_O-	do list	
Is thi	is variational EM or CAVI?	11
Coi	ntents	
Esti 6.1 6.2 6.3 6.4	mation of I-probit models using variational inference Relevant distributions	1 2 3 9 11
A.1 A.2	Proof of Lemma 6.1	15 15 18 21
Bibliography		21
List of Figures		22
st of	Tables	23
st of	Theorems	24
st of	Definitions	25
st of	Symbols	26
	Esti 6.1 6.2 6.3 6.4 App A.1 A.2 A.3 ibliog st of st of st of	6.2 Mean field distributions

Haziq Jamil
Department of Statistics
London School of Economics and Political Science
August 16, 2017

## Chapter 6

# Estimation of I-probit models using variational inference

In this chapter we provide the details of the variational algorithm to estimate cate-

gorical I-prior models.

### 6.1 Relevant distributions

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

$$p(\mathbf{y}^*|\mathbf{f}) = \prod_{i=1}^n \prod_{j=1}^m N(f_{ij}, 1)$$

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

$$= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} ||\mathbf{y}^* - \mathbf{f}||^2 \right]$$

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

$$p(\mathbf{w}) = \prod_{i=1}^{n} \prod_{j=1}^{m} p(w_{ij})$$

$$= \prod_{i=1}^{n} \prod_{j=1}^{m} N(0, 1)$$

$$= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{m} w_{ij}^{2} \right]$$

$$= \exp \left[ -\frac{nm}{2} \log 2\pi - \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} \right]$$

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

#### 6.2 Mean field distributions

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

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

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

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

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

$$\log \tilde{q}(y_i^*) = \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \cdot \mathbf{E}_{\mathbf{w},\alpha,\lambda} \left[ -\frac{1}{2} \sum_{k=1}^m (y_{ik}^* - f_{ik})^2 \right] + \text{const.}$$

$$= \mathbb{1}[y_{ij}^* = \max_k y_{ik}^*] \cdot \left[ -\frac{1}{2} \sum_{k=1}^m (y_{ik}^* - \tilde{f}_{ik})^2 \right] + \text{const.}$$

$$\equiv \begin{cases} \prod_{k=1}^m \mathbf{N}(\tilde{f}_{ik}, 1) & \text{if } y_{ij}^* > y_{ik}^*, \forall k \neq j \\ 0 & \text{otherwise} \end{cases}$$

where  $\tilde{f}_{ik} = \mathrm{E}[\alpha_k] + \sum_{l=1}^m h_{\mathrm{E}[\lambda_k]}(x_i, x_l) \, \mathrm{E}[w_{il}]$ , 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 jth component is always largest. It is worth investigating the properties of this distribution, and we now present some relevant definitions and results.

**Definition 6.1** (Conically-truncated multivariate normal distribution). Let  $\mathbf{X} = (X_1, \ldots, X_d)$  be a d-dimensional random variable with pdf defined as

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

for some  $j \in \{1, ..., d\}$ . We denote the distribution of **X** by  $N^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , with  $\boldsymbol{\mu} = (\mu_1, ..., \mu_d)$  and  $\boldsymbol{\Sigma} = \operatorname{diag}(\sigma_1^2, ..., \sigma_d^2)$ . The pdf of **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

$$C = \mathbf{E}_{Z} \left[ \prod_{\substack{i=1\\i\neq j}}^{d} \Phi \left( \frac{\sigma_{j}}{\sigma_{i}} Z + \frac{\mu_{j} - \mu_{i}}{\sigma_{i}} \right) \right]$$

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

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

*Proof.* A derivation of the functional form for the pdf of  $X \sim N^{(j)}(\mu, \Sigma)$  is given. Using

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

$$\int \mathbb{1}[x_{i} < x_{j}, \forall i \neq j] \prod_{i=1}^{d} \mathcal{N}(\mu_{i}, \sigma_{i}^{2}) \, \mathrm{d}x_{1} \cdots \mathrm{d}x_{d}$$

$$= \int \mathbb{1}[x_{i} < x_{j}, \forall i \neq j] \prod_{i=1}^{d} \left[\frac{1}{\sigma_{i}} \phi\left(\frac{x_{i} - \mu_{i}}{\sigma}\right)\right] \, \mathrm{d}x_{1} \cdots \mathrm{d}x_{d}$$

$$= \int \mathbb{1}[x_{i} < x_{j}, \forall i \neq j] \frac{1}{\sigma_{j}} \phi\left(\frac{x_{j} - \mu_{j}}{\sigma_{j}}\right) \prod_{\substack{i=1\\i\neq j}}^{d} \left[\frac{1}{\sigma_{i}} \phi\left(\frac{x_{i} - \mu_{i}}{\sigma_{i}}\right)\right] \, \mathrm{d}x_{1} \cdots \, \mathrm{d}x_{d}$$

$$= \int \prod_{\substack{i=1\\i\neq j}}^{d} \Phi\left(\frac{x_{j} - \mu_{i}}{\sigma_{i}}\right) \frac{1}{\sigma_{j}} \phi\left(\frac{x_{j} - \mu_{j}}{\sigma_{j}}\right) \, \mathrm{d}x_{j}$$

$$= \int \prod_{\substack{i=1\\i\neq j}}^{d} \Phi\left(\frac{\sigma_{j}z_{j} + \mu_{j} - \mu_{i}}{\sigma_{i}}\right) \phi(z_{j}) \, \mathrm{d}z_{j}$$
(by using the standardisation  $z_{j} = (x_{j} - \mu_{j})/\sigma_{j}$ )
$$= \mathbb{E}\left[\prod_{\substack{i=1\\i\neq j}}^{d} \Phi\left(\frac{\sigma_{j}}{\sigma_{i}}Z_{j} + \frac{\mu_{j} - \mu_{i}}{\sigma_{i}}\right)\right]$$

the proof follows directly.

**Lemma 6.1.** Let  $X \sim N^{(j)}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with pdf  $p(\mathbf{x})$  as defined in Definition 6.1. Then

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

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

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

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

where  $C = \mathbb{E}\left[\prod_{i \neq j} \Phi_i\right]$ , 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  $Z \sim N(0,1)$ , 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  $\{1, \ldots, m\}$ . Thus, we have that the distribution of  $(y_{i1}^*, \ldots, y_{im}^*)$  is  $N^{(y_i)}(\boldsymbol{\mu}_i, \mathbf{I}_m)$ , where  $\boldsymbol{\mu}_i = (\tilde{f}_{i1}, \ldots, \tilde{f}_{im})$ . The expectation is given by

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

where

$$\phi_{ik}(Z) = \phi(Z + \tilde{f}_{iy_i} - \tilde{f}_{ik})$$

$$\Phi_{ik}(Z) = \Phi(Z + \tilde{f}_{iy_i} - \tilde{f}_{ik})$$

$$C_i = E_Z \left[ \prod_{\substack{i=1\\i\neq j}}^d \Phi\left(Z + \tilde{f}_{iy_i} - \tilde{f}_{ik}\right) \right]$$

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

$$E_{Z}\left[\phi_{ik}(Z)\prod_{l\neq k,j}\Phi_{il}(Z)\right] = \int \phi_{ik}(z)\prod_{l\neq k,j}\Phi_{il}(z)\phi(z)\,dz, \quad \forall k\neq y_{i}$$

$$C_{i} = E_{Z}\left[\prod_{l\neq j}\Phi_{il}(Z)\right] = \int \prod_{l\neq j}\Phi_{il}(z)\phi(z)\,dz$$

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

 $|\tilde{q}(\mathbf{w})|$ 

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

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

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

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

Let  $\mathbf{A}_j = \mathrm{E}[\mathbf{H}_{\lambda_i}^2] + \mathbf{I}_n$  and  $\mathbf{a}_j = \mathrm{E}[\mathbf{H}_{\lambda_j}](\mathrm{E}[\mathbf{y}_j^*] - \mathrm{E}[\alpha_j]\mathbf{1}_n)$ . 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  $j = 1, \ldots, m$ ,

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

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  $\mathrm{E}[\mathbf{w}_{j}\mathbf{w}_{j}^{\top}] = \mathbf{A}_{j}^{-1}(\mathbf{I}_{n} + \mathbf{a}_{j}\mathbf{a}_{j}^{\top}\mathbf{A}_{j}^{-1}) =: \widetilde{\mathbf{W}}_{j}$ .

 $\tilde{q}(\lambda)$ 

For j = 1, ..., m,

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

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

$$= -\frac{1}{2} \sum_{j=1}^m \left( \lambda_j^2 \operatorname{tr} \left( \mathbf{H}^2 \mathbf{E}[\mathbf{w}_j \mathbf{w}_j^\top] \right) - 2\lambda_j (\mathbf{E}[\mathbf{y}_j^*] - \mathbf{E}[\alpha_j] \mathbf{1}_n)^\top \mathbf{H} \mathbf{E}[\mathbf{w}_j] \right) + \text{const.}$$

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

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

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

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

 $\tilde{q}(\alpha)$ 

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

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

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

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

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

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

## 6.3 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

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

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

$$= \mathrm{E}\left[ \log \prod_{i=1}^{n} \prod_{j=1}^{m} p(y_i | y_{ij}^*) \right] + \mathrm{E}\left[\log p(\mathbf{y}^* | \mathbf{f})\right] + \mathrm{E}\left[\log p(\mathbf{w})\right] + \mathrm{E}\left[\log p(\lambda)\right] + \mathrm{E}\left[\log p(\lambda)\right]$$

$$- \mathrm{E}\left[\log q(\mathbf{y}^*)\right] - \mathrm{E}\left[\log q(\mathbf{w})\right] - \mathrm{E}\left[\log q(\lambda)\right] - \mathrm{E}\left[\log q(\lambda)\right]$$

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

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

**Lemma 6.2.** 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  $\Sigma$ ,

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

Terms involving distributions of  $y^*$ 

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

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

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

Terms involving distributions of w

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

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

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

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

Terms involving distribution of  $q(\lambda)$ 

$$-\operatorname{E}\left[\log q(\lambda)\right] = \sum_{j=1}^{m} \mathcal{H}\left(q(\lambda_{j})\right)$$
$$= \sum_{j=1}^{m} \left(\frac{1}{2}(1 + \log 2\pi) - \frac{1}{2}\log c_{j}\right)$$
$$= \frac{m}{2}(1 + \log 2\pi) - \frac{1}{2}\sum_{j=1}^{m}\log c_{j}$$

or if using single  $\lambda$ 

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

Terms involving distribution of  $q(\alpha)$ 

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

or if using single  $\alpha$ 

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

The lower bound

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

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

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

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

## 6.4 The variational 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.

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

1. Is this variational EM... or CAVI?

```
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
                    for j = 1, \ldots, m do
                              \tilde{\mathbf{w}}_{i}^{(0)} \leftarrow \mathbf{0}_{n}
   3:
                             \begin{split} &\tilde{\alpha}_{j}^{(0)} \leftarrow \mathrm{N}(0,1) \\ &\tilde{\lambda}_{j}^{(0)} \leftarrow \mathrm{N}(0,1) \\ &\tilde{\lambda}_{j}^{sq(0)} \leftarrow (\tilde{\lambda}_{j}^{(0)})^{2} \qquad \triangleright \text{ this is } \mathrm{E}[\lambda_{j}^{2}] \end{split}
   4:
                             \mathbf{H}_{\lambda_{j}}^{(0)} \leftarrow \tilde{\lambda}_{j}^{(0)} \mathbf{H} \\ \mathbf{H}_{\lambda_{j}}^{sq(0)} \leftarrow \tilde{\lambda}_{j}^{sq(0)} \mathbf{H}^{2}
  9:
 10: end procedure
11: procedure UPDATE FOR \hat{\mathbf{f}} (time t)
                    for j = 1, ..., m do
\tilde{\mathbf{f}}_{j}^{(t+1)} \leftarrow \tilde{\alpha}_{j}^{(t)} \mathbf{1}_{n} + \mathbf{H}_{\lambda_{j}} \tilde{\mathbf{w}}_{j}^{(t)}
13:
                    \mathbf{end} \mathbf{for} \\ \tilde{\mathbf{f}}^{(t+1)} \leftarrow \big(\tilde{\mathbf{f}}_1^{(t+1)}, \dots, \tilde{\mathbf{f}}_m^{(t+1)}\big)^\top
14:
 15:
16: end procedure
17: procedure UPDATE FOR y_{ij}^* (time t)
                    for i = 1, \ldots, n do
18:
                             \begin{array}{l} j \leftarrow y_i \\ C_i^{(t+1)} \leftarrow \prod_{k \neq j} \Phi \left( (\tilde{f}_{ij}^{(t+1)} - \tilde{f}_{ik}^{(t+1)}) / \sqrt{2} \right) \end{array}
19:
20:
                              for k = 1, ..., j - 1, j + 1, ..., m do
21:
                                       D_{ik} \leftarrow \mathrm{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]
22:
                                        \tilde{y}_{ik}^{*(t+1)} \leftarrow \tilde{\tilde{f}}_{ik}^{(t+1)} - D_{ik}/C_i^{(t+1)}
23:
                              end for \tilde{y}_{ij}^{*(t+1)} \leftarrow \tilde{f}_{ij}^{(t+1)} - \sum_{k \neq j} \left( \tilde{y}_{ik}^{*(t+1)} - \tilde{f}_{ik}^{(t+1)} \right)
24:
25:
                    end for
26:
27: end procedure
```

```
28: procedure Update for \mathbf{w}_i (time t)
                         \begin{aligned} \mathbf{for} \ j &= 1, \dots, m \ \mathbf{do} \\ \tilde{\mathbf{y}}_{j}^{*(t+1)} &\leftarrow (\tilde{y}_{1j}^{(t+1)}, \dots, \tilde{y}_{nj}^{(t+1)})^{\top} \\ \mathbf{A}_{j} &\leftarrow \mathbf{H}_{\lambda_{j}}^{sq(t)} + \mathbf{I}_{n} \end{aligned}
30:
31:
                                    \mathbf{a}_{j} \leftarrow \mathbf{H}_{\lambda}(\tilde{\mathbf{y}}_{j}^{*(t+1)} - \tilde{\alpha}_{j}^{(t)} \mathbf{1}_{n})
\tilde{\mathbf{w}}_{j}^{(t+1)} \leftarrow \mathbf{A}_{j}^{-1} \mathbf{a}_{j}
\widetilde{\mathbf{W}}_{j}^{(t+1)} \leftarrow \mathbf{A}_{j}^{-1} (\mathbf{I}_{n} + \mathbf{a}_{j} \mathbf{a}_{j}^{\top} \mathbf{A}_{j}^{-1})
\operatorname{logdet}_{j} \mathbf{A}_{j}^{(t+1)} \leftarrow \operatorname{log}|\mathbf{A}_{j}|
33:
34:
35:
                          end for
 36:
37: end procedure
38: procedure UPDATE FOR \lambda (time t)
                        for j = 1, ..., m do
c_j^{(t+1)} \leftarrow \operatorname{tr}\left(\mathbf{H}^2 \widetilde{\mathbf{W}}_j\right)
d_j \leftarrow (\widetilde{\mathbf{y}}_j^{*(t+1)} - \widetilde{\alpha}_j^{(t)} \mathbf{1}_n)^{\top} \mathbf{H} \widetilde{\mathbf{w}}_j^{(t+1)}
\widetilde{\lambda}_j^{(t+1)} \leftarrow d_j / c_j^{(t+1)}
\widetilde{\lambda}_j^{sq(t+1)} \leftarrow 1 / c_j^{(t)} + (d_i / c_i^{(t+1)})^2
 40:
 41:
 42:
 43:
                          end for
 44:
                         if single \lambda then \forall j
\tilde{\lambda}_{j}^{(t+1)} \leftarrow \sum_{j} d_{j} / \sum_{j} c_{j}^{(t+1)}
 45:
 46:
                                      \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}
 47:
                          end if
 48:
                          call Update Kernel Matrices
50: end procedure
51: procedure UPDATE KERNEL MATRICES (time t)
                         \begin{aligned} & \mathbf{for} \ j = 1, \dots, m \ \mathbf{do} \\ & \mathbf{H}_{\lambda_j}^{(t+1)} \leftarrow \tilde{\lambda}_j^{(t+1)} \mathbf{H} \\ & \mathbf{H}_{\lambda_j}^{sq(t+1)} \leftarrow \tilde{\lambda}_j^{sq(t+1)} \mathbf{H}^2 \end{aligned}
52:
53:
54:
                          end for
56: end procedure
```

```
57: procedure UPDATE FOR \alpha (time t)
                              if single \alpha then
 58:
                                          \tilde{\alpha}^{(t+1)} \leftarrow \frac{1}{nm} \sum_{i=1}^{m} \sum_{i=1}^{n} \left( \tilde{y}_{ij}^{*(t+1)} - \tilde{\lambda}_{j}^{(t+1)} \mathbf{H}_{i} \tilde{\mathbf{w}}_{j}^{(t+1)} \right)
59:
 60:
                              else
                                          for j = 1, ..., m do
\tilde{\alpha}_j^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n \left( \tilde{y}_{ij}^{*(t+1)} - \tilde{\lambda}_j^{(t+1)} \mathbf{H}_i \tilde{\mathbf{w}}_j^{(t+1)} \right)
61:
 62:
 63:
 64:
                              end if
65: end procedure
66: procedure Calculate lower bound (time t)
                            \mathcal{L}^{(t)} \leftarrow \frac{1}{2} \left( nm - \log nm + 3(1 + \log 2\pi) \right) - \frac{1}{2} \left( \operatorname{logdetA}^{(t)} + \operatorname{tr} \widetilde{\mathbf{W}}^{(t)} + \sum_{i=1}^{2} \log c_{i}^{(t)} \right) + C_{i}^{(t)} + C_{i
\sum_{i=1}^{n} \log C_i^{(t)}
68: end procedure
69: procedure The VB-EM ALGORITHM
 70:
                             while \mathcal{L}^{(t+1)} - \mathcal{L}^{(t)} > \delta or t < t_{max} do
 71:
                                           call Update for \mathbf{y}^*
 72:
 73:
                                           call Update for w
                                           call Update for \lambda
 74:
                                           call Update for \alpha
 75:
                                           call Calculate lower bound
 76:
                                           t \leftarrow t + 1
 77:
                             end while
 78:
 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(\underset{k=1}{\operatorname{arg max}} \hat{y}_{1k}^*, \dots, \underset{k=1}{\operatorname{arg max}} \hat{y}_{nk}^*\right)
                                                                                                                                                                                                                                                                    ▷ predicted classes
 82: for i = 1, ..., n do
                              for j = 1, \ldots, m do
 83:
                                          return \hat{p}_{ij} \leftarrow \prod_{\substack{k=1\\k\neq j}}^m \Phi\left(\frac{\hat{y}_{ij}^* - \hat{y}_{ik}^*}{\sqrt{2}}\right) \triangleright predicted probabilities
 84:
                              end for
 85:
 86: end for
```

## Appendix A

## Appendix for I-probit

## A.1 Proof of Lemma 6.1

*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

$$\begin{split} \mathbf{E}[x_i] &= C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_i \prod_{k=1}^d \frac{1}{\sigma_k} \phi\left(\frac{x_k - \mu_k}{\sigma_k}\right) \mathrm{d}x_1 \cdots \mathrm{d}x_d \\ &= C^{-1} \iint \mathbb{1}[x_i < x_j] \frac{x_i}{\sigma_i} \phi\left(\frac{x_i - \mu_i}{\sigma_i}\right) \prod_{k \neq i,j} \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) \mathrm{d}x_i \, \mathrm{d}x_j \\ &= C^{-1} \iint \mathbb{1}[\sigma_i z_i + \mu_i < \sigma_j z_j + \mu_j] (\sigma_i z_i + \mu_i) \phi(z_i) \prod_{k \neq i,j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) \, \mathrm{d}z_i \, \mathrm{d}z_j \\ &= \mu_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i) / \sigma_i] \phi(z_i) \prod_{k \neq i,j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) \, \mathrm{d}z_i \, \mathrm{d}z_j \\ &+ \sigma_i C^{-1} \iint \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i) / \sigma_i] z_i \phi(z_i) \prod_{k \neq i,j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) \, \mathrm{d}z_i \, \mathrm{d}z_j \\ &= \mu_i C^{-1} \int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) \, \mathrm{d}z_j \\ &+ \sigma_i C^{-1} \int \mathbb{1}[z_i < (\sigma_j z_j + \mu_j - \mu_i) / \sigma_i] z_i \phi(z_i) \prod_{k \neq i,j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \phi(z_j) \, \mathrm{d}z_i \, \mathrm{d}z_j \end{split}$$

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,

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

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

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

The expectation for the jth component is

$$\begin{split} \mathbf{E}[X_j] &= C^{-1} \int \cdots \int \mathbb{1}[x_k < x_j, \forall k \neq j] \cdot x_j \prod_{k=1}^d \frac{1}{\sigma_k} \phi\left(\frac{x_k - \mu_k}{\sigma_k}\right) \, \mathrm{d}x_1 \cdots \, \mathrm{d}x_d \\ &= C^{-1} \int x_j \prod_{k \neq j} \Phi\left(\frac{x_j - \mu_k}{\sigma_k}\right) \cdot \frac{1}{\sigma_j} \phi\left(\frac{x_j - \mu_j}{\sigma_j}\right) \, \mathrm{d}x_j \\ &= C^{-1} \int (\sigma_j z_j + \mu_j) \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot \phi(z_j) \, \mathrm{d}z_j \\ &= \mu_j C^{-1} \int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot \phi(z_j) \, \mathrm{d}z_j \\ &+ \sigma_j C^{-1} \int \prod_{k \neq j} \Phi\left(\frac{\sigma_j z_j + \mu_j - \mu_k}{\sigma_k}\right) \cdot z_j \phi(z_j) \, \mathrm{d}z_j \\ &= \mu_j + \sigma_j C^{-1} \, \mathbf{E}\left[Z_j \prod_{k \neq j} \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right)\right] \\ &= \mu_j + \sigma_j \sum_{\substack{i=1 \\ i \neq j}}^d \sigma_i C^{-1} \, \mathbf{E}\left[\phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_i}{\sigma_i}\right) \prod_{\substack{k=1 \\ k \neq i, j}}^d \Phi\left(\frac{\sigma_j Z_j + \mu_j - \mu_k}{\sigma_k}\right)\right] \\ &= \mu_j - \sigma_j \sum_{\substack{i=1 \\ i \neq j}} \left(\mathbf{E}[X_i] - \mu_i\right) \end{split}$$

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

(ii) The differential entropy is given by

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

$$= -\operatorname{E} \left[ -\log C - \frac{d}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^{d} \log \sigma_{i}^{2} - \frac{1}{2} \sum_{i=1}^{d} \left(\frac{x_{i} - \mu_{i}}{\sigma_{i}}\right)^{2} \right]$$

$$= \log C + \frac{d}{2} \log 2\pi + \frac{1}{2} \sum_{i=1}^{d} \log \sigma_{i}^{2} + \frac{1}{2} \sum_{i=1}^{d} \frac{1}{\sigma_{i}^{2}} \operatorname{E}[x_{i} - \mu_{i}]^{2}.$$

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

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

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

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

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  $m \in \{\mathbb{N} \mid m > 1\}$ . 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

$$g'_{m+1}(z) = \sigma_{m+1}\phi_{m+1}g_m(z) + g'_m(z)\Phi_{m+1}$$

$$= \sigma_{m+1}\phi_{m+1} \prod_{\substack{k=1\\k\neq j}}^m \Phi_k + \sum_{\substack{i=1\\i\neq j}}^m \left[\sigma_i\phi_i \prod_{\substack{k=1\\k\neq i,j}}^m \Phi_k\right] \Phi_{m+1}$$

$$= \sigma_{m+1}\phi_{m+1} \prod_{\substack{k=1\\k\neq j,m+1}}^{m+1} \Phi_k + \sum_{\substack{i=1\\i\neq j}}^m \left[\sigma_i\phi_i \prod_{\substack{k=1\\k\neq i,j}}^{m+1} \Phi_k\right]$$

$$= \sum_{\substack{i=1\\i\neq j}}^{m+1} \left[\sigma_i\phi_i \prod_{\substack{k=1\\k\neq i,j}}^{m+1} \Phi_k\right]$$

$$= g'_{m+1}(z).$$

Thus, by induction and linearity of expectations, the proof is complete.

## A.2 Proof for ...

**Lemma A.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  $\Sigma$ ,

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

(iii) p is distribution of the upper-tail of a univariate, one-sided normal distribution

truncated at zero with mean  $\mu$  and variance 1,

$$\mathcal{H}(p) = \frac{1}{2}\log 2\pi + \frac{1}{2}\left(E[x^2] + \mu^2 - 2\mu E[x]\right) + \log \Phi(\mu)$$

(iv) p is distribution of the **lower-tail** of a univariate, one-sided normal distribution truncated at zero with mean  $\mu$  and variance 1,

$$\mathcal{H}(p) = \frac{1}{2} \log 2\pi + \frac{1}{2} (E x^2 + \mu^2 - 2\mu E x) + \log (1 - \Phi(\mu))$$

Proof.

Case (i): 
$$-\log p(x) = \frac{1}{2}\log 2\pi + \frac{1}{2}\log \sigma^2 + \frac{1}{2}(x-\mu)^2$$
. Then

$$\mathcal{H}(p) = E_x \left[ \frac{1}{2} \log 2\pi + \frac{1}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (x - \mu)^2 \right]$$
$$= \frac{1}{2} \log 2\pi + \frac{1}{2} \log \sigma^2 + \frac{1}{2\sigma^2} E(x - \mu)^{2^*\sigma^2}$$
$$= \frac{1}{2} (1 + \log 2\pi) + \frac{1}{2} \log \sigma^2$$

Case (ii):  $-\log p(x) = \frac{d}{2}\log 2\pi + \frac{1}{2}\log |\Sigma| + \frac{1}{2}(x-\mu)^{\top}\Sigma^{-1}(x-\mu)$ . Then

$$\mathcal{H}(p) = \frac{d}{2} \log 2\pi + \frac{1}{2} \log |\Sigma| + \frac{1}{2} \operatorname{E}_x \left[ (x - \mu)^{\top} \Sigma^{-1} (x - \mu) \right]$$

$$= \frac{d}{2} \log 2\pi + \frac{1}{2} \log |\Sigma| + \frac{1}{2} \operatorname{tr} \left( \Sigma^{-1} \operatorname{E}_x \left[ (x - \mu)(x - \mu)^{\top} \right]^{\top} \right)$$

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

For the next two cases, we state the following properties of a truncated normal distribution without proof.

**Lemma A.3.** Let  $x \sim N(\mu, \sigma^2)$  with x lying in the interval (a, b). Then we say that x follows a truncated normal distribution, and

(i) the mean of x (conditional on a < x < b) is

$$E[x] = \mu + \sigma \frac{\phi(\alpha) - \phi(\beta)}{Z},$$

(ii) the variance of x (conditional on a < x < b) is

$$\operatorname{Var}[x] = \sigma^2 \left[ 1 + \frac{\alpha \phi(\alpha) - \beta \phi(\beta)}{Z} - \left( \frac{\phi(\alpha) - \phi(\beta)}{Z} \right)^2 \right], and$$

(iii) the entropy of the pdf of x (conditional on a < x < b) is

$$\mathcal{H} = \frac{1}{2}\log 2\pi + \frac{1}{2}\log \sigma^2 + \log Z + \frac{1}{2} + \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{2Z},$$

where  $\alpha = (a - \mu)/\sigma$ ,  $\beta = (b - \mu)/\sigma$ , and  $Z = \Phi(\beta) - \Phi(\alpha)$ , and  $\phi$  and  $\Phi$  are the pdf and cdf of a standard normal distribution respectively.

In the special case when  $\sigma = 1$  (the case we are interested in), then with some manipulation, one arrives at the following expression for the entropy of the pdf p of a truncated normal distribution:

$$\mathcal{H}(p) = \frac{1}{2}\log 2\pi + \frac{1}{2}\log \sigma^2 + \log Z + \frac{1}{2}\left(1 + \frac{\alpha\phi(\alpha) - \beta\phi(\beta)}{Z}\right)$$

$$= \frac{1}{2}\log 2\pi + \log Z + \frac{1}{2}\left(\operatorname{Var}[x] + \left(\frac{\phi(\alpha) - \phi(\beta)}{Z}\right)^2\right)$$

$$= \frac{1}{2}\log 2\pi + \log Z + \frac{1}{2}\left(\operatorname{E}[x^2] - \operatorname{E}^2[x] + (\operatorname{E}[x] - \mu)^2\right)$$

$$= \frac{1}{2}\log 2\pi + \log Z + \frac{1}{2}\left(\operatorname{E}[x^2] + \mu^2 - 2\mu\operatorname{E}[x]\right)$$

We now continue with the proof.

Case (iii): Using Lemma A.3 with a = 0,  $b = +\infty$ , and  $\sigma = 1$ , we get that  $Z = 1 - \Phi(-\mu) = \Phi(\mu)$ . Therefore, the entropy of p is given by

$$\mathcal{H}(p) = \frac{1}{2} \log 2\pi + \frac{1}{2} \left( \mathbb{E}[x^2] + \mu^2 - 2\mu \, \mathbb{E}[x] \right) + \log \Phi(\mu)$$

Case (iv): Again, using Lemma A.3 with  $a = -\infty$ , b = 0, and  $\sigma = 1$ , we get that  $Z = \Phi(-\mu) = 1 - \Phi(\mu)$ . Therefore, the entropy of p is given by

$$\mathcal{H}(p) = \frac{1}{2}\log 2\pi + \frac{1}{2}\left(\mathbb{E}[x^2] + \mu^2 - 2\mu\,\mathbb{E}[x]\right) + \log\left(1 - \Phi(\mu)\right)$$

## A.3 Distribution of $\tilde{q}(\mathbf{y}^*)$ for binary case

Case:  $y_i = 1$ 

$$\begin{split} \log \tilde{q}(y_i^*) &= \mathbb{1}[y_i^* \geq 0] \cdot \mathbf{E}_{\mathbf{w},\alpha,\lambda} \left[ -\frac{1}{2} (y_i^* - \alpha - \lambda \mathbf{H}_i \mathbf{w})^2 \right] + \text{const.} \\ &= \mathbb{1}[y_i^* \geq 0] \cdot \left[ -\frac{1}{2} \left( y_i^{*2} - 2 \, \mathbf{E}_{\mathbf{w},\alpha,\lambda} [\alpha + \lambda \mathbf{H}_i \mathbf{w}] y_i \right) \right] + \text{const.} \\ &= \mathbb{1}[y_i^* \geq 0] \left[ -\frac{1}{2} (y_i^* - \tilde{\eta}_i)^2 \right] + \text{const.} \\ &\equiv \begin{cases} \mathbf{N}(\tilde{\eta}_i, 1) & \text{if } y_i^* \geq 0 \\ 0 & \text{if } y_i^* < 0 \end{cases} \end{split}$$

where

$$\tilde{\eta}_i = E \alpha + E \lambda \mathbf{H}_i E \mathbf{w}$$

by independence of  $q(\mathbf{w})$ ,  $q(\alpha)$  and  $q(\lambda)$ .  $\tilde{q}(y_i^*)$  is recognised as being the upper-tail of a one-sided normal distribution truncated at zero. The mean is

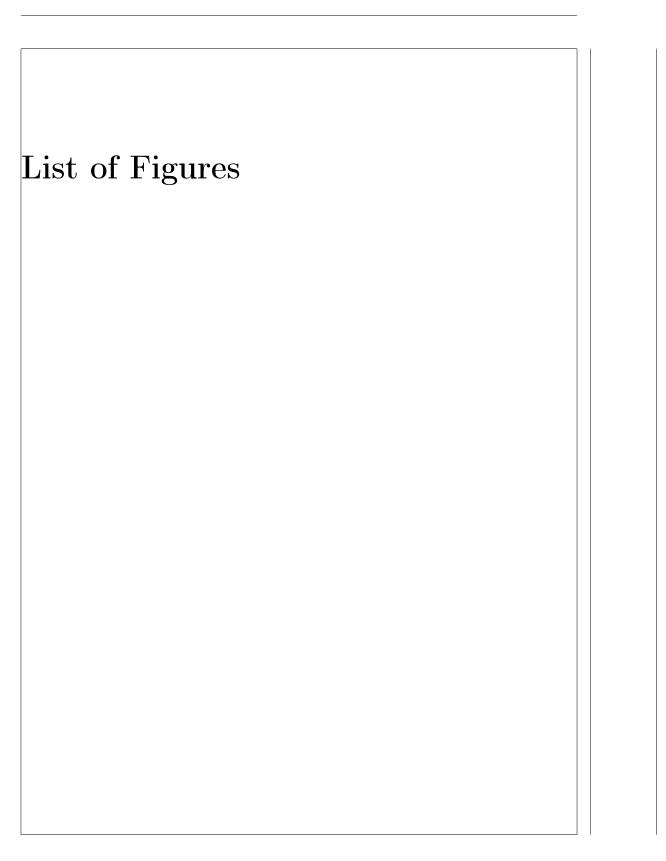
$$E[y_i^*|y_i^* \ge 0] = \tilde{\eta}_i + \frac{\phi(\tilde{\eta}_i)}{\Phi(\tilde{\eta}_i)}$$

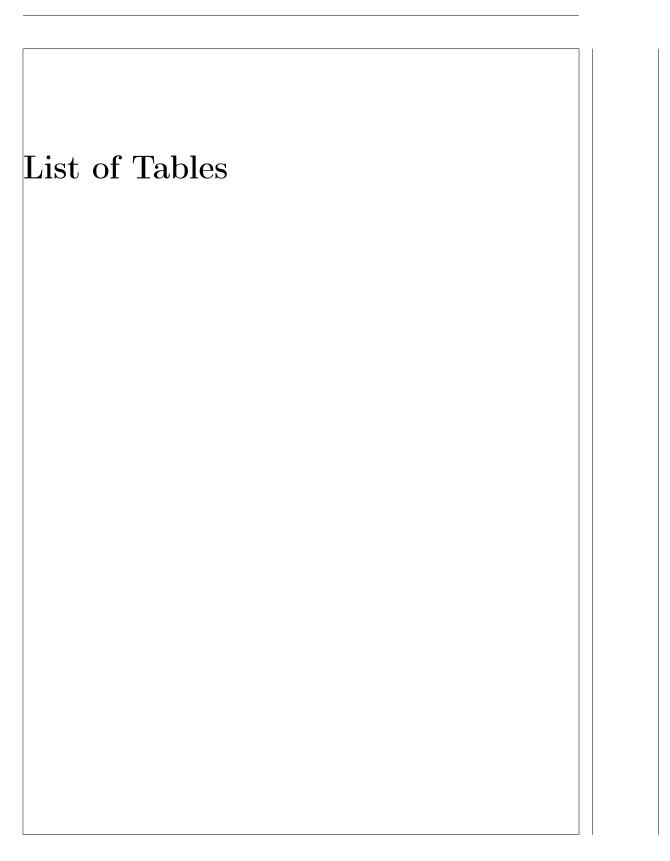
where  $\phi$  and  $\Phi$  are, respectively, the pdf and cdf of a standard normal distribution.

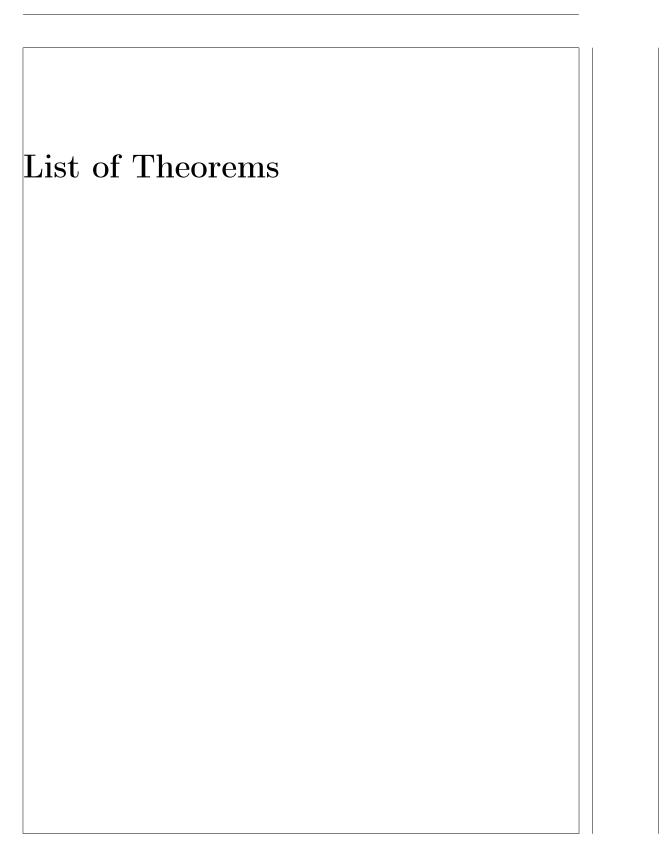
Case:  $y_i = 0$ 

Following the same argument, we can deduce that  $q(y_i^*)$  in this case would be the lower-tail of a one-sided normal distribution truncated at zero. The mean is

$$E[y_i^*|y_i^*<0] = \tilde{\eta}_i + \frac{\phi(\tilde{\eta}_i)}{\Phi(\tilde{\eta}_i) - 1}.$$







# List of Definitions

0.1	Definition (Conically-truncated multivariate normal distribution)	3
6.2	Definition (Differential entropy)	9

# List of Symbols

 $N_p(\mu, \Sigma)$  p-dimensional multivariate normal distribution with mean vector  $\mu$  and covariance  $\Sigma$ .

 $\sim$  Is distributed as.

 $\otimes$  The tensor product.