neal1999 | Neal, Radford M. (1999). "Regression and Classification using Gaussian Process Priors". In: *Bayesian Statistics.* 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. |
| nobile1998h ybrid | Nobile, Agostino (1998). "A hybrid Markov chain for the Bayesian analysis of the multinomial probit model". In: *Statistics and Computing* 8.3, pp. 229–242. |
| rasmussen20 06gaussian | Rasmussen, Carl Edward and Christopher K I Williams (2006). *Gaussian Processes for Machine Learning.* The MIT Press. |
| scholkopf20 02learning | Schölkopf, Bernhard and Alexander J Smola (2002). *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond.* MIT Press. |