

A FLEXIBLE APPROACH FOR EXAMINING CATEGORICAL DATA: ENHANCED LINEAR REGRESSION

Emily Jones

Division of Data Science, University of Toronto, Toronto, ON M5S 1A1, Canada

Abstract:

When analyzing ordinal data, particularly when the proportional odds assumption of the ordered logit model is not met, researchers often face challenges. The proportional odds assumption implies equal odds ratios across all outcome levels, which real-world data often does not satisfy. In such cases, some practitioners persist with the ordered logit model, while others turn to generalized linear models that may not be the most suitable choice. This paper introduces a novel methodology for calibrating category breakpoints, effectively removing the need for the equal distance assumption. By introducing M-1 breakpoints, the model gains additional degrees of freedom, offering a flexible solution for cases where the proportional odds assumption is untenable.

Keywords: ordinal data, ordered logit model, proportional odds assumption, generalized linear model, category breakpoints, odds ratio, and calibration.

Introduction

When the outcome variable is ordinal, meaning the outcome is ordered and discrete, the most popular model for the analysis is the ordered logit model. The outcome variable ($Y = Y_{(m)}, m = 1, 2, \dots, M$) takes M potential values. The ordered logit model requires data to meet the ordered odds assumption, which requires the odds ratio to stay the same for all the outcome levels. The odds ratio is defined as $P(Y > Y_{(m)})/P(Y \leq Y_{(m)}), m =$

$1, 2, \dots, M - 1$.

Unfortunately, the data in real life rarely meets the proportional odds assumption. Some practitioners choose to ignore the assumption and continue to use the model. Some practitioners choose to use the generalized linear model which is not necessarily suitable for the purpose. One of the main objections to the generalized linear model is that it requires the assumption that the numerical distance between each set of subsequent categories is equal, if we assign the M outcome variables naturally to ($Y = Y_{(m)} = m, m = 1, 2, \dots, M$).

In this paper, we propose a methodology to calibrate the breakpoint between the categories, which implicitly removes the equal distance assumption. The introduction of the M-1 breakpoints will add M-1 degree of freedom. Depending on the dataset, in many cases, those parameters are worth-while addition to the model.

In this work, we will illustrate the application of the calibration method with the National Bridge Inventory (NBI) data. The NBI condition rating is an important data source on bridge conditions nationwide. In our study, we apply the model to three outcome variables, Deck, Superstructure and Substructure. The independent variables are location, age, etc. We used the R packages and excel to perform the analysis.

In [1], Pan Lu, Hao Wang and Denver Tolliver compared the ordinary linear model with the ordinal logit model. The breakpoints of the ordinary linear model in [1] is not calibrated. They found the ordinal logit

model performed better. Predicting the future condition rating of bridges has been a topic for many researchers [2 - 13]; Many types of models has been experimented by the researchers in the literature, including straight-line extrapolation, linear regression, Markovian, nonlinear regression, logistic regression models, artificial neural networks, Bayesian network, Monte Carlo methods, and data mining-based algorithms.

To assess the overall performance of the model, we use three measures, the accuracy, Mean Absolute Error (MAE) and the balance. The accuracy rate is defined as percentage of accurate predictions. MAE measures the distance between true categories and predicted categories. Balance measures the overrating and underrating at the overall level. Ideally, the predicted categories reflect the true categories at the aggregate level without bias.

The Model and the Performance Metric

Let $S = \{(X_n, Y_n) | n = 1, 2, \dots, N\}$ be the training data set of N observations, where Y_n is the outcome variable for the observation n , and $X_n = (x_{n,1}, x_{n,2}, \dots, x_{n,p})$ is the input variable and $x_{n,j}$, $j = 1, 2, \dots, J$ are individual predictors.

The linear regression model is applied as usual. The $\tilde{Y}_n = L(X_n)$ be the (continuous) prediction. The usual way of discretized the projection is to have the break point in the middle of the outcome ($Y_{(m)}, m = 1, 2, \dots, M$), that is $(0.5 * (Y_{(j)} + Y_{(j+1)}), j = 1, 2, \dots, M - 1)$. We use the name Model1 as name of the model and \bar{Y}_i for the prediction of model.

$$\begin{aligned} \bar{Y}_i &= \text{Model1}(X_i) \\ &= \begin{cases} Y_{(1)}, & \text{if } L(X_i) < 0.5 * (Y_{(1)} + Y_{(2)}) \\ Y_{(j)}, & \text{if } 0.5 * (Y_{(j-1)} + Y_{(j)}) \leq L(X_i) < 0.5 * (Y_{(j)} + Y_{(j+1)}), j = 2, \dots, M - 1 \\ Y_{(M)}, & \text{if } 0.5 * (Y_{(M-1)} + Y_{(M)}) \leq L(X_i) \end{cases} \end{aligned}$$

We propose the following method to calibrate the breakpoint. The model with the new break point will be referred as Model2.

The following process is used to calibrate the breakpoints $\{B_{(m)}, m = 1, 2, \dots, M - 1\}$.

For each $j = 1, 2, \dots, M - 1$,

$\{\text{Count of } i \text{ such that } L(X_i) > B_{(m)}\} = \{\text{Count of } i \text{ such that } (Y_m) = Y_{(m)}\}$.

In practice, $B_{(j)}$ would be calculated recursively from $M - 1$ to 1. The new model, with calibrated break points, is

$$\begin{aligned} \bar{Y}_i &= \text{Model2}(X_i) \\ &= \begin{cases} Y_{(1)}, & \text{if } L(X_i) < B_{(1)} \\ Y_{(j)}, & \text{if } B_{(j-1)} \leq L(X_i) < B_{(j)}, j = 2, \dots, M - 1 \\ Y_{(M)}, & \text{if } B_{(M-1)} \leq L(X_i) \end{cases} \end{aligned}$$

For reference, we would use the $\bar{Y}_i = \text{Model3}(X_i)$ as the prediction of the third model, the ordered logit model. The theoretical definition of the model is readily available in the literature, and the implementation of the model is available in the R package. The R package used in this research is VGAM.

To access the overall performance of the models, we use the following three measures, accuracy, MAE and balance.

Let $\bar{Y}_i = M(X_i)$ be the projected outcome of the input variable X_i . We define

The Accuracy Rate = $\{\sum_{n=1}^N \delta(Y_n, \bar{Y}_n)\} / N$,

where $\delta(Y_n, \bar{Y}_n) = 1$ if $Y_n = \bar{Y}_n$ and $\delta(Y_n, \bar{Y}_n) = 0$ if $Y_n \neq \bar{Y}_n$.

$MAE = \{\sum_{n=1}^N Abs(Y_n - \bar{Y}_n)\}/N$.

And

$In-Balance = Abs\{\{\sum_{n=1}^N (\bar{Y}_n - Y_n) * (\bar{Y}_n > Y_n)\} - \{\sum_{n=1}^N (Y_n - \bar{Y}_n) * (\bar{Y}_n < Y_n)\}\}/N$

As defined above, a better model should have a higher Accuracy Rate, smaller MAE, and smaller In-Balance.

The Data set and the Result

We will use the NBI data to test and illustrate the methodology [14], [15]. NBI is the most comprehensive database on bridges in the United States. It has more than 100 fields in the database, collection information such as condition of the deck, structures, location, years build, traffic volume, and engineering attributes, such as span, high. In this study, we will study the condition of deck, superstructure and sub-structure

Table 1: Condition ratings used in the National Bridge Inventory (NBI)

Code	Meaning	Description
9	Excellent	As new
8	Very Good	No problems noted.
7	Good Some	Minor problems.
6	Satisfactory	Structural elements show some minor deterioration.
5	Fair	All primary structural elements are sound but may have minor section loss, cracking, spalling or scour
4	Poor	Advanced section loss, deterioration, spalling or scour.
3	Serious	Loss of section, deterioration, spalling or scour has seriously affected primary structural components. Local failures are possible. Fatigue cracks in steel or shear cracks in concrete may be present.
2	Critical	Advanced deterioration of primary structural elements. Fatigue cracks in steel or shear cracks in concrete may be present or scour may have removed substructure support. Unless closely monitored it may be necessary to close the bridge until corrective action is taken.
1	Imminent Failure	Major deterioration or section loss present in critical structural components or obvious vertical or horizontal movement affecting structure stability. The bridge is closed to traffic but with corrective action may put back in light service.
0	Failed	Out of service, beyond corrective action.

Source: United States Department of Transportation. Recording and Coding Guide for the Structure Inventory and Appraisal of the Nation's Bridges. Washington, D.C., 1995, page 38.

Among the fields in the NBI, we pick the following variables to be the explanatory variable. The age of the bridge and the annual daily traffic per lane and the bridge material type are critical variables. We also included age squared to capture if there is any convexity in the age variable.

Table 2: Description of variables used in analysis.

Name of variable	Description of variable
Reconstruction	Reconstruction record: Yes, No (binary variable)
Bridge Material Type	Structure materials: Steel, Concrete, Timber (dummy variable)
District Highway districts:	Bismarck, Devils Lake, Dickinson, Grand Forks, Minot, Valley City, Williston, Fargo (dummy Variable)
Age Bridge age:	Inspection year-construction year or inspection year-reconstruction year (continuous variable)
Age2	Bridge age squared (continuous variable)
ADT	Annual daily traffic per lane (continuous variable)

The Result

We performed the regression on three outputs, deck, superstructure and sub-structure with all 3 models. Model1 is the ordinary linear model with breakpoint in the middle. Model 2 is an ordinary linear model with calibrated breakpoints. Model 3 is the ordinal logit model. All three outputs take potential values of 0, 1, 2, ...9.

In Table 3, Table 4 and Table 5 below, we have the detailed cross table of comparison of the Model 1 and Model 2. The difference between the two models is the break points. In the table, the vertical counts are the original categories the horizontal count are the predicted categories. The diagonal would be the count the projected and the original categories did not change, that they are the same. Looking at the Model 2, in the total column and the total row, the count in each of the categories are very close, in most of the case, they are exactly the same. This is implied by the calibration method of Model 2.

Table 3: Cross table for Deck, predicted vs original, Model 1 and Model 2

Model 2	Original										
Deck Predicted	0	1	2	3	4	5	6	7	8	9	Total
0	0	0	0	0	3	1	0	1	0	0	5
1	0	0	0	0	1	0	0	0	0	0	1
2	0	0	1	1	2	3	1	0	0	0	8
3	0	0	0	2	4	12	2	1	0	0	21
4	0	0	3	4	24	36	27	3	3	0	100
5	0	1	1	10	41	132	132	46	6	0	369
6	1	0	3	3	16	109	245	211	49	4	641
7	0	0	0	2	11	73	217	448	147	9	907
8	0	0	0	0	1	3	16	186	327	45	578
9	0	0	0	0	0	0	1	11	46	69	127
Total	1	1	8	22	103	369	641	907	578	127	2757
Accuracy Rate 45% MAE 70% In-Balance 0.54%											
Model 1	Original										
Deck Predicted	0	1	2	3	4	5	6	7	8	9	Total
0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0

	4	0	0	0	0	0	0	0	0	0	0	0
	5	0	0	5	12	60	109	78	15	8	0	287
	6	1	1	3	8	34	200	364	292	62	4	969
	7	0	0	0	2	9	60	189	463	219	13	955
	8	0	0	0	0	0	0	10	137	275	83	505
	9	0	0	0	0	0	0	0	0	14	27	41
	Total	1	1	8	22	103	369	641	907	578	127	2757
Accuracy Rate 45% MAE 66% In-Balance 1.49%												

Table 4: Cross table for Superstructure, predicted vs original, Model 1 and Model 2

Model 2	Original											
Super Predicted	0	1	2	3	4	5	6	7	8	9	Total	
0	0	0	0	0	1	1	2	0	0	0	4	
1	0	0	0	0	1	0	0	0	0	0	1	
2	0	0	0	3	2	2	1	0	0	0	8	
3	0	0	2	1	6	10	3	0	0	0	22	
4	0	0	0	4	21	37	25	13	3	0	103	
5	0	1	4	10	37	115	136	54	12	0	369	
6	1	0	2	1	23	127	228	212	41	4	639	
7	0	0	0	3	12	73	219	437	149	13	906	
8	0	0	0	0	0	4	27	180	328	39	578	
9	0	0	0	0	0	0	0	11	45	71	127	
Total	1	1	8	22	103	369	641	907	578	127	2757	
Accuracy Rate 44% MAE 73% In-Balance 0.69%												
Model 1	Original											
Super Predicted	0	1	2	3	4	5	6	7	8	9	Total	
0	0	0	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	0	
4	0	0	2	4	15	16	5	4	1	0	47	
5	0	0	3	13	48	116	107	41	12	0	340	
6	1	1	3	3	31	178	313	287	53	4	874	
7	0	0	0	2	9	55	187	377	135	13	778	
8	0	0	0	0	0	4	29	191	349	69	642	
9	0	0	0	0	0	0	0	7	28	41	76	
Total	1	1	8	22	103	369	641	907	578	127	2757	
Accuracy Rate 44% MAE 69% In-Balance 3.48%												

Table 5: Cross table for Superstructure, predicted vs original, Model 1 and Model 2

Model 2	Original											
---------	----------	--	--	--	--	--	--	--	--	--	--	--

Sub	0	1	2	3	4	5	6	7	8	9	Total	
Predicted	0	0	0	0	2	2	0	1	0	0	5	
	1	0	0	0	1	0	0	0	0	0	1	
	2	0	0	0	3	3	1	1	0	0	8	
	3	0	0	1	1	3	11	3	2	0	21	
	4	0	0	2	3	23	37	25	8	2	100	
	5	0	1	3	10	38	108	137	63	9	369	
	6	1	0	2	2	22	128	227	210	45	641	
	7	0	0	0	3	11	79	222	434	146	907	
	8	0	0	0	0	0	3	26	184	332	578	
	9	0	0	0	0	0	0	0	5	44	78	127
	Total	1	1	8	22	103	369	641	907	578	127	2757
	Accuracy											
Rate				44%	MAE		73%	In-Balance		0.54%		
Model 1	Original											
Sub	0	1	2	3	4	5	6	7	8	9	Total	
Predicted	0	0	0	0	0	0	0	0	0	0	0	
	1	0	0	0	0	0	0	0	0	0	0	
	2	0	0	0	0	0	0	0	0	0	0	
	3	0	0	0	0	0	1	0	0	0	1	
	4	0	0	3	7	37	55	36	16	3	157	
	5	0	1	5	11	43	170	242	125	24	623	
	6	1	0	0	3	18	110	239	331	55	760	
	7	0	0	0	1	5	32	111	316	195	673	
	8	0	0	0	0	0	1	13	117	278	60	469

	9	0	0	0	0	0	0	2	23	49	74
Total	1	1	8	22	103	369	641	907	578	127	2757
Accuracy											
Rate 39% MAE 77% In-Balance 31.45%											

Table 6 below has the calibrated break point for each of the output, Deck, Superstructure and the sub structure. We see that the calibrated breakpoints tend to shift to the center. In Table 7, the summary performance is provided. Overall, ordinal logit model is best for super structure predication. The ordinary linear model performed better for deck and sub structure. And calibration of the breakpoint did not impact the performance much for all three put projection, but the in-balance is reduced by design.

Table 6: The Calibrated Breakpoints

	Between n 1 and 2	Between n 2 and 3	Between n 3 and 4	Between n 4 and 5	Between n 5 and 6	Between n 6 and 7	Between n 7 and 8	Between n 8 and 9	Between n 9 and
Model 1 Breakpoint , mid	0.5	1.5	2.5	3.5	4.5	5.5	6.5	7.5	8.5
Model 2 Breakpoint B(m)									
Deck	4.73	4.75	4.81	4.87	5.19	5.78	6.38	7.32	8.12
Super Structure	4.01	4.03	4.16	4.42	4.93	5.67	6.39	7.52	8.36
Sub Structure	3.48	3.58	3.75	3.97	4.44	5.12	5.98	7.26	8.35

Table 7: Summary performance measures

	Deck			Super structure			Sub structure		
	Model 1	Model 2	Model 3	Model 1	Model 2	Model 3	Model 1	Model 2	Model 3
Accuracy	45%	45%	21%	44%	44%	51%	39%	44%	22%
MAE	66%	70%	117%	69%	73%	62%	77%	73%	124%
In-Balance	1.49%	0.54%	110.26%	3.48%	0.69%	15.23%	31.45%	0.54%	115.67%

The Conclusion

Even for the ordinal data, the linear regression method would still provide a better fit of the data than the ordered logit model, depending on the nature of the data. The additional degree of freedom, which calibrates the breakpoint between the outcome categories, would improve the fit, but not necessarily always the case. Conceptually, the calibrated breakpoint provides additional flexibility. But as expected, the balance of the fit would improve, the balance measures the overall level of the overestimate and the underestimate of the fitted data. It would be preferable that the so that all over characteristics are preserved through the modeling process.

References

- Pan Lu, Hao Wang and Denver Tolliver, 'Prediction of Bridge Component Ratings Using Ordinal Logistic Regression Model', *Mathematical Problems in Engineering*, Volume 2019, Article ID 9797584, <https://doi.org/10.1155/2019/9797584>
- G. Bu, J. Lee, H. Guan, Y.-C. Loo, and M. Blumenstein, "Long-term performance of bridge elements using integrated deterioration method incorporating elman neural network," *Applied Mechanics and Materials*, vol. 204-208, pp. 1980–1987, 2012.
- O. Thomas and J. Sobanjo, "Comparison of markov chain and semi-markov models for crack deterioration on flexible pavements," *Journal of Infrastructure Systems*, vol. 19, no. 2, pp. 186–195, 2013.
- M. Alsharqawi, T. Zayed, and S. Abu Dabous, "Integrated condition rating and forecasting method for bridge decks using Visual Inspection and Ground Penetrating Radar," *Automation in Construction*, vol. 89, pp. 135–145, 2018.
- S. T. Ariaratnam, A. El-Assaly, and Y. Yang, "Assessment of infrastructure inspection needs using logistic models," *Journal of Infrastructure Systems*, vol. 7, no. 4, pp. 160–165, 2001.
- D. Tolliver and P. Lu, "Modeling bridge condition levels in the United States," *Journal of Civil Engineering and Architecture*, vol. 6, no. 4, pp. 415–432, 2012.
- S.M. Madanat, M. G. Karlaftis, and P. S. McCarthy, "Probabilistic infrastructure deterioration models with panel data," *Journal of Infrastructure Systems*, vol. 3, no. 1, pp. 4–9, 1997.
- Y.-H. Huang, "Artificial neural network model of bridge deterioration," *Journal of Performance of Constructed Facilities*, vol. 24, no. 6, pp. 597–602, 2010.
- J. Ruwanpura, S. T. Ariaratnam, and A. El-Assaly, "Prediction models for sewer infrastructure utilizing rule-based simulation," *Civil Engineering and Environmental Systems*, vol. 21, no. 3, pp. 169–185, 2004.
- O. Thomas and J. Sobanjo, "Semi-Markov models for the deterioration of bridge elements," *Journal of Infrastructure Systems*, vol. 22, no. 3, Article ID 04016010, 2016.

- G. Bu, J. H. Lee, H. Guan, Y. C. Loo, and M. Blumenstein, "Implementation of Elman neural networks for enhancing reliability of integrated bridge deterioration model," *Australian Journal of Structural Engineering*, vol. 15, no. 1, pp. 51–63, 2014.
- G. Nani, I. Mensah, and T. Adjei-Kumi, "Duration estimation model for bridge construction projects in Ghana," *Journal of Engineering, Design and Technology*, vol. 15, no. 6, pp. 754–777, 2017.
- M. Ahmed, O. Moselhi, and A. Bhowmick, "Integration of NDE measurements and current practice in bridge deterioration modeling," in *Proceedings of the 33rd International Symposium on Automation and Robotics in Construction, ISARC 2016*, pp. 341–349, USA, July 2016.
- FHWA, *Information Public Disclosure of National Bridge Inventory (NBI) Data*, 2007, <http://www.fhwa.dot.gov/bridge/nbi/20070517.cfm>.
- FHWA, *Recording and Coding Guide for The Structure Inventory and Appraisal of The Nation's Bridges*, Federal Highway Administration, U.S. Department of Transportation, 1995.