

## To-do list

1. But isn't identifiability resolved by setting all covariances to zero? . . . . .	3
2. Citation? . . . . .	4
3. Add section for Laplace's method . . . . .	5
4. Add section for MCMC . . . . .	5
5. Add section for variational inference . . . . .	5
6. This isn't the variational EM... CAVI? . . . . .	14
7. For future work, can consider estimating the covariances/correlations across choices. More suitable for social science data. . . . .	20
8. Describe the model when there are only two alternatives. . . . .	21
9. It's definitely possible to extend to multiple scale parameters. It's just a matter of algebra.	21
10. A different model using thresholds, but it might be possible to model these and estimate using variational inference. . . . .	21

# Multinomial probit I-prior models

Haziq Jamil

*Department of Statistics*

*London School of Economics and Political Science*

May 23, 2017

## Abstract

Extension of I-prior models to multi-class classification with estimation using variational inference.

**Keywords:** some, keywords, go, here

## 1 Introduction

Consider multinomial response variables  $y_1, \dots, y_n$ , where each response  $y_i$  takes on one of the values  $\{1, \dots, m\}$  from a set of  $m$  possible values. We model each response as following a categorical distribution (a special case of the multinomial distribution)

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

with probability mass function (pmf)

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

such that  $p_{ij} \geq 0$  for each  $j$  and  $\sum_{j=1}^m p_{ij} = 1$ . It might also be convenient to think of the responses  $y_i$  as comprising of a binary vector of length  $m$ , with a single ‘1’ at the position corresponding to the value that  $y_i$  takes. That is,

$$y_i = (y_{i1}, \dots, y_{im})$$

with

$$y_{ik} = \begin{cases} 1 & \text{if } y_i = k \\ 0 & \text{if } y_i \neq k \end{cases}$$

and  $\sum_{j=1}^m y_{ij} = 1$ . In this formulation, each  $y_{ij}$  is distributed as Bernoulli with probability  $p_{ij}$ .

Suppose for each observation  $y_i$  there is an associated  $p$ -dimensional vector of covariates  $x_i = (x_{i1}, \dots, x_{ip})$  belonging to some set  $\mathcal{X}$ . We would like to model the multinomial outcomes  $y_i$  based on these vectors of covariates. Assume that for each  $y_i = (y_{i1}, \dots, y_{im})$ , there exists a

corresponding continuous, underlying, latent variable  $y_i^* = (y_{i1}^*, \dots, y_{im}^*)$  such that

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

In other words,

$$y_{ij} = \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*].$$

We consider modelling the  $n$  latent variables in each class  $j = 1, \dots, m$  according to the regression problem

$$\begin{aligned} y_{ij}^* &= f_j(x_i) + \epsilon_{ij} \\ \epsilon_{ij} &\stackrel{\text{iid}}{\sim} \text{N}(0, \sigma_j^2) \\ i &= 1, \dots, n \end{aligned} \tag{1}$$

with  $f_j : \mathcal{X} \rightarrow \mathbb{R}$  being a regression function belonging to some reproducing kernel Hilbert space of functions  $\mathcal{F}_j$  having the reproducing kernel  $h_{\lambda_j} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Here,  $\lambda_j$  is the scale parameter for the reproducing kernel, so that  $h_{\lambda_j}(\cdot) = \lambda_j h(\cdot)$ . One advantage of the multinomial probit model is the ability to also model the correlations across choices, such that for each  $j, k \in \{1, \dots, m\}$  and  $j \neq k$ ,

$$\text{Corr}[\epsilon_{ij}, \epsilon_{ik}] = \frac{\sigma_{jk}}{\sigma_j \sigma_k}.$$

This setting is suitable for modelling multinomial data where the independence axiom is not desired. Such cases arise frequently in economics and social science. The famous Red-Bus-Blue-Bus example is often used to illustrate independence of irrelevant alternatives (IIA). Suppose commuters face the decision between taking cars and red busses. The addition of blue busses to commuters' choice should not affect the choice between cars or busses (assuming commuters are indifferent about the colour of the bus). Yet, if IIA is imposed, then the three choices are distinct, and the fact that red and blue busses are substitutable is ignored.

The IIA assumption is realised by fixing  $\sigma_{jk} = 0$ ,  $j \neq k$ , which is clearly a simplification of the model (and as we will see later, benefits us in the algebra when deriving some distributional results). As this may not be suitable for certain modelling purposes, we might want to consider how this assumption can be relaxed, but we leave this for future work.

Back to model (1): We wish to model the function  $f_j$  as having an I-prior. Denoting  $f_{ij} = f_j(x_i)$  as the evaluation of the function  $f_j(\cdot)$  at  $x_i$ , and also  $\mathbf{f}_j = (f_{1j}, \dots, f_{nj})^\top$  as the vector containing all  $n$  evaluations pertaining to the  $j$ th alternative, an I-prior on  $f_j$  is

$$\mathbf{f}_j \sim \text{N}(\mathbf{f}_j^0, \mathcal{I}_j)$$

where  $\mathcal{I}_j$  is the  $n \times n$  Fisher covariance kernel for the regression function  $f_j$  in model (1), which has  $(r, s)$  entries given by

$$\mathcal{I}_j(f_j(x_r), f_j(x_s)) = \sigma_j^{-2} \sum_{k=1}^n \sum_{l=1}^n h_{\lambda_j}(x_r, x_k) h_{\lambda_j}(x_l, x_s)$$

and  $\mathbf{f}_j^0$  is a vector of prior means. By concatenating the vectors  $\mathbf{f}_1, \dots, \mathbf{f}_m$  into the vector  $\mathbf{f}$  of length  $nm$ , it is easy to see that

$$\mathbf{f} \sim \mathcal{N}(\mathbf{f}^0, \mathbf{\Omega})$$

where  $\mathbf{f}^0 = (\mathbf{f}_1^0, \dots, \mathbf{f}_m^0)^\top$  and  $\mathbf{\Omega} = \text{diag}(\mathcal{I}_1, \dots, \mathcal{I}_m)$ . This stems from the fact that the error components are independent across choices (IIA), and as such,  $\text{Cov}[\mathbf{f}_j, \mathbf{f}_k] = \mathbf{0}$ ,  $j \neq k$ . In the more general case where  $\sigma_{jk} \neq 0$ , then extra care must be taken to ensure the covariances are represented in the I-prior covariance matrix.

We make two further simplifications to the model. Firstly, the choice model as stated is not identified in scale, i.e. multiplication of the latent variables by a positive constant does not make any difference to the outcome. This is a well known identification issue<sup>1</sup> with the multinomial probit model, and this is typically overcome by setting some restrictions. In our case, we set all  $\sigma_j^2 = 1$ .

Secondly, we assume that the intercept functions (prior means) are constants, so that  $\mathbf{f}_j^0 = \alpha_j \mathbf{1}_n$  for  $j = 1, \dots, m$ , and the  $\alpha_j$ s are just additional hyperparameters to be estimated. With all of these in mind, the I-prior model simplifies to

$$\begin{aligned} y_{ij}^* &= \alpha_j + \sum_{k=1}^n h_{\lambda_j}(x_i, x_k) w_{kj} + \epsilon_{ij} \\ \epsilon_{ij} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1) \\ w_{kj} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1) \end{aligned}$$

for each  $j \in \{1, \dots, m\}$ . Define  $f_{ij} = \alpha_j + \sum_{k=1}^n h_{\lambda_j}(x_i, x_k) w_{kj}$ , so that each  $y_{ij}^* | f_{ij} \sim \mathcal{N}(f_{ij}, 1)$ . The probit link is seen as follows:

$$\begin{aligned} p_{ij} &= \mathbb{P} \left[ y_{ij}^* = \max_k y_{ik}^* \right] \\ &= \mathbb{P} [y_{ij}^* > y_{ik}^* : \forall k \neq j] \\ &= \mathbb{P} [f_{ij} + \epsilon_{ij} > f_{ik} + \epsilon_{ik} : \forall k \neq j] \\ &= \mathbb{P} [\epsilon_{ik} - \epsilon_{ij} \leq f_{ij} - f_{ik} : \forall k \neq j] \\ &= \int \cdots \int \mathbb{1} [\epsilon_{ik} \leq \epsilon_{ij} + f_{ij} - f_{ik} : \forall k \neq j] \prod_{k=1}^m [\phi(\epsilon_{ik}) d\epsilon_{ik}] \\ &= \int \prod_{\substack{k=1 \\ k \neq j}}^m \Phi(\epsilon_{ij} + f_{ij} - f_{ik}) \phi(\epsilon_{ij}) d\epsilon_{ij} \\ &= \mathbb{E}_Z \left[ \prod_{\substack{k=1 \\ k \neq j}}^m \Phi(Z + f_{ij} - f_{ik}) \right] \end{aligned} \tag{2}$$

where  $Z$  is a standard normal random variable, and  $\phi(\cdot)$  and  $\Phi(\cdot)$  are the pdf and cdf of a standard normal distribution respectively. It is well known that for  $m > 3$  this has no closed form expression, which makes the probit model unattractive from a likelihood maximisation

1. But isn't identifiability resolved by setting all covariances to zero?

<sup>1</sup>In the unrestricted case for a model with  $m$  alternatives, there would be  $m(m+1)/2$  variance components to estimate. However, in general, only  $m(m-1)/2$  can be freely estimated.

The above describes  $m$  regression functions being estimated for each class, and each of the  $m$  regression functions estimated using an I-prior. It is possible for all  $m$  regression functions<sup>2</sup> to share a common I-prior  $\mathbf{f}_1, \dots, \mathbf{f}_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\alpha \mathbf{1}_n, \mathcal{I})$ , so that only one set of intercept and RKHS scale parameters need to be estimated (instead of  $m$  sets, one for each class). There is also flexibility in using the same covariance kernel for instance, but different intercepts for the  $m$  I-priors, or vice-versa. The probit I-prior model can be represented by the following DAG.

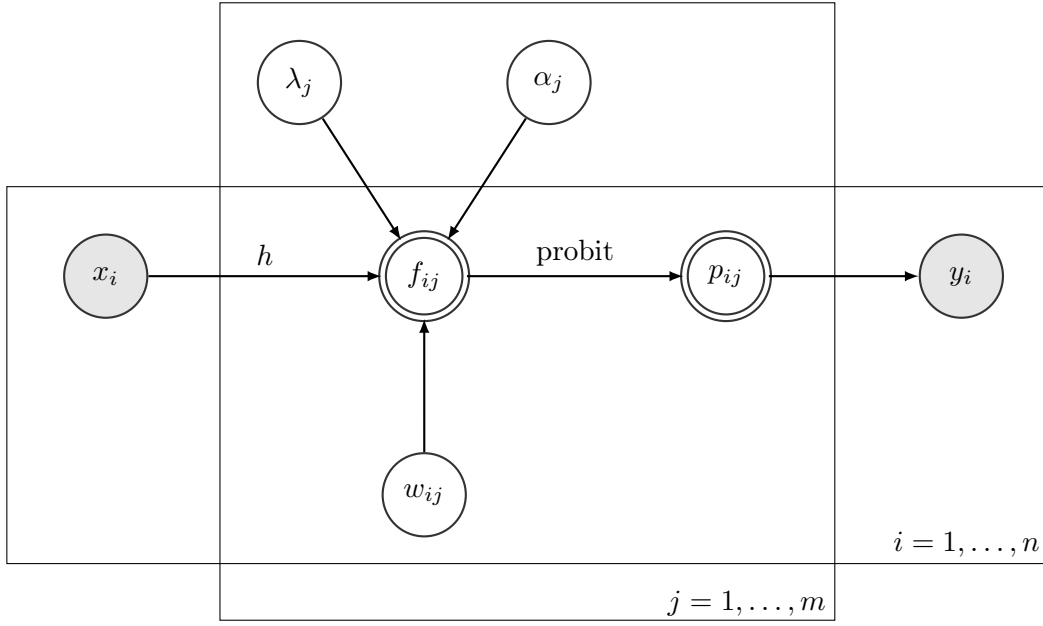


Figure 1: A DAG of the probit I-prior model. Observed nodes are shaded, while double-lined nodes represented known or calculable quantities. The latent variables  $y_{ij}^*$  have been marginalised and absorbed into the probit. The  $w_{ij}$ s are the standard normal random-effects associated with the I-prior. There are at most  $m$  sets of intercept ( $\alpha_j$ ) and scale ( $\lambda_j$ ) parameters to estimate.

## 2 Estimation

The parameters to estimate in the probit I-prior model are  $\boldsymbol{\theta} = (\alpha_1, \dots, \alpha_m, \lambda_1, \dots, \lambda_m)$ . The likelihood function  $L(\cdot)$  for  $\boldsymbol{\theta}$  is obtained by integrating out the I-prior from the multinomial likelihood, as follows:

$$\begin{aligned} L(\boldsymbol{\theta}) &= \int \prod_{i=1}^n \prod_{j=1}^m p(y_i | f_{ij}) p(f_{ij} | \alpha_j, \lambda_j) d\mathbf{f} \\ &= \int \dots \int \prod_{i=1}^n \prod_{j=1}^m p_{ij}^{\mathbb{1}[y_i=j]} p(\mathbf{f}_j | \alpha_j, \lambda_j) d\mathbf{f}_1 \dots d\mathbf{f}_m. \end{aligned}$$

<sup>2</sup>It is also possible to reparameterise the model (anchoring on one latent variable as the reference class and working with the latent differences) so that only  $m - 1$  I-priors are required.

From (2), we know that  $p_{ij}$  is related to the  $f_{ij}$  via the integral involving the CDF and PDF of a standard normal. Thus, the intractable integral above presents a practical challenge which makes estimation via direct maximisation of the likelihood difficult to accomplish. Several approximations are considered, and discussed below.

## 2.1 Laplace approximation

Add section for Laplace's method

## 2.2 Markov Chain Monte Carlo

Add section for MCMC

## 2.3 Variational inference

Add section for variational inference

## 3 Estimation using variational inference

The variational approximation to the probit I-prior model is made tractable by the inclusion of the latent variables  $y_{ij}^*$  in the analysis as an intermediary step. The unknown quantities to estimate via variational inference are  $y_{ij}^*$ ,  $w_{ij}$ ,  $\lambda_j$  and  $\alpha_j$ .

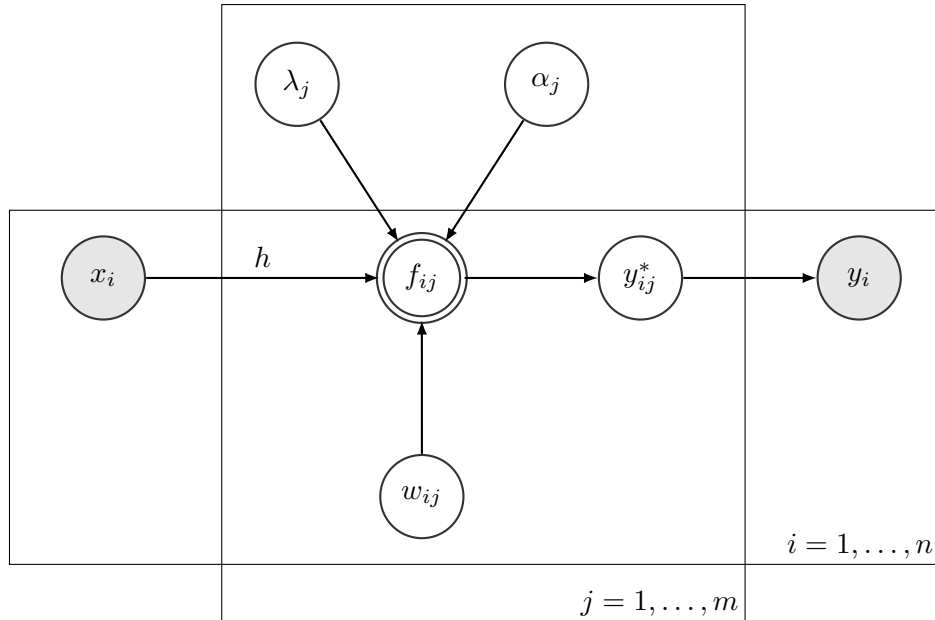


Figure 2: A DAG of the variational probit I-prior model.

### 3.1 The joint likelihood

$$\begin{aligned} p(\mathbf{y}, \mathbf{y}^*, \mathbf{w}, \lambda, \alpha) &= p(\mathbf{y}|\mathbf{y}^*, \mathbf{w}, \alpha, \lambda) p(\mathbf{y}^*, \mathbf{f}, \mathbf{w}, \alpha, \lambda) \\ &= p(\mathbf{y}|\mathbf{y}^*) p(\mathbf{y}^*|\mathbf{f}) p(\mathbf{w}) p(\lambda) p(\alpha) \end{aligned}$$

### 3.2 Relevant distributions

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

$$\begin{aligned} p(\mathbf{y}^*|\mathbf{f}) &= \prod_{i=1}^n \prod_{j=1}^m \mathcal{N}(f_{ij}, 1) \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m (y_{ij}^* - f_{ij})^2 \right] \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \|\mathbf{y}^* - \mathbf{f}\|^2 \right] \end{aligned}$$

$$f_{ij} = \alpha_j + \sum_{k=1}^n h_{\lambda_j}(x_i, x_k) w_{kj}$$

$$\begin{aligned} p(\mathbf{w}) &= \prod_{i=1}^n \prod_{j=1}^m p(w_{ij}) \\ &= \prod_{i=1}^n \prod_{j=1}^m \mathcal{N}(0, 1) \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^m w_{ij}^2 \right] \\ &= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \mathbf{w}^\top \mathbf{w} \right] \end{aligned}$$

$$p(\lambda, \alpha) \propto \text{const.}$$

### 3.3 Mean field approximation

$$\begin{aligned} p(\mathbf{y}^*, \mathbf{w}, \alpha, \lambda | \mathbf{y}) &\equiv q(\mathbf{y}^*) q(\mathbf{w}) q(\lambda) q(\alpha) \\ &\equiv \prod_{i,j} q(y_{ij}^*) q(\mathbf{w}) q(\lambda) q(\alpha) \end{aligned}$$

The first line is by assumption, while the second line follows from an induced factorisation, as we will see later. Denote by  $\tilde{q}$  the distributions which minimise the KL divergence (maximises the lower bound). Then, for each of  $\xi \in \{\mathbf{y}^*, \mathbf{w}, \alpha, \lambda\}$ ,  $\tilde{q}$  satisfies

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

### 3.3.1 $\tilde{q}(\mathbf{y}^*)$

In this subsection, we use the notation  $y_i^* = (y_{i1}^*, \dots, y_{im}^*)$  to denote the vector of length  $m$  containing the latent variables for response  $i$ . The joint distribution for  $\mathbf{y}^* = (y_1^*, \dots, y_n^*)^\top$  is a product of the distribution for each of the components  $y_i^*$  - this is a consequence of the independence structure across observations. Therefore, we can consider the variational density for each  $y_i^*$  separately.

Consider the case where  $y_i$  takes one particular value  $j \in \{1, \dots, m\}$ . The mean-field density  $q(y_i^*)$  for each  $i = 1, \dots, n$  is found to be

$$\begin{aligned} \log \tilde{q}(y_i^*) &= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \cdot \mathbb{E}_{\mathbf{w}, \alpha, \lambda} \left[ -\frac{1}{2} \sum_{k=1}^m (y_{ik}^* - f_{ik})^2 \right] + \text{const.} \\ &= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \cdot \left[ -\frac{1}{2} \sum_{k=1}^m (y_{ik}^* - \tilde{f}_{ik})^2 \right] + \text{const.} \\ &\equiv \begin{cases} \prod_{k=1}^m \mathcal{N}(\tilde{f}_{ik}, 1) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

where  $\tilde{f}_{ik} = \mathbb{E}[\alpha_k] + \sum_{l=1}^m h_{\mathbb{E}[\lambda_k]}(x_i, x_l) \mathbb{E}[w_{il}]$ , and expectations are taken under the optimal mean-field distribution  $\tilde{q}$ . The distribution for  $q(y_i^*)$  is a truncated  $m$ -variate normal distribution such that the  $j$ th component is always largest. It is worth investigating the properties of this distribution, and we now present some relevant definitions and results.

**Definition 1** (Conically-truncated multivariate normal distribution). *Let  $\mathbf{X} = (X_1, \dots, X_d)$  be a  $d$ -dimensional random variable with pdf defined as*

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

for some  $j \in \{1, \dots, d\}$ . We denote the distribution of  $\mathbf{X}$  by  $\mathcal{N}^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , with  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$  and  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ . 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:

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

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

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



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

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

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

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

the proof follows directly. □

**Lemma 1.** Let  $X \sim N^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with pdf  $p(\mathbf{x})$  as defined in Definition 1. Then

(i) The expectation  $\mathbb{E}[\mathbf{X}] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_d])$  is given by

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

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

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

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

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

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

As we know,  $y_i$  takes on any one value from the set  $\{1, \dots, m\}$ . Thus, we have that the distribution of  $(y_{i1}^*, \dots, y_{im}^*)$  is  $\mathcal{N}^{(y_i)}(\boldsymbol{\mu}_i, \mathbf{I}_m)$ , where  $\boldsymbol{\mu}_i = (\tilde{f}_{i1}, \dots, \tilde{f}_{im})$ . The expectation is given by

$$\mathbb{E}[y_{ik}^*] = \begin{cases} \tilde{f}_{ik} - C_i^{-1} \mathbb{E}_Z \left[ \phi_{ik}(Z) \prod_{l \neq k, y_i} \Phi_{il}(Z) \right] & \text{if } k \neq y_i \\ \tilde{f}_{iy_i} - \sum_{k \neq y_i} (\mathbb{E}[y_{ik}^*] - \tilde{f}_{ik}) & \text{if } k = y_i \end{cases}$$

where

$$\begin{aligned}\phi_{ik}(Z) &= \phi(Z + \tilde{f}_{iy_i} - \tilde{f}_{ik}) \\ \Phi_{ik}(Z) &= \Phi(Z + \tilde{f}_{iy_i} - \tilde{f}_{ik}) \\ C_i &= \mathbb{E}_Z \left[ \prod_{\substack{i=1 \\ i \neq j}}^d \Phi \left( Z + \tilde{f}_{iy_i} - \tilde{f}_{ik} \right) \right]\end{aligned}$$

and  $Z \sim \mathcal{N}(0, 1)$  with PDF and CDF  $\phi(\cdot)$  and  $\Phi(\cdot)$  respectively. In order to calculate these expectations, we need to compute the following integrals:

$$\begin{aligned}\mathbb{E}_Z \left[ \phi_{ik}(Z) \prod_{l \neq k, j} \Phi_{il}(Z) \right] &= \int \phi_{ik}(z) \prod_{l \neq k, j} \Phi_{il}(z) \phi(z) dz, \quad \forall k \neq y_i \\ C_i &= \mathbb{E}_Z \left[ \prod_{l \neq j} \Phi_{il}(Z) \right] = \int \prod_{l \neq j} \Phi_{il}(z) \phi(z) dz\end{aligned}$$

Since these are functions of a Gaussian pdf, these can be computed rather efficiently using quadrature methods.

### 3.3.2 $\tilde{q}(\mathbf{w})$

For each  $j = 1, \dots, m$ , denote  $\mathbf{y}_j^* = (y_{1j}^*, \dots, y_{nj}^*)^\top$  as the vector of length  $n$  containing all latent observations for each class. Then,

$$\log \tilde{q}(\mathbf{w}) = \mathbb{E}_{\mathbf{y}^*, \alpha, \lambda} \left[ -\frac{1}{2} \sum_{j=1}^m \|\mathbf{y}_j^* - \alpha_j \mathbf{1}_n - \mathbf{H}_{\lambda_j} \mathbf{w}_j\|^2 - \frac{1}{2} \sum_{j=1}^m \|\mathbf{w}_j\|^2 \right] + \text{const.}$$

$$\begin{aligned}
&= -\frac{1}{2} \sum_{j=1}^m \mathbb{E}_{\mathbf{y}^*, \alpha, \lambda} \left[ \mathbf{w}_j^\top \mathbf{H}_{\lambda_j}^2 \mathbf{w}_j + \mathbf{w}_j^\top \mathbf{w}_j - 2(\mathbf{y}_j^* - \alpha_j \mathbf{1}_n)^\top \mathbf{H}_{\lambda_j} \mathbf{w}_j \right] + \text{const.} \\
&= -\frac{1}{2} \sum_{j=1}^m \left( \mathbf{w}_j^\top (\mathbb{E}[\mathbf{H}_{\lambda_j}^2] + \mathbf{I}_n) \mathbf{w}_j - 2(\mathbb{E}[\mathbf{y}_j^*] - \mathbb{E}[\alpha_j] \mathbf{1}_n)^\top \mathbb{E}[\mathbf{H}_{\lambda_j}] \mathbf{w}_j \right) + \text{const.}
\end{aligned}$$

Let  $\mathbf{A}_j = \mathbb{E}[\mathbf{H}_{\lambda_j}^2] + \mathbf{I}_n$  and  $\mathbf{a}_j = \mathbb{E}[\mathbf{H}_{\lambda_j}](\mathbb{E}[\mathbf{y}_j^*] - \mathbb{E}[\alpha_j] \mathbf{1}_n)$ . Then, using the fact that

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

we see the  $\log \tilde{q}(\mathbf{w})$  is a sum of quadratic terms in  $\mathbf{w}_j$ , and we recognise this as the kernel of the product of independent multivariate normal densities. Therefore, for each  $j = 1, \dots, m$ ,

$$\tilde{q}(\mathbf{w}_j) \equiv \mathcal{N}(\mathbf{A}_j^{-1} \mathbf{a}_j, \mathbf{A}_j^{-1}),$$

and  $\tilde{q}(\mathbf{w}) = \prod_{j=1}^m \tilde{q}(\mathbf{w}_j)$ . Because of this induced factorisation, we can obtain mean-field densities for each  $\mathbf{w}_j$  separately. For convenience later in deriving the lower bound, we note that the second moment of  $\tilde{q}(\mathbf{w}_j)$  is equal to  $\mathbb{E}[\mathbf{w}_j \mathbf{w}_j^\top] = \mathbf{A}_j^{-1} (\mathbf{I}_n + \mathbf{a}_j \mathbf{a}_j^\top \mathbf{A}_j^{-1}) =: \tilde{\mathbf{W}}_j$ .

### 3.3.3 $\tilde{q}(\lambda)$

For  $j = 1, \dots, m$ ,

$$\begin{aligned}
\log \tilde{q}(\lambda_j) &= \mathbb{E}_{\mathbf{y}^*, \mathbf{w}, \alpha} \left[ -\frac{1}{2} \sum_{j=1}^m \|\mathbf{y}_j^* - \alpha_j \mathbf{1}_n - \lambda_j \mathbf{H} \mathbf{w}_j\|^2 \right] + \text{const.} \\
&= -\frac{1}{2} \sum_{j=1}^m \mathbb{E}_{\mathbf{y}^*, \mathbf{w}, \alpha} \left[ \lambda_j^2 \mathbf{w}_j^\top \mathbf{H}^2 \mathbf{w}_j - 2\lambda_j (\mathbf{y}_j^* - \alpha_j \mathbf{1}_n)^\top \mathbf{H} \mathbf{w}_j \right] + \text{const.} \\
&= -\frac{1}{2} \sum_{j=1}^m \left( \lambda_j^2 \text{tr}(\mathbf{H}^2 \mathbb{E}[\mathbf{w}_j \mathbf{w}_j^\top]) - 2\lambda_j (\mathbb{E}[\mathbf{y}_j^*] - \mathbb{E}[\alpha_j] \mathbf{1}_n)^\top \mathbf{H} \mathbb{E}[\mathbf{w}_j] \right) + \text{const.}
\end{aligned}$$

By completing the squares, we recognise this is as the kernel of the product of independent univariate normal densities. Thus, each  $\lambda_j \sim \mathcal{N}(d_j/c_j, 1/c_j)$ , where

$$c_j = \text{tr}(\mathbf{H}^2 \mathbb{E}[\mathbf{w}_j \mathbf{w}_j^\top]) \quad \text{and} \quad d_j = (\mathbb{E}[\mathbf{y}_j^*] - \mathbb{E}[\alpha_j] \mathbf{1}_n)^\top \mathbf{H} \mathbb{E}[\mathbf{w}_j].$$

Supposing we use the same covariance kernel (and therefore scale parameter) for each regression class, the distribution for  $\lambda$  is easily seen as

$$\lambda \sim \mathcal{N}\left(\frac{\sum_{j=1}^m d_j}{\sum_{j=1}^m c_j}, \frac{1}{\sum_{j=1}^m c_j}\right).$$

### 3.3.4 $\tilde{q}(\alpha)$

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

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

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

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

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

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

## 3.4 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

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

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

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

**Lemma 2.** 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  $\Sigma$ ,

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

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

$$\begin{aligned} \mathbb{E} [\log p(\mathbf{y}^* | \mathbf{f})] - \mathbb{E} [\log q(\mathbf{y}^*)] &= \sum_{i=1}^n \sum_{j=1}^m \mathbb{E} [\log p(y_{ij}^* | f_{ij})] + \sum_{i=1}^n \mathcal{H}(q(y_i^*)) \\ &= \sum_{i=1}^n \sum_{j=1}^m \left( -\frac{1}{2} \log 2\pi - \frac{1}{2} \mathbb{E}[y_{ij}^* - f_{ij}]^2 \right) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^m \left( \frac{1}{2} \log 2\pi + \frac{1}{2} \mathbb{E}[y_{ij}^* - f_{ij}]^2 \right) + \sum_{i=1}^n \log C_i \end{aligned}$$

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

$$\begin{aligned} \mathbb{E} [\log p(\mathbf{w})] - \mathbb{E} [\log q(\mathbf{w})] &= \sum_{j=1}^m \left( \mathbb{E} [\log p(\mathbf{w}_j)] - \mathbb{E} [\log q(\mathbf{w}_j)] \right) \\ &= \sum_{j=1}^m \left( -\frac{n}{2} \log 2\pi - \frac{1}{2} \mathbb{E}[\mathbf{w}_j^\top \mathbf{w}_j] + \mathcal{H}(q(\mathbf{w}_j)) \right) \\ &= \sum_{j=1}^m \left( -\frac{n}{2} \log 2\pi - \frac{1}{2} \text{tr}(\mathbb{E}[\mathbf{w}_j \mathbf{w}_j^\top]) + \frac{n}{2} (1 + \log 2\pi) - \frac{1}{2} \log |\mathbf{A}_j| \right) \\ &= \frac{nm}{2} - \frac{1}{2} \sum_{j=1}^m \left( \text{tr} \widetilde{\mathbf{W}}_j + \log |\mathbf{A}_j| \right) \end{aligned}$$

### 3.4.3 Terms involving distribution of $q(\lambda)$

$$\begin{aligned} -\mathbb{E} [\log q(\lambda)] &= \sum_{j=1}^m \mathcal{H}(q(\lambda_j)) \\ &= \sum_{j=1}^m \left( \frac{1}{2} (1 + \log 2\pi) - \frac{1}{2} \log c_j \right) \end{aligned}$$

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

or if using single  $\lambda$

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

#### 3.4.4 Terms involving distribution of $q(\alpha)$

$$\begin{aligned} - \mathbb{E} [\log q(\alpha)] &= \sum_{j=1}^m \mathcal{H}(q(\alpha_j)) \\ &= \frac{m}{2}(1 + \log 2\pi - \log n) \end{aligned}$$

or if using single  $\alpha$

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

#### 3.4.5 The lower bound

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

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

### 3.5 Prediction

For a new data point  $x_{\text{new}}$ , we calculate the predicted latent values  $\tilde{f}_{\text{new}} = (\tilde{f}_{\text{new},1}, \dots, \tilde{f}_{\text{new},m})$  for each of the classes, using the variational estimates of the posterior means for the unknown quantities (denoted with tildes), as follows:

$$\tilde{f}_{\text{new},j} = \tilde{\alpha}_j + \sum_{k=1}^n h_{\tilde{\lambda}_j}(x_{\text{new}}, x_k) \tilde{w}_{kj}, \quad j = 1, \dots, m.$$

The predicted class is equal to

$$y_{\text{new}} = \arg \max_j \tilde{f}_{\text{new},j}.$$

To get the fitted probabilities for each class, the following integrals needs to be computed:

$$\begin{aligned}\tilde{p}_{\text{new},j} &= \mathbb{E}_Z \left[ \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( Z + \tilde{f}_{\text{new},j} - \tilde{f}_{\text{new},k} \right) \right] \\ &= \int \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( z + \tilde{f}_{\text{new},j} - \tilde{f}_{\text{new},k} \right) \phi(z) \, dz\end{aligned}$$

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

### 3.6 The variational Bayes EM algorithm

Since there is a cyclic dependence of the parameters on each other, we employ a sequential update algorithm. In what follows, a tilde on the parameters indicate that these are the expectations of the parameters given the optimal factorised distributions  $\tilde{q}$  derived earlier.

6. This isn't the variational EM... CAVI?

- STEP 1: Update  $\tilde{\mathbf{y}}^{*(t+1)}$  given  $\tilde{\mathbf{w}}^{(t)}$ ,  $\tilde{\lambda}^{(t)}$ , and  $\tilde{\alpha}^{(t)}$
- STEP 2: Update  $\tilde{\mathbf{w}}^{(t+1)}$  given  $\tilde{\mathbf{y}}^{*(t+1)}$ ,  $\tilde{\lambda}^{(t)}$ , and  $\tilde{\alpha}^{(t)}$
- STEP 3: Update  $\tilde{\lambda}^{(t+1)}$  given  $\tilde{\mathbf{y}}^{*(t+1)}$ ,  $\tilde{\mathbf{w}}^{(t+1)}$ , and  $\tilde{\alpha}^{(t)}$
- STEP 4: Update  $\tilde{\alpha}^{(t+1)}$  given  $\tilde{\mathbf{y}}^{*(t+1)}$ ,  $\tilde{\mathbf{w}}^{(t+1)}$ , and  $\tilde{\lambda}^{(t+1)}$

---

#### Algorithm 1 VB-EM algorithm for the probit I-prior model

---

```

1: procedure INITIALISE
2:   for  $j = 1, \dots, m$  do
3:      $\tilde{\mathbf{w}}_j^{(0)} \leftarrow \mathbf{0}_n$ 
4:      $\tilde{\alpha}_j^{(0)} \leftarrow \mathcal{N}(0, 1)$ 
5:      $\tilde{\lambda}_j^{(0)} \leftarrow \mathcal{N}(0, 1)$ 
6:      $\tilde{\lambda}_j^{sq(0)} \leftarrow (\tilde{\lambda}_j^{(0)})^2 \quad \triangleright \text{this is } \mathbb{E}[\lambda_j^2]$ 
7:      $\mathbf{H}_{\lambda_j}^{(0)} \leftarrow \tilde{\lambda}_j^{(0)} \mathbf{H}$ 
8:      $\mathbf{H}_{\lambda_j}^{sq(0)} \leftarrow \tilde{\lambda}_j^{sq(0)} \mathbf{H}^2$ 
9:   end for
10: end procedure

11: procedure UPDATE FOR  $\tilde{\mathbf{f}}$  (time  $t$ )
12:   for  $j = 1, \dots, m$  do
13:      $\tilde{\mathbf{f}}_j^{(t+1)} \leftarrow \tilde{\alpha}_j^{(t)} \mathbf{1}_n + \mathbf{H}_{\lambda_j} \tilde{\mathbf{w}}_j^{(t)}$ 
14:   end for
15:    $\tilde{\mathbf{f}}^{(t+1)} \leftarrow (\tilde{\mathbf{f}}_1^{(t+1)}, \dots, \tilde{\mathbf{f}}_m^{(t+1)})^\top$ 
16: end procedure
```

---

---

```

17: procedure UPDATE FOR  $y_{ij}^*$  (time  $t$ )
18:   for  $i = 1, \dots, n$  do
19:      $j \leftarrow y_i$ 
20:      $C_i^{(t+1)} \leftarrow \prod_{k \neq j} \Phi \left( (\tilde{f}_{ij}^{(t+1)} - \tilde{f}_{ik}^{(t+1)}) / \sqrt{2} \right)$ 
21:     for  $k = 1, \dots, j-1, j+1, \dots, m$  do
22:        $D_{ik} \leftarrow \mathbb{E}_Z \left[ \phi_k(Z + \tilde{f}_{ij}^{(t+1)} - \tilde{f}_{ik}^{(t+1)}) \prod_{l \neq k, j} \Phi_l(Z + \tilde{f}_{ij}^{(t+1)} - \tilde{f}_{ik}^{(t+1)}) \right]$ 
23:        $\tilde{y}_{ik}^{*(t+1)} \leftarrow \tilde{f}_{ik}^{(t+1)} - D_{ik} / C_i^{(t+1)}$ 
24:     end for
25:      $\tilde{y}_{ij}^{*(t+1)} \leftarrow \tilde{f}_{ij}^{(t+1)} - \sum_{k \neq j} (\tilde{y}_{ik}^{*(t+1)} - \tilde{f}_{ik}^{(t+1)})$ 
26:   end for
27: end procedure

28: procedure UPDATE FOR  $\mathbf{w}_j$  (time  $t$ )
29:   for  $j = 1, \dots, m$  do
30:      $\tilde{\mathbf{y}}_j^{*(t+1)} \leftarrow (\tilde{y}_{1j}^{(t+1)}, \dots, \tilde{y}_{nj}^{(t+1)})^\top$ 
31:      $\mathbf{A}_j \leftarrow \mathbf{H}_{\lambda_j}^{sq(t)} + \mathbf{I}_n$ 
32:      $\mathbf{a}_j \leftarrow \mathbf{H}_\lambda(\tilde{\mathbf{y}}_j^{*(t+1)} - \tilde{\alpha}_j^{(t)} \mathbf{1}_n)$ 
33:      $\tilde{\mathbf{w}}_j^{(t+1)} \leftarrow \mathbf{A}_j^{-1} \mathbf{a}_j$ 
34:      $\tilde{\mathbf{W}}_j^{(t+1)} \leftarrow \mathbf{A}_j^{-1} (\mathbf{I}_n + \mathbf{a}_j \mathbf{a}_j^\top \mathbf{A}_j^{-1})$ 
35:      $\log \det \mathbf{A}_j^{(t+1)} \leftarrow \log |\mathbf{A}_j|$ 
36:   end for
37: end procedure

38: procedure UPDATE FOR  $\lambda$  (time  $t$ )
39:   for  $j = 1, \dots, m$  do
40:      $c_j^{(t+1)} \leftarrow \text{tr} \left( \mathbf{H}^2 \tilde{\mathbf{W}}_j \right)$ 
41:      $d_j \leftarrow (\tilde{\mathbf{y}}_j^{*(t+1)} - \tilde{\alpha}_j^{(t)} \mathbf{1}_n)^\top \mathbf{H} \tilde{\mathbf{w}}_j^{(t+1)}$ 
42:      $\tilde{\lambda}_j^{(t+1)} \leftarrow d_j / c_j^{(t+1)}$ 
43:      $\tilde{\lambda}_j^{sq(t+1)} \leftarrow 1 / c_j^{(t)} + (d_j / c_j^{(t+1)})^2$ 
44:   end for
45:   if single  $\lambda$  then  $\forall j$ 
46:      $\tilde{\lambda}_j^{(t+1)} \leftarrow \sum_j d_j / \sum_j c_j^{(t+1)}$ 
47:      $\tilde{\lambda}_j^{sq(t+1)} \leftarrow 1 / \sum_j c_j^{(t+1)} + \left( \sum_j d_j / \sum_j c_j^{(t+1)} \right)^2$ 
48:   end if
49:   call UPDATE KERNEL MATRICES
50: end procedure

```

---



---

```

51: procedure UPDATE KERNEL MATRICES (time  $t$ )
52:   for  $j = 1, \dots, m$  do
53:      $\mathbf{H}_{\lambda_j}^{(t+1)} \leftarrow \tilde{\lambda}_j^{(t+1)} \mathbf{H}$ 
54:      $\mathbf{H}_{\lambda_j}^{sq(t+1)} \leftarrow \tilde{\lambda}_j^{sq(t+1)} \mathbf{H}^2$ 
55:   end for
56: end procedure

57: procedure UPDATE FOR  $\alpha$  (time  $t$ )
58:   if single  $\alpha$  then
59:      $\tilde{\alpha}^{(t+1)} \leftarrow \frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n (\tilde{y}_{ij}^{*(t+1)} - \tilde{\lambda}_j^{(t+1)} \mathbf{H}_i \tilde{\mathbf{w}}_j^{(t+1)})$ 
60:   else
61:     for  $j = 1, \dots, m$  do
62:        $\tilde{\alpha}_j^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n (\tilde{y}_{ij}^{*(t+1)} - \tilde{\lambda}_j^{(t+1)} \mathbf{H}_i \tilde{\mathbf{w}}_j^{(t+1)})$ 
63:     end for
64:   end if
65: end procedure

66: procedure CALCULATE LOWER BOUND (time  $t$ )
67:    $\mathcal{L}^{(t)} \leftarrow \frac{1}{2} (nm - \log nm + 3(1 + \log 2\pi)) - \frac{1}{2} \left( \log \det \mathbf{A}^{(t)} + \text{tr } \widetilde{\mathbf{W}}^{(t)} + \sum_{i=1}^2 \log c_i^{(t)} \right) +$ 
68:      $\sum_{i=1}^n \log C_i^{(t)}$ 
69: end procedure

69: procedure THE VB-EM ALGORITHM
70:    $t \leftarrow 0$ 
71:   while  $\mathcal{L}^{(t+1)} - \mathcal{L}^{(t)} > \delta$  or  $t < t_{max}$  do
72:     call UPDATE FOR  $\mathbf{y}^*$ 
73:     call UPDATE FOR  $\mathbf{w}$ 
74:     call UPDATE FOR  $\lambda$ 
75:     call UPDATE FOR  $\alpha$ 
76:     call CALCULATE LOWER BOUND
77:      $t \leftarrow t + 1$ 
78:   end while
79: end procedure

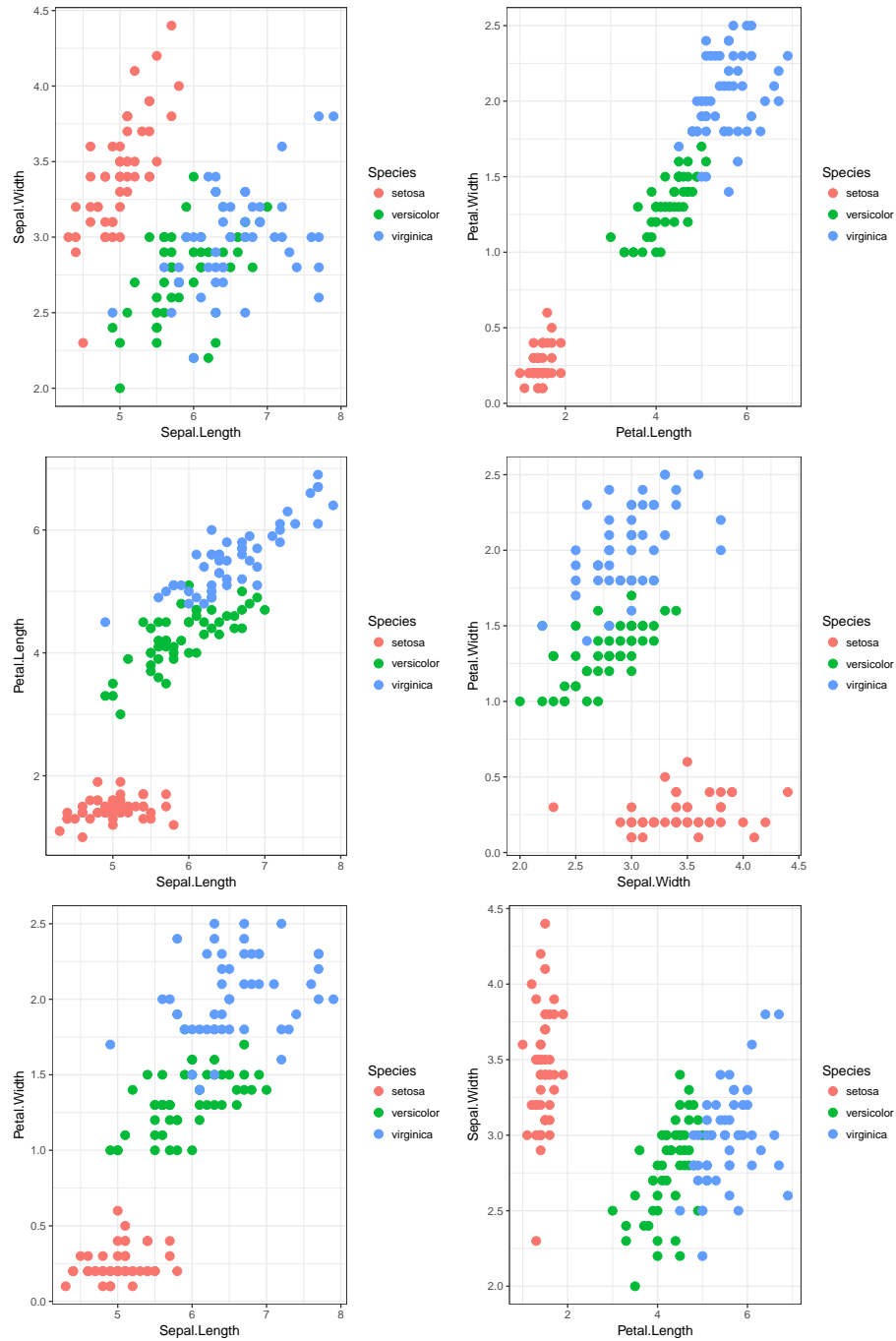
80: return  $(\hat{\mathbf{y}}^*, \hat{\mathbf{w}}, \hat{\lambda}, \hat{\alpha}) \leftarrow (\tilde{\mathbf{y}}^{*(t)}, \tilde{\mathbf{w}}^{(t)}, \tilde{\lambda}^{(t)}, \tilde{\alpha}^{(t)})$   $\triangleright$  converged parameter estimates
81: return  $(\hat{y}_1, \dots, \hat{y}_n) \leftarrow \left( \arg \max_{k=1}^m \hat{y}_{1k}^*, \dots, \arg \max_{k=1}^m \hat{y}_{nk}^* \right)$   $\triangleright$  predicted classes
82: for  $i = 1, \dots, n$  do
83:   for  $j = 1, \dots, m$  do
84:     return  $\hat{p}_{ij} \leftarrow \prod_{\substack{k=1 \\ k \neq j}}^m \Phi \left( \frac{\tilde{y}_{ij}^* - \tilde{y}_{ik}^*}{\sqrt{2}} \right)$   $\triangleright$  predicted probabilities
85:   end for
86: end for

```

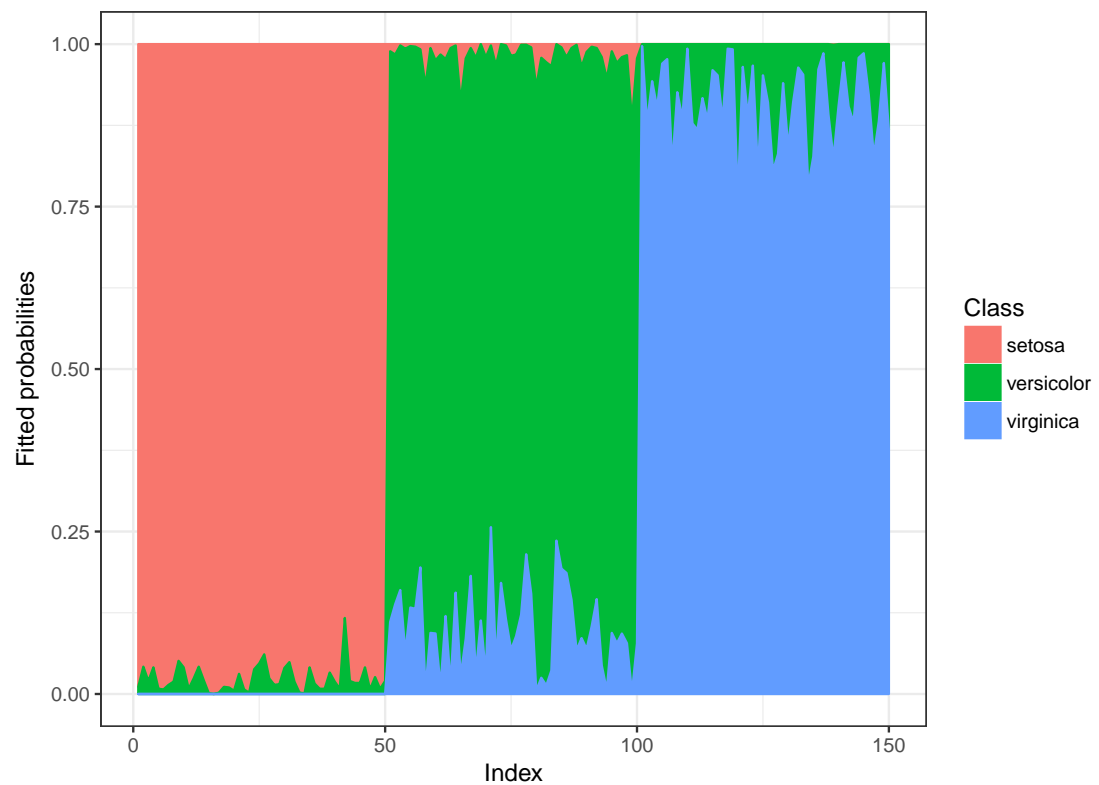
---

## 4 Examples

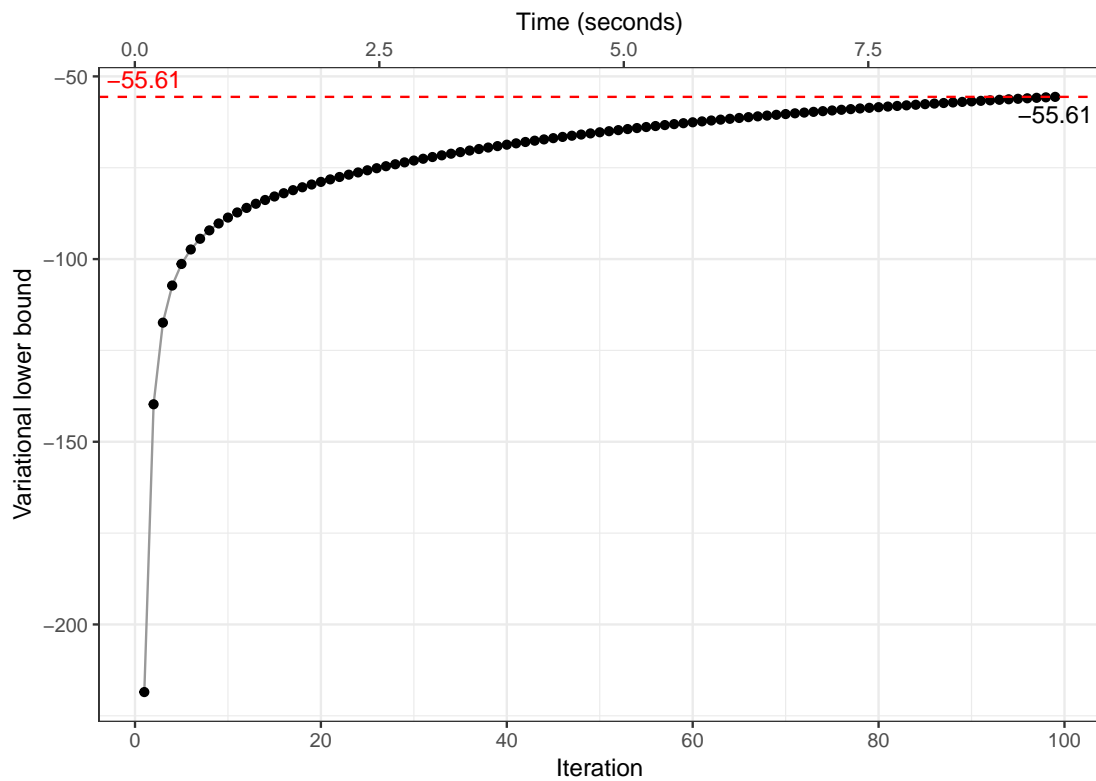
### 4.1 Iris data set



```
R> (mod <- iprobit_mult(y, X, silent = TRUE))
## Lower bound value = -55.60715
## Iterations = 100
##
##          Class = 1 Class = 2 Class = 3
## Intercept -0.21708  1.20009 -1.13551
## lambda    -0.35226 -0.35226 -0.35226
R> plot(mod)
```



```
R> iplot_lb(mod)
```

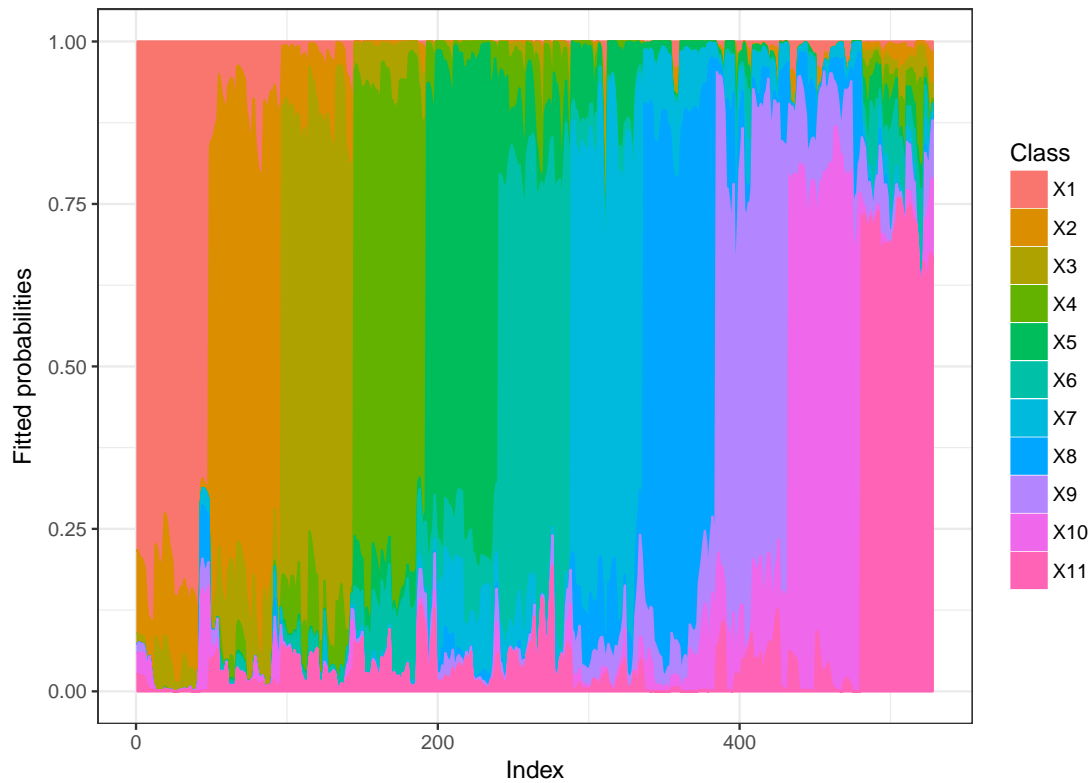


## 4.2 Vowel recognition data set

```
## class x.1 x.2 x.3 x.4 x.5 x.6 x.7 x.8 x.9 x.10
## 1 1 -3.639 0.418 -0.670 1.779 -0.168 1.627 -0.388 0.529 -0.874 -0.814
## 2 2 -3.327 0.496 -0.694 1.365 -0.265 1.933 -0.363 0.510 -0.621 -0.488
## 3 3 -2.120 0.894 -1.576 0.147 -0.707 1.559 -0.579 0.676 -0.809 -0.049
## 4 4 -2.287 1.809 -1.498 1.012 -1.053 1.060 -0.567 0.235 -0.091 -0.795
## 5 5 -2.598 1.938 -0.846 1.062 -1.633 0.764 0.394 -0.150 0.277 -0.396
## 6 6 -2.852 1.914 -0.755 0.825 -1.588 0.855 0.217 -0.246 0.238 -0.365
```

```
R> set.seed(123)
R> (mod <- iprobit_mult(vow.tr$class, vow.tr[, -1], kernel = "FBM", silent = TRUE))

## Lower bound value = -736.8918
## Iterations = 100
##
##          Class = 1 Class = 2 Class = 3 Class = 4 Class = 5 Class = 6
## Intercept -0.11514  0.13838  0.04304  0.07129  0.21767  0.46536
## lambda    -0.13430 -0.13430 -0.13430 -0.13430 -0.13430 -0.13430
##          Class = 7 Class = 8 Class = 9 Class = 10 Class = 11
## Intercept  0.40117 -0.3387  0.47458 -0.06605  0.67874
## lambda    -0.13430 -0.1343 -0.13430 -0.13430 -0.13430
R> plot(mod)
```



```
R> predict(mod, X.test = vow.ts[, -1], y.test = vow.ts[, 1])
## Test error rate: 41 %
```

	Error rates	
	Training	Test
k-Nearest neighbours	NA	44
Linear regression	48	67
Linear discriminant analysis	32	56
Neural network	NA	45
FDA/BRUTO	6	44
FDA/MARS	13	39
I-probit (FBM-0.5)	0	41

## 5 Discussion

### 5.1 Estimating a scaled probit model

For future work, can consider estimating the covariances/correlations across choices. More suitable for social science data.

## 5.2 Reduction to binary models

Describe the model when there are only two alternatives.

## 5.3 Multiple scale parameters with multiple kernels

It's definitely possible to extend to multiple scale parameters. It's just a matter of algebra.

## 5.4 Modelling ordinal data

A different model using thresholds, but it might be possible to model these and estimate using variational inference.

# A Proofs

## A.1 Proof of Lemma 1

*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

$$\begin{aligned}
E[x_i] &= C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_i \prod_{k=1}^d \frac{1}{\sigma_k} \phi\left(\frac{x_k - \mu_k}{\sigma_k}\right) dx_1 \cdots dx_d \\
&= C^{-1} \iint \mathbb{1}[x_i < x_j] \frac{x_i}{\sigma_i} \phi\left(\frac{x_i - \mu_i}{\sigma_i}\right) \prod_{k \neq i, j} \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) dx_i dx_j \\
&= C^{-1} \iint \mathbb{1}[\sigma_i z_i + \mu_i < \sigma_j z_j + \mu_j] (\sigma_i z_i + \mu_i) \phi(z_i) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) dz_i dz_j \\
&= \mu_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i] \phi(z_i) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) dz_i dz_j \\
&\quad + \sigma_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i] z_i \phi(z_i) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) dz_i dz_j \\
&= \mu_i C^{-1} \int \overbrace{\prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right)}^C \phi(z_j) dz_j \\
&\quad + \sigma_i C^{-1} \int \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i)/\sigma_i] z_i \phi(z_i) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) dz_i dz_j
\end{aligned}$$

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,

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

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

$$\begin{aligned} \mathbb{E}[X_i] &= \mu_i - \sigma_i C^{-1} \int \phi\left(\frac{\sigma_j z_j + \mu_j - \mu_i}{\sigma_i}\right) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) dz_j \\ &= \mu_i - \sigma_i C^{-1} \mathbb{E} \left[ \phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_i}{\sigma_i}\right) \prod_{k \neq i, j} \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right) \right]. \end{aligned}$$

The expectation for the  $j$ th component is

$$\begin{aligned} \mathbb{E}[X_j] &= C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_j \prod_{k=1}^d \frac{1}{\sigma_k} \phi\left(\frac{x_k - \mu_k}{\sigma_k}\right) dx_1 \cdots dx_d \\ &= C^{-1} \int x_j \prod_{k \neq j} \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) \cdot \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) dx_j \\ &= C^{-1} \int (\sigma_j z_j + \mu_j) \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot \phi(z_j) dz_j \\ &= \mu_j C^{-1} \int \overbrace{\prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right)}^C \cdot \phi(z_j) dz_j \\ &\quad + \sigma_j C^{-1} \int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot z_j \phi(z_j) dz_j \\ &= \mu_j + \sigma_j C^{-1} \mathbb{E} \left[ Z_j \prod_{k \neq j} \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right) \right] \\ &= \mu_j + \sigma_j \sum_{\substack{i=1 \\ i \neq j}}^d \sigma_i C^{-1} \mathbb{E} \left[ \phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right) \right] \\ &= \mu_j - \sigma_j \sum_{i \neq j} (\mathbb{E}[X_i] - \mu_i) \end{aligned}$$

where we have made use of Lemma 3 in the second last step of the above.

(ii) The differential entropy is given by

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

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

□

**Lemma 3.** Let  $Z \sim \mathcal{N}(0, 1)$ . Then for all  $m \in \{\mathbb{N} \mid m > 1\}$  and  $(\mu, \sigma) \in \mathbb{R} \times \mathbb{R}^+$ ,

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

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

*Proof.* Use the fact that for any differentiable function  $g$ ,  $\mathbb{E}[Zg(Z)] = \mathbb{E}[g'(Z)]$ , 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:

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

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

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

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

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

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.  $\square$

## References

- Albert, J. and S. Chib (1993). “Bayesian Analysis of Binary and Polychotomous Response Data”. In: *Journal of the American Statistical Association* 88.422, pp. 669–679. ISSN: 0162-1459. DOI: 10.2307/2290350.
- Bolduc, D. (1999). “A practical technique to estimate multinomial probit models in transportation”. In: *Transportation Research Part B: Methodological* 33.1, pp. 63–79. ISSN: 01912615. DOI: 10.1016/S0191-2615(98)00028-9.
- Dow, J. K. and J. W. Endersby (2004). “Multinomial probit and multinomial logit: A comparison of choice models for voting research”. In: *Electoral Studies* 23.1, pp. 107–122. ISSN: 02613794. DOI: 10.1016/S0261-3794(03)00040-4.
- Girolami, M. and S. Rogers (2006). “Variational Bayesian Multinomial Probit Regression with Gaussian Process Priors”. In: *Neural Computation* 18.8, pp. 1790–1817. ISSN: 0899-7667. DOI: 10.1093/bioinformatics/btm535.
- Natarajan, R., C. E. McCulloch, and N. M. Kiefer (1995). “Maximum likelihood for the multinomial probit model”. In: *Unpublished manuscript*.
- Rogers, S. and M. Girolami (2007). “Multi-class Semi-supervised Learning With The  $\epsilon$ -truncated Multinomial Probit Gaussian Process”. In: *JLMR: Workshop and Conference Proceedings* 2006, pp. 17–32.
- Sheffi, Y., R. Hall, and C. Daganzo (1982). “On the estimation of the multinomial probit model”. In: *Transportation Research Part A: General* 16.5-6, pp. 447–456. ISSN: 01912607. DOI: 10.1016/0191-2607(82)90071-1.