

CS3390 Assignment 1

Problem 2

Gautam Singh (CS21BTECH11018)
Jaswanth Beere (BM21BTECH11007)

CONTENTS

1	Summary	1
2	Differences From Other Methods	1
3	Maximum Likelihood Parameter Estimation	2
	References	4

This document summarises and demonstrates the methods described in [1] using the dataset in [2].

1 SUMMARY

Data is said to be *ordinal* if it can be grouped into ordered categories. These categories can be thought of as intervals of a function of some underlying, unknown continuous random variable, where these intervals can be of any distance and non-continuous. Two models for performing regression on ordinal data, namely the *proportional odds* and *proportional hazards* model, are described in [1] along with their practical use cases. Suppose that the response variable belongs to k ordered categories with probabilities $\pi_j(\mathbf{x})$. In both cases, the models are of the form

$$\text{link}\{\gamma_j(\mathbf{x})\} = \theta_j - \boldsymbol{\beta}^\top \mathbf{x} \quad (1)$$

which is a general linear model where

- 1) $\gamma_j(\mathbf{x}) = \sum_{i=1}^j \pi_i(\mathbf{x})$ denotes the probability that the response variable falls in category at most j for inputs \mathbf{x} .
- 2) $\boldsymbol{\beta}$ and θ_j , $1 \leq j < k$ are the parameters of the model to be estimated. Usually, θ_j are referred to as *cut points*.
- 3) link is a monotone function mapping $(0, 1)$ to $(-\infty, \infty)$. The link function should be selected based on ease of interpretation for the application.

Also studied in this paper are nonlinear models of the form

$$\text{link}\{\gamma_j(\mathbf{x}_i)\} = \frac{\theta_j - \boldsymbol{\beta}^\top \mathbf{x}_i}{\tau_i} \quad (2)$$

which are useful for applications where the covariates \mathbf{x} may not have the same “variance”. Here, τ_i is called the *scale* of the i th row of a dataset, and is given by

$$\log \tau_i = \boldsymbol{\tau}^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \quad (3)$$

where $\boldsymbol{\tau}$ is a vector of parameters of the model.

In both cases, it turns out that numerical methods of solving these models, in particular reweighted least squares, converges to the maximum likelihood estimate. The paper also contrasts the performance of ordinal regression with application-specific qualitative tests and alternative models for various healthcare-related applications, illustrating that models such as (1) are easy to compute and interpret.

2 DIFFERENCES FROM OTHER METHODS

It is tempting to consider the ordinal categories as classes and use a classification model. However, in cases where the categories are models of underlying continuous random variables and processes, it is better to use ordinal regression. Using classification methods is amenable only when the proportion of classes in the sample dataset and population dataset are equal, which is not always possible. Classification in such cases can lead to indicators with poor generalization performance.

On the other hand, regression is suitable for applications where the response variable is real-valued. Obviously, applying regression directly for ordinal data will not work, since the data is discrete and grouped into categories. Further, it is difficult to obtain information about the underlying model or

process which puts the response variable into these intervals. Therefore, ordinal regression makes use of a link function which makes the model amenable to the same tools as for linear regression, such as least squares and gradient descent.

3 MAXIMUM LIKELIHOOD PARAMETER ESTIMATION

Suppose that the categorical responses in a dataset of n inputs are distributed as $\{n_1, n_2, \dots, n_k\}$. Define for $1 \leq j \leq k$,

$$R_j \triangleq \sum_{i=1}^j n_i, \quad Z_j \triangleq \frac{R_j}{n}. \quad (4)$$

Using the definitions of γ_j described earlier, define

$$\phi_j \triangleq \log\left(\frac{\gamma_j}{\gamma_{j+1} - \gamma_j}\right) = \text{logit}\left(\frac{\gamma_j}{\gamma_{j+1}}\right) \quad (5)$$

$$g(\phi_j) \triangleq \log(1 + \exp(\phi_j)) = \log\left(\frac{\gamma_{j+1}}{\gamma_{j+1} - \gamma_j}\right). \quad (6)$$

From (5), we have

$$\frac{\partial \phi_j}{\partial \gamma_j} = \left(\frac{\gamma_{j+1} - \gamma_j}{\gamma_j}\right) \frac{(\gamma_{j+1} - \gamma_j) - (-\gamma_j)}{(\gamma_{j+1} - \gamma_j)^2} \quad (7)$$

$$= \frac{\gamma_{j+1}}{\gamma_j(\gamma_{j+1} - \gamma_j)} \quad (8)$$

$$\frac{\partial \phi_j}{\partial \gamma_{j+1}} = \left(\frac{\gamma_{j+1} - \gamma_j}{\gamma_j}\right) \frac{-\gamma_j}{(\gamma_{j+1} - \gamma_j)^2} \quad (9)$$

$$= \frac{-1}{\gamma_{j+1} - \gamma_j} \quad (10)$$

$$\Rightarrow \frac{\partial \phi_j}{\partial \gamma_j} = -\frac{\gamma_{j+1}}{\gamma_j} \frac{\partial \phi_j}{\partial \gamma_{j+1}} \quad (11)$$

Noting that $\gamma_k = 1$, the likelihood can be written as

$$L = \prod_{j=1}^k \pi_j^{n_j} \quad (12)$$

$$= \gamma_1^{R_1} \prod_{j=1}^{k-1} (\gamma_{j+1} - \gamma_j)^{R_{j+1} - R_j} \quad (13)$$

$$= \gamma_1^{R_1} \prod_{j=1}^{k-1} \gamma_{j+1}^{R_{j+1}} \left(\frac{1}{\gamma_{j+1}^{R_j}}\right) \left(\frac{\gamma_{j+1} - \gamma_j}{\gamma_{j+1}}\right)^{R_{j+1} - R_j} \quad (14)$$

$$= \prod_{j=1}^{k-1} \left(\frac{\gamma_j}{\gamma_{j+1}}\right)^{R_j} \left(\frac{\gamma_{j+1} - \gamma_j}{\gamma_{j+1}}\right)^{R_{j+1} - R_j} \quad (15)$$

$$= \prod_{j=1}^{k-1} \left(\frac{\gamma_j}{\gamma_{j+1} - \gamma_j}\right)^{R_j} \left(\frac{\gamma_{j+1} - \gamma_j}{\gamma_{j+1}}\right)^{R_{j+1}}. \quad (16)$$

Therefore, using (4), the log-likelihood is

$$l = n \sum_{i=1}^{k-1} (Z_i \phi_i - Z_{i+1} g(\phi_i)). \quad (17)$$

To find the maximum likelihood estimate of θ_j , β and τ , [1] recommends the use of the Newton-Raphson Method with Fisher scoring. This method converges rapidly even for poor initial estimates, provided that the θ_j are monotone increasing.

Define

$$\beta^* \triangleq (\theta_1 \quad \theta_2 \quad \dots \quad \theta_{k-1} \quad \beta_1 \quad \dots \quad \beta_p)^\top \quad (18)$$

and denote the entire parameter vector as

$$\psi \triangleq \begin{pmatrix} \beta^* \\ \tau \end{pmatrix}. \quad (19)$$

Define

$$\mathbf{X}_j^* \triangleq \begin{pmatrix} \mathbf{e}_j \\ \mathbf{x} \end{pmatrix} \quad (20)$$

$$\mathbf{U} \triangleq \bar{\mathbf{x}} - \mathbf{x} \quad (21)$$

$$w \triangleq \exp(\tau^\top \mathbf{U}). \quad (22)$$

Then, (2) can be rewritten as

$$Y_j = \text{link}(\gamma_j) = \beta^{*\top} \mathbf{X}_j^* \exp(\tau^\top \mathbf{U}). \quad (23)$$

Here, \mathbf{e}_j^\top denotes the j th standard basic vector in \mathbb{R}^{k-1} .

From (23), we have

$$\frac{\partial Y_j}{\partial \beta_r^*} = X_{jr}^* \exp(\tau^\top \mathbf{U}) = w X_{jr}^* \quad (24)$$

$$\frac{\partial Y_j}{\partial \tau_r} = U_r \left(\beta^{*\top} \mathbf{X}_j^* \exp(\tau^\top \mathbf{U}) \right) = Y_j U_r \quad (25)$$

Using the chain rule,

$$\frac{\partial l}{\partial \beta_r^*} = \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} \frac{\partial \phi_j}{\partial \beta_r^*} \quad (26)$$

$$= \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} \left(\frac{\partial \phi_j}{\partial \gamma_j} \frac{\partial \gamma_j}{\partial \beta_r^*} + \frac{\partial \phi_j}{\partial \gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial \beta_r^*} \right) \quad (27)$$

$$= \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} \left(\frac{\partial \phi_j}{\partial \gamma_j} \frac{\partial \gamma_j}{\partial Y_j} \frac{\partial Y_j}{\partial \beta_r^*} + \frac{\partial \phi_j}{\partial \gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} \frac{\partial Y_{j+1}}{\partial \beta_r^*} \right) \quad (28)$$

$$= w \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} \frac{\partial \phi_j}{\partial \gamma_j} \left(\frac{\partial \gamma_j}{\partial Y_j} X_{jr}^* - \frac{\gamma_j}{\gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} X_{j+1,r}^* \right) \quad (29)$$

$$= w \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} V_j^{-1} q_{jr} \quad (30)$$

where

$$V_j \triangleq \frac{\partial \gamma_j}{\partial \phi_j} \quad (31)$$

$$q_{jr} \triangleq \frac{\partial \gamma_j}{\partial Y_j} X_{jr}^* - \frac{\gamma_j}{\gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} X_{j+1,r}^*. \quad (32)$$

From (5), (6) and (17), we have,

$$\frac{\partial l}{\partial \phi_j} = n \left(Z_j - Z_{j+1} \frac{\partial g(\phi_j)}{\partial \phi_j} \right) \quad (33)$$

$$= n \left(Z_j - Z_{j+1} \frac{\gamma_j}{\gamma_{j+1}} \right). \quad (34)$$

Thus, using (24) and (32),

$$\frac{\partial^2 l}{\partial \beta_r^* \partial \phi_j} = -nw Z_{j+1} \frac{\gamma_{j+1} \frac{\partial \gamma_j}{\partial Y_j} X_{jr}^* - \gamma_j \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} X_{j+1,r}^*}{\gamma_{j+1}^2} \quad (35)$$

$$= -nw \frac{Z_{j+1}}{\gamma_{j+1}} q_{jr}. \quad (36)$$

Notice that V_j and q_{jr} are independent of β_s^* , $s \neq r$. Differentiating l twice, from (30) and (36),

$$\frac{\partial^2 l}{\partial \beta_r^* \partial \beta_s^*} = -nw^2 \sum_{j=1}^{k-1} V_j^{-1} \frac{Z_{j+1}}{\gamma_{j+1}} q_{jr} q_{js} \quad (37)$$

$$\Rightarrow E \left[\frac{\partial^2 l}{\partial \beta_r^* \partial \beta_s^*} \right] = -nw^2 \sum_{j=1}^{k-1} V_j^{-1} q_{jr} q_{js} \quad (38)$$

since by (4), $E[Z_{j+1}] = \gamma_{j+1}$.

Using (25), (11) and the chain rule in a similar manner, we obtain

$$\frac{\partial l}{\partial \tau_r} = U_r \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} \frac{\partial \phi_j}{\partial \gamma_j} \left(\frac{\partial \gamma_j}{\partial Y_j} Y_j - \frac{\gamma_j}{\gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} Y_{j+1} \right) \quad (39)$$

$$= U_r \sum_{j=1}^{k-1} \frac{\partial l}{\partial \phi_j} V_j^{-1} q_j \quad (40)$$

where we define

$$q_j \triangleq \frac{\partial \gamma_j}{\partial Y_j} Y_j - \frac{\gamma_j}{\gamma_{j+1}} \frac{\partial \gamma_{j+1}}{\partial Y_{j+1}} Y_{j+1}. \quad (41)$$

Thus, the expected second derivative turns out to be

$$E \left[\frac{\partial^2 l}{\partial \tau_r \partial \tau_s} \right] = -n U_r U_s \sum_{j=1}^{k-1} V_j^{-1} q_j^2. \quad (42)$$

Similarly, the mixed expected second derivatives are

$$E \left[\frac{\partial^2 l}{\partial \beta_r^* \partial \tau_s} \right] = -nw U_s \sum_{j=1}^{k-1} V_j^{-1} q_j q_{jr}. \quad (43)$$

Thus, the negative expectation of the Hessian matrix $\mathbf{A}_l = -E[\mathbf{H}_l]$ is symmetric and from (38) and (42), has nonnegative diagonal entries. Therefore, it is positive semidefinite.

Applying a Taylor series expansion for $\nabla_\psi l$ around ψ_n using the expected Hessian matrix, we get

$$\nabla l(\psi) = \nabla l(\psi_n) - \mathbf{A}_l(\psi - \psi_n) + \dots \quad (44)$$

Using a linear approximation, and taking $\nabla l(\psi) = \mathbf{0}$ at $\psi = \psi_{n+1}$, we get the update equation

$$\mathbf{A}_l(\psi_{n+1} - \psi_n) = \nabla l(\psi_n) \quad (45)$$

$$\Rightarrow \psi_{n+1} = \psi_n + \mathbf{A}_l^{-1} \nabla l(\psi_n). \quad (46)$$

where the entries of \mathbf{A}_l are defined in (38), (42) and (43).

REFERENCES

- [1] P. McCullagh, “Regression models for ordinal data,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 42, no. 2, pp. 109–142, 1980, Accessed Oct. 2, 2023. [Online]. Available: <http://www.jstor.org/stable/2984952>
- [2] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis, “Wine Quality,” UCI Machine Learning Repository, 2009, Accessed: Oct. 2, 2023. DOI: <https://doi.org/10.24432/C56S3T>.