

# Bayesian Ordinal Regression

April 28, 2022

## 1 Introduction

Ordinal regression is a type of regression analysis used for predicting an ordinal variable, i.e., a variable whose values exists on an arbitrary scale where only the relative ordering between different values is significant. In [1], Gaussian Process model was presented for the ordinal regression models. The following sections describe this model.

## 2 Latent variable model for ordinal regression

Assume we have a training set  $\mathcal{D}$  of  $n$  (IID) observations,  $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, n\}$ , where  $\mathbf{x} \in \mathbb{R}^d$  denotes an input vector (covariates) of dimension  $\mathcal{D}$  and  $y$  denotes an ordinal response variable on a scale  $0, \dots, K$ . Let  $Z$  be a latent variable that underlies the generation of the ordinal responses

$$Z = f(\mathbf{x}) + \epsilon \quad (1)$$

where  $f(\mathbf{x})$  is a zero-mean Gaussian process with covariance function  $\mathcal{K}(\mathbf{x}, \mathbf{x}')$  and  $\epsilon$  is zero mean Gaussian noise with variance  $\sigma_\epsilon^2$ . I.e.,

$$p(Z | f(\mathbf{x}), \sigma_\epsilon^2) = \mathcal{N}(f(\mathbf{x}), \sigma_\epsilon^2). \quad (2)$$

The response variable  $y$  results from an “incomplete measurements” of  $Z$ ,

$$y = \begin{cases} 1, & Z \leq \eta_1, \\ 2, & \eta_1 \leq Z \leq \eta_2 \\ \vdots & \\ K & \eta_{K-1} \leq Z. \end{cases} \quad (3)$$

Defining  $\eta_0 = -\infty$  and  $\eta_K = \infty$ , the above can be summarized as  $y = k$  if (and only if)  $\eta_{k-1} \leq Z \leq \eta_k$ .

Denote with  $\boldsymbol{\theta}$  the parameters vector including the thresholds,  $\{\eta_1, \dots, \eta_{K-1}\}$ , the noise level  $\sigma_\epsilon^2$  and the covariance function parameters. The probability that  $y$  equals  $k$

$$\begin{aligned} P_r(y = k|f(\mathbf{x}), \boldsymbol{\theta}) &= P_r(\eta_{k-1} \leq Z \leq \eta_k | f(\mathbf{x}), \sigma_\epsilon^2) = \\ &= \int_{\eta_{k-1}}^{\eta_k} p(z|f(\mathbf{x}), \sigma_\epsilon^2) dz = \Phi\left(\frac{\eta_k - f(\mathbf{x})}{\sigma_\epsilon}\right) - \Phi\left(\frac{\eta_{k-1} - f(\mathbf{x})}{\sigma_\epsilon}\right) \end{aligned} \quad (4)$$

where  $\Phi$  is the cumulative distribution function of the Gaussian distribution. For  $k = 1$  and  $k = K$  we obtain,

$$\begin{aligned} P_r(y = 0|f(\mathbf{x}), \boldsymbol{\theta}) &= \Phi\left(\frac{\eta_k - f(\mathbf{x})}{\sigma_\epsilon}\right) \\ P_r(y = K|f(\mathbf{x}), \boldsymbol{\theta}) &= 1 - \Phi\left(\frac{\eta_{k-1} - f(\mathbf{x})}{\sigma_\epsilon}\right) \end{aligned}$$

The conditional distribution of  $y$  is now given by

$$\begin{aligned} p(y|f(\mathbf{x}), \boldsymbol{\theta}) &= \prod_{i=1}^K P_r(y = k|f(\mathbf{x}), \boldsymbol{\theta})^{[y=k]} = \\ &= \prod_{k=1}^K \left[ \Phi\left(\frac{\eta_k - f(\mathbf{x})}{\sigma_\epsilon}\right) - \Phi\left(\frac{\eta_{k-1} - f(\mathbf{x})}{\sigma_\epsilon}\right) \right]^{[y=k]} = \\ &= \Phi\left(\frac{\eta_y - f(\mathbf{x})}{\sigma_\epsilon}\right) - \Phi\left(\frac{\eta_{y-1} - f(\mathbf{x})}{\sigma_\epsilon}\right) \end{aligned}$$

using the Iverson bracket  $[y = k]$ . Then, the likelihood of the ordinal model can now be stated as

$$p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta}) = \prod_{i=1}^n p(y_i|f(\mathbf{x}_i), \boldsymbol{\theta}), \quad (5)$$

where  $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]^T$ .

### 3 Full Bayesian treatment

In the full Bayesian treatment, we must assign priors for both to  $\mathbf{f}$  and  $\boldsymbol{\theta}$ . The posterior probability can then be written as

$$\begin{aligned} p(\mathbf{f}, \boldsymbol{\theta}|\mathcal{D}) &= \frac{p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f}, \boldsymbol{\theta})}{p(\mathcal{D})} \stackrel{f, \boldsymbol{\theta} \text{ iid}}{=} \\ &= \frac{p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f})p(\boldsymbol{\theta})}{p(\mathcal{D})} \end{aligned} \quad (6)$$

where  $p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f})p(\boldsymbol{\theta})d\mathbf{f}d\boldsymbol{\theta}$ . The posterior predictive distribution is given by

$$p(y_*|\mathbf{x}_*, \mathcal{D}) = \int p(y_*|f(\mathbf{x}_*), \boldsymbol{\theta})p(\mathbf{f}, \boldsymbol{\theta}|\mathcal{D})d\mathbf{f}d\boldsymbol{\theta} \quad (7)$$

Computing the posterior distribution is analytically intractable and most often Monte Carlo methods are used to obtain approximations.

### 3.1 Prior specification

The prior probability for  $\mathbf{f}$  is a multivariate Gaussian

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, K) \quad (8)$$

where  $K$  is  $n \times n$  covariance matrix whose  $ij$ -element equals  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ .

To set priors to the thresholds,  $\{\eta_1, \dots, \eta_{K-1}\}$  that enforce increasing order and positivity we use the following definition

$$\eta_j = \eta_1 + \sum_{l=2}^j \log \Delta_l, \quad j = 2, \dots, K-1. \quad (9)$$

Now we can assign normal prior over the parameters  $\{\eta_1, \log \Delta_2, \dots, \log \Delta_{K-1}\}$ .

## 4 Partial Bayesian treatment

In a full Bayesian treatment the parameters  $\boldsymbol{\theta}$  must be integrated over the  $\boldsymbol{\theta}$ -space. An alternative solution is to find a point estimate for  $\boldsymbol{\theta}$ . This results in a Bayesian framework conditional on the parameters  $\boldsymbol{\theta}$

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \frac{p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f})}{p(\mathcal{D}|\boldsymbol{\theta})}. \quad (10)$$

A point estimate for  $\boldsymbol{\theta}$  can be either computed by maximizing the evidence  $p(\mathcal{D}|\boldsymbol{\theta})$  (ML estimator) or by maximizing the posterior  $p(\boldsymbol{\theta}|\mathcal{D})$  (MAP estimator), where  $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$ . The evidence (which is required for the estimating  $\boldsymbol{\theta}$ ) is given by a high dimensional integral,

$$p(\mathcal{D}|\boldsymbol{\theta}) = \int p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f})d\mathbf{f}, \quad (11)$$

which is analytically intractable. A popular approach is to approximate the posterior,  $p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta})$  as a Gaussian (e.g. using Laplace approximation), and then the evidence can be calculated using explicit formula. The MLE or MAP estimation of  $\boldsymbol{\theta}$  can then be obtained using gradient-based optimization methods.

#### 4.1 Laplace Approximation

In this section, we develop the Laplace approximation of the posterior  $p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta})$  at the maximum a-posteriori (MAP) estimate. Denotes with  $\Psi(\mathbf{f})$  the (unnormalized) negative log of the posterior

$$\Psi(\mathbf{f}) = -\log p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f}) \quad (12)$$

The Laplace approximation refers to using the Taylor expansion of  $\Psi(\mathbf{f})$  at the MAP point and retaining the terms up to the second order. This is equivalent to approximate the posterior with the following Gaussian distribution

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \mathcal{N}(\hat{\mathbf{f}}, H^{-1}) \quad (13)$$

where  $\hat{\mathbf{f}}$  is the MAP estimate of the posterior and  $H$  is the Hessian matrix of  $\Psi(\mathbf{f})$  evaluated at  $\hat{\mathbf{f}}$

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \Psi(\mathbf{f}) \quad (14)$$

$$H = \frac{\partial^2 \Psi(\mathbf{f})}{\partial \mathbf{f} \partial \mathbf{f}^T} \Big|_{\mathbf{f}=\hat{\mathbf{f}}}. \quad (15)$$

Using Eq. (5),(8) we obtain the following expression for  $\Psi(\mathbf{f})$

$$\Psi(\mathbf{f}) = \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta}) + \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} + \frac{n}{2} \log 2\pi + \frac{1}{2} \log |K| \quad (16)$$

where

$$\begin{aligned} \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta}) &= -\log p(y_i|f(\mathbf{x}_i), \boldsymbol{\theta}) \\ &= -\log \left[ \Phi \left( \frac{\eta_{y_i} - f(\mathbf{x}_i)}{\sigma_\epsilon} \right) - \Phi \left( \frac{\eta_{y_i-1} - f(\mathbf{x}_i)}{\sigma_\epsilon} \right) \right] \end{aligned} \quad (17)$$

In Appendix A we show that the Hessian matrix is given by

$$H = \Lambda + K^{-1}$$

where  $\Lambda$  is diagonal matrix whose  $ii$ -elements is  $\frac{\partial^2 \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i^2}$  given as in (31). We further show (Appendix A) that  $H$  is positive define, hence the optimization problem involve in finding the MAP estimate is convex and have a unique solution. The resulting Gaussian approximation is given by

$$p(\mathbf{f}|\mathcal{D}, \boldsymbol{\theta}) = \mathcal{N}(\hat{\mathbf{f}}, (\Lambda_{\text{MAP}} + K^{-1})^{-1}), \quad (18)$$

where  $\Lambda_{\text{MAP}}$  denotes  $\Lambda$  at the MAP estimate. The Second order Taylor expansion of  $\Psi(\mathbf{f})$  at the MAP estimate is given by

$$\hat{\Psi}(\mathbf{f}) = \Psi(\hat{\mathbf{f}}) + \frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^T (\Lambda_{\text{MAP}} + K^{-1})(\mathbf{f} - \hat{\mathbf{f}}) \quad (19)$$

Table 1: Algorithm for model adaptation using the MAP approach with Laplace approximation

Initialization	choose a favorite gradient-descent optimization package.
	select the starting point $\boldsymbol{\theta}$ for the optimization package.
Looping	While the optimization package requests evidence/gradient evaluation at $\boldsymbol{\theta}$
	(1) Find the MAP estimate by solving the convex optimization problem (14)
	(2) Evaluate the <b>negative</b> log-evidence (20) at the MAP estimate
	(3) Calculate the gradients with respect to $\boldsymbol{\theta}$ (Appendix B).
	(4) feed the evidence and gradients to the optimization package.
Exit	Return the optimal $\boldsymbol{\theta}$ found by optimization package

## 4.2 Model Adaptation

To obtain a point estimate for  $\boldsymbol{\theta}$  we must compute the evidence  $p(\mathcal{D}|\boldsymbol{\theta})$  as defined in (11). Using (12) we obtain

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\theta}) &= \int p(\mathcal{D}|\mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f}) d\mathbf{f} \\ &= \int \exp -\Psi(\mathbf{f}). \end{aligned}$$

Replacing  $\Psi(\mathbf{f})$  with the approximation  $\hat{\Psi}(\mathbf{f})$  (19) we obtain,

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\theta}) &\approx \int \exp -\hat{\Psi}(\mathbf{f}) = \\ &= \exp \left[ -\Psi(\hat{\mathbf{f}}) \right] \int \exp \left[ -\frac{1}{2}(\mathbf{f} - \hat{\mathbf{f}})^T (\Lambda_{\text{MAP}} + K^{-1})(\mathbf{f} - \hat{\mathbf{f}}) \right] d\mathbf{f} \\ &= \exp \left[ -\Psi(\hat{\mathbf{f}}) \right] (2\pi)^{n/2} |\Lambda_{\text{MAP}} + K^{-1}|^{-1/2} = \\ &= \end{aligned}$$

Hence, the log evidence is given by

$$\begin{aligned} \log p(\mathcal{D}|\boldsymbol{\theta}) &\approx - \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) - \frac{1}{2} \hat{\mathbf{f}}^T K^{-1} \hat{\mathbf{f}} - \frac{1}{2} \log |K| |\Lambda_{\text{MAP}} + K^{-1}| + C \\ &= - \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) - \frac{1}{2} \hat{\mathbf{f}}^T K^{-1} \hat{\mathbf{f}} - \frac{1}{2} \log |I + K \Lambda_{\text{MAP}}| + C \end{aligned} \quad (20)$$

The gradient of (20) with respect to the hyper-parameters  $\boldsymbol{\theta}$  can be derived analytically (Appendix B). The outline of algorithm for model adaptation is described in Table 1.

## 4.3 Prediction

At the optimal hyper-parameters we inferred, denoted as  $\hat{\boldsymbol{\theta}}$ , let us take a test case  $\mathbf{x}_*$  for which the ordinal response variable  $y_*$  is unknown. To predict  $y_*$  we first predict

the latent variable  $f_* \triangleq f(\mathbf{x}_*)$ . The posterior predictive distribution  $p(f_*|\mathbf{x}_*\mathcal{D}, \hat{\boldsymbol{\theta}})$  can be written as

$$p(f_*|\mathbf{x}_*, \mathcal{D}, \hat{\boldsymbol{\theta}}) = \int p(f_*|\mathbf{x}_*, \mathbf{f})p(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})d\mathbf{f}$$

where  $p(\mathbf{f}|\mathcal{D}, \hat{\boldsymbol{\theta}})$  is the posterior approximation (18) and  $p(f_*|\mathbf{x}_*\mathbf{f})$  is the prior conditional distribution of  $f_*$  given  $\mathbf{f}$ . Since the posterior and the prior are Gaussian, the posterior predictive distribution is also Gaussian

$$p(f_*|\mathbf{x}_*\mathcal{D}, \hat{\boldsymbol{\theta}}) = \mathcal{N}(\mu_x, \sigma_x^2).$$

Furthermore, the mean  $\mu_x$  and the variance  $\sigma_x^2$  satisfies

$$\mu_x = \mathbf{k}_*^T K^{-1} \boldsymbol{\mu}_p \quad (21)$$

$$\sigma_x^2 = \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K^{-1} - K^{-1} \Sigma_p K^{-1}) \mathbf{k}_* \quad (22)$$

where  $\mathbf{k}_* = [\mathcal{K}(\mathbf{x}_1, \mathbf{x}_*), \mathcal{K}(\mathbf{x}_2, \mathbf{x}_*), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_*)]$ , and  $\boldsymbol{\mu}_p, \Sigma_p$  are the mean vector and covariance matrix of the posterior. Using (18) and the matrix inverse lemma we obtain

$$\begin{aligned} \mu_x &= \mathbf{k}_*^T K^{-1} \hat{\mathbf{f}}, \\ \sigma_x^2 &= \mathcal{K}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (K + \Lambda_{\text{MAP}}^{-1})^{-1} \mathbf{k}_*. \end{aligned}$$

The predictive distribution over ordinal targets  $y_x$  is

$$\begin{aligned} p(y_x|\mathbf{x}_*, \mathcal{D}, \hat{\boldsymbol{\theta}}) &= \int p(y_x|\mathbf{x}_*, f_*, \hat{\boldsymbol{\theta}})p(f_*|\mathbf{x}_*, \mathcal{D}, \hat{\boldsymbol{\theta}}) = \\ &= \Phi\left(\frac{\eta_{y_x} - \mu_x}{\sqrt{\sigma_\epsilon^2 + \sigma_x^2}}\right) - \Phi\left(\frac{\eta_{y_x-1} - \mu_x}{\sqrt{\sigma_\epsilon^2 + \sigma_x^2}}\right). \end{aligned} \quad (23)$$

## Appendix A - Hessian Matrix

The Hessian matrix of  $\Psi(\mathbf{f})$  (16) is given by

$$\begin{aligned} H &= \frac{\partial^2 \Psi}{\partial \mathbf{f} \partial \mathbf{f}^T} = \frac{\partial}{\partial \mathbf{f} \partial \mathbf{f}^T} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta}) + \frac{\partial^2}{\partial \mathbf{f} \partial \mathbf{f}^T} \frac{1}{2} \mathbf{f}^T K^{-1} \mathbf{f} \\ &= \frac{\partial^2}{\partial \mathbf{f} \partial \mathbf{f}^T} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta}) + K^{-1} = \Lambda + K^{-1} \end{aligned} \quad (24)$$

where  $\Lambda$  is the Hessian matrix of  $\sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta})$ . The  $ij$ -element of  $\Lambda$  is given by

$$\Lambda_{ij} = \frac{\partial^2}{\partial f(\mathbf{x}_i) \partial f(\mathbf{x}_j)} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta}) = \begin{cases} \frac{\partial^2 \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta})}{\partial f(\mathbf{x}_i)^2} & i = j \\ 0 & i \neq j \end{cases}. \quad (25)$$

That is,  $\Lambda$  is a diagonal matrix whose  $ii$  element equals  $\frac{\partial^2 \ell(y_i, f(\mathbf{x}_i), \boldsymbol{\theta})}{\partial f(\mathbf{x}_i)^2}$ . For better clarity we use  $f_i = f(\mathbf{x}_i)$  hereinafter. Using Eq. (17) we obtain,

$$\begin{aligned}\Lambda_{ii} &= \frac{\partial^2 \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i^2} = \\ &= -\frac{\partial^2}{\partial f_i^2} \log \left[ \Phi \left( \frac{\eta_{y_i} - f_i}{\sigma_\epsilon} \right) - \Phi \left( \frac{\eta_{y_i-1} - f_i}{\sigma_\epsilon} \right) \right], \\ &= -\frac{\partial^2}{\partial f_i^2} \log [\Phi(z_1^i) - \Phi(z_2^i)]\end{aligned}\quad (26)$$

where  $z_1^i = \frac{\eta_{y_i} - f_i}{\sigma_\epsilon}$  and  $z_2^i = \frac{\eta_{y_i-1} - f_i}{\sigma_\epsilon}$ . The first derivative of  $\log [\Phi(z_1^i) - \Phi(z_2^i)]$  with respect to  $f_i$  is given by

$$\begin{aligned}\frac{\partial}{\partial f_i} \log [\Phi(z_1^i) - \Phi(z_2^i)] &= \frac{\Phi'(z_1^i) - \Phi'(z_2^i)}{[\Phi(z_1^i) - \Phi(z_2^i)]} \\ &= -\frac{1}{\sigma_\epsilon} \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{[\Phi(z_1^i) - \Phi(z_2^i)]}\end{aligned}\quad (27)$$

The second derivative with respect to  $f_i$  is given by

$$\begin{aligned}\frac{\partial^2}{\partial f_i^2} \log [\Phi(z_1^i) - \Phi(z_2^i)] &= -\frac{1}{\sigma_\epsilon} \frac{\partial}{\partial f_i} \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{[\Phi(z_1^i) - \Phi(z_2^i)]} \\ &= -\frac{1}{\sigma_\epsilon} \frac{\partial}{\partial f_i} \frac{h(f_i)}{k(f_i)} = -\frac{1}{\sigma_\epsilon} \frac{h'(f_i)k(f_i) - h(f_i)k'(f_i)}{k(f_i)^2}\end{aligned}\quad (28)$$

where

$$\begin{aligned}h(f_i) &= \mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1) \\ h'(f_i) &= \frac{\partial}{\partial f_i} [\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)] \\ &= -\mathcal{N}(z_1^i, 0, 1) z_1^i \frac{\partial z_1^i}{\partial f_i} + \mathcal{N}(z_2^i, 0, 1) z_2^i \frac{\partial z_2^i}{\partial f_i} \\ &= \frac{1}{\sigma_\epsilon} [z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)] \\ k(f_i) &= \Phi(z_1^i) - \Phi(z_2^i) \\ k'(f_i) &= \mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial f_i} - \mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial f_i}\end{aligned}$$

Thus we obtain,

$$\begin{aligned}h'(f_i)k(f_i) - h(f_i)k'(f_i) &= \frac{1}{\sigma_\epsilon} [z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)] [\Phi(z_1^i) - \Phi(z_2^i)] \\ &\quad + \frac{1}{\sigma_\epsilon} [\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)]^2\end{aligned}\quad (29)$$

and hence,

$$\begin{aligned} \frac{\partial^2}{\partial f_i^2} \log [\Phi(z_1^i) - \Phi(z_2^i)] &= -\frac{1}{\sigma_\varepsilon^2} \left[ \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \right] \\ &\quad - \frac{1}{\sigma_\varepsilon^2} \left[ \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \right]^2 \end{aligned} \quad (30)$$

Substituting Eq. (30) in Eq. (26) we obtain,

$$\begin{aligned} \Lambda_{ii} &= \frac{1}{\sigma_\varepsilon^2} \left[ \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \right] + \\ &\quad + \frac{1}{\sigma_\varepsilon^2} \left[ \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \right]^2 \end{aligned} \quad (31)$$

We further note that the first term in Eq. (31) is positive. The constraint  $\eta_{y_i} > \eta_{y_i-1}$  impose  $z_1^i > z_2^i$  and therefore  $\Phi(z_1^i) - \Phi(z_2^i) > 0$ . Additionally we can show that  $z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1) > 0$  (TBD). Hence,  $\Lambda$  is a positive-definite matrix and hence  $H$  (24) is also positive definite (sum of positive definite matrices).

## Appendix B

Evidence maximization is equivalent to finding the minimizer for the negative log evidence which according to (20) can be written as

$$g(\boldsymbol{\theta}) \triangleq -\log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) + \frac{1}{2} \hat{\mathbf{f}}^T K^{-1} \hat{\mathbf{f}} + \frac{1}{2} \log |I + K \Lambda_{\text{MAP}}|. \quad (32)$$

The hyper-parameter vector  $\boldsymbol{\theta}$  include the threshold parameters  $\{\eta_1, \log \Delta_2, \dots, \log \Delta_{K-1}\}$ , the noise variance  $\log \sigma_\varepsilon$  and the kernel parameters. The loss function  $\ell$  and the covariance matrix  $K$  are function of  $\boldsymbol{\theta}$ , but  $\hat{\mathbf{f}}$  and therefore  $\Lambda_{\text{MAP}}$  are also implicitly functions of  $\boldsymbol{\theta}$ , since when  $\boldsymbol{\theta}$  changes, the optimum of the posterior  $\hat{\mathbf{f}}$  also changes. Thus

$$\frac{\partial g(\boldsymbol{\theta})}{\partial \theta_i} = \frac{\partial g(\boldsymbol{\theta})}{\partial \theta_i} \Big|_{\text{explicit}} + \sum_{i=1}^n \frac{\partial g(\boldsymbol{\theta})}{\partial \hat{f}_i} \frac{\partial \hat{f}_i}{\partial \theta_j} \quad (33)$$



by the chain rule. The explicit term is given by

$$\begin{aligned}
\frac{\partial g(\boldsymbol{\theta})}{\partial \theta_j} \Big|_{\text{explicit}} &= \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) - \frac{1}{2} \hat{\mathbf{f}}^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} \hat{\mathbf{f}} + \\
&\quad + \frac{1}{2} \text{trace} \left[ (I + K \Lambda_{\text{MAP}})^{-1} \frac{\partial}{\partial \theta_j} K \Lambda_{\text{MAP}} \right] \\
&= \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) - \frac{1}{2} \hat{\mathbf{f}}^T K^{-1} \frac{\partial K}{\partial \theta_j} K^{-1} \hat{\mathbf{f}} \\
&\quad + \frac{1}{2} \text{trace} \left[ (\Lambda_{\text{MAP}}^{-1} + K)^{-1} \frac{\partial K}{\partial \theta_j} \right] \\
&\quad + \frac{1}{2} \text{trace} \left[ \Lambda_{\text{MAP}}^{-1} (\Lambda_{\text{MAP}}^{-1} + K)^{-1} K \frac{\partial}{\partial \theta_j} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} \right] =
\end{aligned} \tag{34}$$

where we used Eq. (19) for the expression of  $\Lambda_{\text{MAP}}$ . When evaluating the remaining term from Eq. (33), we utilize the fact that  $\hat{\mathbf{f}}$  is the maximum of the posterior so that  $\partial \Psi(\mathbf{f}) / \partial \mathbf{f} = \mathbf{0}$  for  $\mathbf{f} = \hat{\mathbf{f}}$  where  $\Psi(\mathbf{f})$  is defined in Eq. (16). Thus the implicit derivatives of the first two term in Eq. (32) vanish, leaving only

$$\begin{aligned}
\frac{\partial g(\boldsymbol{\theta})}{\partial \hat{f}_i} &= \frac{1}{2} \frac{\partial \log |I + K \Lambda_{\text{MAP}}|}{\partial \hat{f}_i} = \frac{1}{2} \text{trace} \left[ (I + K \Lambda_{\text{MAP}})^{-1} K \frac{\partial \Lambda_{\text{MAP}}}{\partial \hat{f}_i} \right] \\
&= \frac{1}{2} \text{trace} \left[ (K^{-1} + \Lambda_{\text{MAP}})^{-1} \frac{\partial \Lambda_{\text{MAP}}}{\partial \hat{f}_i} \right] = \frac{1}{2} [(K^{-1} + \Lambda_{\text{MAP}})^{-1}]_{ii} \frac{\partial^3 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^3}
\end{aligned} \tag{35}$$

In order to evaluate the derivative  $\partial \hat{\mathbf{f}} / \partial \theta_j$ , we differentiate the fixed-point equation  $\hat{\mathbf{f}} = -K \sum_{i=1}^n \partial \ell(y_i, f_i, \boldsymbol{\theta}) / \partial \mathbf{f} \Big|_{\mathbf{f}=\hat{\mathbf{f}}}$  (Rasmussen et al.) to obtain

$$\begin{aligned}
\frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} &= -\frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial \mathbf{f}} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} - K \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \\
&= -\frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial \mathbf{f}} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} - K \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \left[ \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \right] \\
&\quad - K \underbrace{\frac{\partial}{\partial \hat{\mathbf{f}}} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}}}_{\Lambda_{\text{MAP}}} \frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} = -\frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \\
&\quad - K \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \left[ \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \right] - K \Lambda_{\text{MAP}} \frac{\partial \hat{\mathbf{f}}}{\partial \theta_j}
\end{aligned}$$

Thus we obtain

$$\begin{aligned}
\frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} &= -(I + K\Lambda_{\text{MAP}})^{-1} \left\{ \frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} + K \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \left[ \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \right] \right\} \\
\frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} &= -\frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} - K \frac{\partial}{\partial \hat{\mathbf{f}}} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} \frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} \\
&= -\frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}} - K\Lambda_{\text{MAP}} \frac{\partial \hat{\mathbf{f}}}{\partial \theta_j} \\
&= -(I + K\Lambda_{\text{MAP}})^{-1} \frac{\partial K}{\partial \theta_j} \sum_{i=1}^n \frac{\partial \ell(y_i, f_i, \boldsymbol{\theta})}{\partial f_i} \Big|_{\mathbf{f}=\hat{\mathbf{f}}}
\end{aligned} \tag{36}$$

The desired derivatives are obtained by plugging Eq. (34-36) to Eq. (33). We Note that the following gradients are required for solving the optimization problem

$$\frac{\partial K}{\partial \theta_j}, \frac{\partial}{\partial \theta_j} \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}), \frac{\partial}{\partial \theta_j} \frac{\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i}, \frac{\partial}{\partial \theta_j} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2}, \frac{\partial^3 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^3}$$

Using Eq. (26) The derivative  $\partial^3 \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) / \partial f_i^3$  in (35) becomes

$$\frac{\partial^3 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^3} = -\frac{\partial^3}{\partial f_i^3} \log [\Phi(z_1^i) - \Phi(z_2^i)] \tag{37}$$

where  $z_1^i = \frac{\eta_{y_i} - f_i}{\sigma_\epsilon}$  and  $z_2^i = \frac{\eta_{y_i} - 1 - f_i}{\sigma_\epsilon}$ . Using the result in (30) we obtain

$$\frac{\partial^3 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^3} = \frac{\partial}{\partial f_i} \frac{1}{\sigma_\epsilon^2} [v_0^2 + v_1] = \frac{1}{\sigma_\epsilon^2} \left( 2v_0 \frac{\partial v_0}{\partial f_i} + \frac{\partial v_1}{\partial f_i} \right) \tag{38}$$

where

$$v_p^i = \frac{(z_1^i)^p \mathcal{N}(z_1^i, 0, 1) - (z_2^i)^p \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \tag{39}$$

Using Eq. (30) The term  $\partial v_0 / \partial f_i$  can be expressed as

$$\frac{\partial v_0^i}{\partial f_i} = \frac{1}{\sigma_\epsilon} [(v_0^i)^2 + v_1^i] \tag{40}$$

The term  $\partial v_1 / \partial f_i$  can be expressed as

$$\frac{\partial v_1^i}{\partial f_i} = \frac{h'(f_i)k(f_i) - h(f_i)k'(f_i)}{k(f_i)^2} \tag{41}$$

where

$$\begin{aligned}
h(f_i) &= z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1) \\
h'(f_i) &= \mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial f_i} - (z_1^i)^2 \mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial f_i} - \mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial f_i} + (z_2^i)^2 \mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial f_i} \\
&= -\frac{1}{\sigma_\epsilon} \mathcal{N}(z_1^i, 0, 1) + \frac{1}{\sigma_\epsilon} (z_1^i)^2 \mathcal{N}(z_1^i, 0, 1) + \frac{1}{\sigma_\epsilon} \mathcal{N}(z_2^i, 0, 1) - \frac{1}{\sigma_\epsilon} (z_2^i)^2 \mathcal{N}(z_2^i, 0, 1) \\
&= -\frac{1}{\sigma_\epsilon} [\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)] + \frac{1}{\sigma_\epsilon} [(z_1^i)^2 \mathcal{N}(z_1^i, 0, 1) - (z_2^i)^2 \mathcal{N}(z_2^i, 0, 1)] \\
k(f_i) &= \Phi(z_1^i) - \Phi(z_2^i) \\
k'(f_i) &= \mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial f_i} - \mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial f_i} \\
&= -\frac{1}{\sigma_\epsilon} [\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)]
\end{aligned}$$

Thus,

$$\begin{aligned}
\frac{\partial v_1}{\partial f_i} &= -\frac{1}{\sigma_\epsilon} \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} + \frac{1}{\sigma_\epsilon} \frac{(z_1^i)^2 \mathcal{N}(z_1^i, 0, 1) - (z_2^i)^2 \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \\
&+ \frac{1}{\sigma_\epsilon} \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \\
&= \frac{1}{\sigma_\epsilon} [-v_0^i + v_2^i + v_1^i v_0^i]
\end{aligned} \tag{42}$$

Substituting Eq. (40,42) to Eq. (38) we obtain

$$\frac{\partial^3 \ell(y_i \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^3} = \frac{1}{\sigma_\epsilon^2} \left( 2v_0^i \frac{\partial v_0^i}{\partial f_i} + \frac{\partial v_1^i}{\partial f_i} \right) = \frac{1}{\sigma_\epsilon^3} [2(v_0^i)^3 + 3v_0^i v_1^i + v_2^i - v_0^i] \tag{43}$$

The derivative of  $\sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta})$  with respect to  $\theta_j$  is given by

$$\frac{\partial}{\partial \theta_j} \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) = -\frac{\partial}{\partial \theta_j} \sum_{i=1}^n \log [\Phi(z_1^i) - \Phi(z_2^i)] = -\sum_{i=1}^n \frac{\mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial \theta_j} - \mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial \theta_j}}{[\Phi(z_1^i) - \Phi(z_2^i)]} \tag{44}$$

The derivative of  $\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) / \partial f_i$  and  $\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) / \partial f_i^2$  with respect to  $\theta_j$  is given by (using Eq. 39)

$$\frac{\partial}{\partial \theta_j} \frac{\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i} = \frac{\partial}{\partial \theta_j} \frac{1}{\sigma_\epsilon} v_0^i \tag{45}$$

$$\frac{\partial}{\partial \theta_j} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} = \frac{\partial}{\partial \theta_j} \frac{1}{\sigma_\epsilon^2} [(v_0^i)^2 + v_1^i] = \frac{1}{\sigma_\epsilon^2} \left( 2v_0^i \frac{\partial v_0^i}{\partial \theta_j} + \frac{\partial v_1^i}{\partial \theta_j} \right) \tag{46}$$

### The derivatives with respect to $\sigma_\epsilon$

For  $\theta_i = \sigma_\epsilon$  we obtain

$$\frac{\partial z_1^i}{\partial \sigma_\epsilon} = -\frac{\eta_{y_i} - f_i}{\sigma_\epsilon^2} = -\frac{z_1^i}{\sigma_\epsilon}, \quad \frac{\partial z_2^i}{\partial \sigma_\epsilon} = -\frac{z_2^i}{\sigma_\epsilon}$$

The derivative of  $\sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta})$  with respect to  $\sigma_\epsilon$  is given by

$$\frac{\partial}{\partial \sigma_\epsilon} \sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) = \frac{1}{\sigma_\epsilon} \sum_{i=1}^n v_1^i \quad (47)$$

To compute the derivative in Eq. (45), (46) with respect to  $\theta_j = \sigma_\epsilon$  we first compute the derivatives  $\partial v_0 / \partial \sigma_\epsilon$  and  $\partial v_1 / \partial \sigma_\epsilon$

$$\begin{aligned} \frac{\partial v_0}{\partial \sigma_\epsilon} &= \frac{1}{\sigma_\epsilon} \frac{[(z_1^i)^2 \mathcal{N}(z_1^i, 0, 1) - (z_2^i)^2 \mathcal{N}(z_2^i, 0, 1)]}{\Phi(z_1^i) - \Phi(z_2^i)} \\ &= + \frac{1}{\sigma_\epsilon} \frac{\mathcal{N}(z_1^i, 0, 1) - \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \frac{\partial}{\partial \theta_j} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} \\ &= \frac{1}{\sigma_\epsilon} (v_2^i + v_0^i v_1^i) \end{aligned} \quad (48)$$

$$\begin{aligned} \frac{\partial v_1}{\partial \sigma_\epsilon} &= -\frac{1}{\sigma_\epsilon} \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \\ &\quad + \frac{1}{\sigma_\epsilon} \frac{(z_1^i)^3 \mathcal{N}(z_1^i, 0, 1) - (z_2^i)^3 \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \\ &\quad + \frac{1}{\sigma_\epsilon} \left( \frac{z_1^i \mathcal{N}(z_1^i, 0, 1) - z_2^i \mathcal{N}(z_2^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \right)^2 \\ &= \frac{1}{\sigma_\epsilon} [-v_1^i + v_2^i + (v_1^i)^2] \end{aligned}$$

Thus we obtain,

$$\frac{\partial}{\partial \sigma_\epsilon} \frac{\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i} = \frac{1}{\sigma_\epsilon^2} (v_2^i + v_0^i v_1^i) \quad (49)$$

$$\frac{\partial}{\partial \sigma_\epsilon} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} = \frac{1}{\sigma_\epsilon^3} [2v_0^i v_2^i + 2(v_0^i)^2 v_1^i - v_1^i + v_3^i + (v_1^i)^2] \quad (50)$$

### The derivatives with respect to $\eta_1$

We note that

$$\eta_i = \begin{cases} -\infty & i = 0 \\ \eta_1 & i = 1 \\ \eta_1 + \sum_{l=1}^i \Delta_l & i \in [2, K-1] \\ \infty & i = K \end{cases} \quad (51)$$

For  $\theta_j = \eta_1$  we obtain,

$$\frac{\partial z_1^i}{\partial \eta_1} = \frac{\partial}{\partial \eta_1} \frac{\eta_{y_i} - f_i}{\sigma_\epsilon} = \begin{cases} \frac{1}{\sigma_\epsilon} & y_i = [1, K-1] \\ 0 & y_i = K \end{cases}$$

$$\frac{\partial z_2^i}{\partial \eta_1} = \frac{\partial}{\partial \eta_1} \frac{\eta_{y_i-1} - f_i}{\sigma_\epsilon} = \begin{cases} 0 & y_i = 1 \\ \frac{1}{\sigma_\epsilon} & y_i = [2, K] \end{cases}$$

We note however that

$$\mathcal{N}(z_1^i, 0, 1) \frac{\partial z_1^i}{\partial f_i} = \frac{1}{\sigma_\epsilon} \mathcal{N}(z_1^i, 0, 1)$$

$$\mathcal{N}(z_2^i, 0, 1) \frac{\partial z_2^i}{\partial f_i} = \frac{1}{\sigma_\epsilon} \mathcal{N}(z_2^i, 0, 1)$$

The derivative of  $\sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta})$  with respect to  $\eta_1$  is given by (using Eq. 44)

$$\frac{\partial}{\partial \eta_1} \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) = -\frac{1}{\sigma_\epsilon} v_0^i \quad (52)$$

To compute the derivative in Eq. (45),(46) with respect to  $\theta_j = \eta_1$  we first compute the derivatives  $\partial v_0 / \partial \eta_1$  and  $\partial v_1 / \partial \eta_1$

$$\frac{\partial v_0^i}{\partial \eta_1} = -\frac{1}{\sigma_\epsilon} [(v_0^i)^2 + v_1^i]$$

$$\frac{\partial v_1^i}{\partial \eta_1} = -\frac{1}{\sigma_\epsilon} (-v_0^i + v_2^i + v_1^i v_0^i)$$

Thus we obtain,

$$\frac{\partial}{\partial \eta_1} \frac{\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i} = -\frac{1}{\sigma_\epsilon^2} [(v_0^i)^2 + v_1^i] \quad (53)$$

$$\frac{\partial}{\partial \eta_1} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} = -\frac{1}{\sigma_\epsilon^3} [2(v_0^i)^3 + 3v_0^i v_1^i - v_0^i + v_2^i] \quad (54)$$

## The derivatives with respect to $\Delta_l$

Using Eq. (51) we obtain

$$\frac{\partial z_1^i}{\partial \Delta_l} = \begin{cases} 0 & y_i \leq l \\ \frac{1}{\sigma_\epsilon} & \text{else} \end{cases}, \quad l = 2, \dots, K-1$$

$$\frac{\partial z_2^i}{\partial \Delta_l} = \begin{cases} 0 & y_i - 1 \leq l \\ \frac{1}{\sigma_\epsilon} & \text{else} \end{cases}, \quad l = 2, \dots, K-1$$

The derivative of  $\sum_{i=1}^n \ell(y_i, \hat{f}_i, \boldsymbol{\theta})$  with respect to  $\Delta_l$  is given by (using Eq. 44)

$$\frac{\partial}{\partial \Delta_l} \ell(y_i, \hat{f}_i, \boldsymbol{\theta}) = \begin{cases} 0 & y_i < l \\ -\frac{1}{\sigma_\epsilon} s_0^i & y_i = l, \quad l = 2, \dots, K-1 \\ -\frac{1}{\sigma_\epsilon} v_0^i & y_i > l \end{cases}$$

where

$$s_p = \frac{(z_1^i)^p \mathcal{N}(z_1^i, 0, 1)}{\Phi(z_1^i) - \Phi(z_2^i)} \quad (55)$$

To compute the derivative in Eq. (45), (46) with respect to  $\theta_j = \eta_1$  we first compute the derivatives  $\partial v_0^i / \partial \eta_1$  and  $\partial v_1^i / \partial \eta_1$

$$\frac{\partial v_0^i}{\partial \Delta_l} = \begin{cases} 0 & y_i < l \\ -\frac{1}{\sigma_\epsilon} [s_1^i + v_0^i s_0^i] & y_i = l, \quad l = 2, \dots, K-1 \\ -\frac{1}{\sigma_\epsilon} [(v_0^i)^2 + v_1^i] & y_i > l \end{cases}$$

$$\frac{\partial v_1^i}{\partial \Delta_l} = \begin{cases} 0 & y_i < l \\ -\frac{1}{\sigma_\epsilon} (-s_0^i + s_2^i + v_1^i s_0^i) & y_i = l, \quad l = 2, \dots, K-1 \\ -\frac{1}{\sigma_\epsilon} (-v_0^i + v_2^i + v_1^i v_0^i) & y_i > l \end{cases}$$

Thus we obtain,

$$\frac{\partial}{\partial \Delta_l} \frac{\partial \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i} = \begin{cases} 0 & y_i < l \\ -\frac{1}{\sigma_\epsilon^2} [v_0^i s_0^i + s_1^i] & y_i = l, \quad l = 2, \dots, K-1 \\ -\frac{1}{\sigma_\epsilon^2} [(v_0^i)^2 + v_1^i] & y_i > l \end{cases}$$

$$\frac{\partial}{\partial \Delta_l} \frac{\partial^2 \ell(y_i, \hat{f}_i, \boldsymbol{\theta})}{\partial f_i^2} = \begin{cases} 0 & y_i < l \\ -\frac{1}{\sigma_\epsilon^3} [2(v_0^i)^2 s_0^i + 2v_0^i s_1^i - s_0^i + s_2^i + v_1^i s_0^i] & y_i = l \\ -\frac{1}{\sigma_\epsilon^3} [2(v_0^i)^3 + 3v_0^i v_1^i - v_0^i + v_2^i] & y_i > l \end{cases}$$

## References

- [1] Chu, Wei, Zoubin Ghahramani, and Christopher KI Williams. "Gaussian processes for ordinal regression." *Journal of machine learning research* 6.7 (2005).