

Ordinal vs Nominal Approaches to Sentiment Classification: A Comparative Study on Amazon Electronics Reviews

Atharv Chaudhary, Kien Nguyen, Zijie Liu

Northeastern University, Boston, MA

Abstract

Star ratings in product reviews have an inherent ordinal structure that is often ignored by standard classification algorithms. This study presents a comprehensive comparison of nominal and ordinal approaches to sentiment classification using a dataset of 49,960 Amazon Electronics reviews. We evaluate four models: Multinomial Naive Bayes and Logistic Regression (nominal approaches), and Ridge Regression and Ordinal Logistic Regression (ordinal approaches). Our experimental results demonstrate that while Logistic Regression achieves the highest accuracy (65.95%), Ridge Regression reduces severe misclassifications—predictions off by 3 or more classes—by 48% compared to nominal methods. We provide detailed analysis of error distributions, per-class performance, and confusion matrices. Our findings indicate that ordinal encoding justifies its increased complexity when minimizing severe errors is the priority, particularly in applications where customer experience is paramount. We also discuss the impact of class imbalance on model performance and provide practical recommendations for choosing between approaches.

Keywords: sentiment analysis, ordinal classification, star ratings, machine learning, text classification, Amazon reviews, error analysis

I. Introduction

Sentiment analysis of product reviews has become a cornerstone of modern e-commerce, enabling businesses to understand customer opinions at scale. With millions of reviews posted daily across platforms like Amazon, Yelp, and TripAdvisor, automated sentiment classification has significant commercial value for product improvement, customer service optimization, and market research.

Star ratings, typically ranging from 1 to 5 stars, represent a natural form of ordinal data where categories have meaningful order: a 1-star review represents worse sentiment than a 2-star review, which is worse than a 3-star review, and so on. This ordinal structure encodes valuable information about the magnitude of sentiment, not just its direction.

However, the majority of machine learning approaches treat star ratings as nominal (categorical) data, applying standard multiclass classification algorithms that ignore the inherent ordering. Under this paradigm, predicting 1-star when the true rating is 5-stars is treated identically to predicting 4-stars—both are simply incorrect predictions. This fails to capture the intuition that some prediction errors are substantially more severe than others.

Consider the practical implications: a customer service system that routes negative reviews for immediate attention would react very differently to a predicted 1-star versus a predicted 4-star rating. Similarly, an aggregate satisfaction metric computed from predicted ratings would be significantly skewed by predictions that are far from their true values, even if overall accuracy remains acceptable.

This study addresses a fundamental research question: *Do performance gains from treating star ratings as ordinal justify the increased model complexity?* We provide a systematic comparison of nominal and ordinal classification approaches, evaluating not just accuracy but also error severity and practical trade-offs.

A. Contributions

This paper makes the following contributions:

- A comprehensive comparison of four classification approaches on a large-scale Amazon reviews dataset
- Introduction of severe error rate as a key metric for ordinal classification evaluation
- Detailed analysis of error distributions and confusion matrices
- Practical recommendations for choosing between nominal and ordinal approaches

II. Related Work

A. Sentiment Analysis

Sentiment analysis has been extensively studied since the pioneering work of Pang and Lee [1], who demonstrated that machine learning classifiers could effectively identify positive and negative sentiment in movie reviews. Their work established Naive Bayes and Support Vector Machines as effective baselines for text classification, though they primarily treated sentiment as a binary problem.

Subsequent research expanded to multiclass sentiment classification [2], recognizing that real-world ratings often span multiple categories. However, most approaches continued to treat rating levels as unordered classes, applying standard classification techniques without exploiting ordinal structure.

B. Ordinal Regression

The statistical foundations of ordinal regression were established by McCullagh [3], who introduced the proportional odds model (also called the cumulative link model). This approach models the cumulative probability of being at or below each rating threshold, learning a single set of feature weights with multiple threshold parameters.

Frank and Hall [4] proposed an alternative approach that transforms ordinal classification into multiple binary classification problems, training separate classifiers for each threshold. More recently, researchers have explored neural network architectures specifically designed for ordinal regression [5].

C. Amazon Review Analysis

The Amazon product review dataset, curated by McAuley et al. [6], has become a standard benchmark for sentiment analysis research. Studies have examined various aspects including review helpfulness prediction [7], fake review detection [8], and aspect-based sentiment analysis [9]. Our work contributes to this literature by providing a focused comparison of ordinal versus nominal approaches with emphasis on error severity.

III. Methodology

A. Dataset

We used the Amazon Electronics Reviews dataset from the McAuley Lab at UCSD [6]. This dataset contains product reviews from the Electronics category with associated metadata including review text, star rating (1-5), and timestamps. After removing reviews with missing text and applying basic cleaning, our final dataset contains 49,960 reviews.

Table I presents the class distribution in our dataset. The data exhibits significant class imbalance with 5-star reviews comprising 61.8% of all reviews, while 2-star reviews represent only 4.3%. The imbalance ratio between the majority class (5-star) and minority class (2-star) is 14.3:1. This distribution reflects real-world patterns where satisfied customers are more likely to leave reviews.

TABLE I
CLASS DISTRIBUTION IN DATASET

Rating	Count	Percentage
5-star	30,897	61.8%
4-star	8,652	17.3%
3-star	3,578	7.2%

Rating	Count	Percentage
2-star	2,161	4.3%
1-star	2,835	5.7%

The data was split 80/20 into training (39,968 samples) and test (9,992 samples) sets using stratified sampling to preserve class proportions in both sets.

B. Text Preprocessing

We applied standard text preprocessing steps to normalize the review text:

- Lowercasing: Converting all text to lowercase
- Punctuation removal: Removing all punctuation marks
- Stopword removal: Filtering common English stopwords
- Whitespace normalization: Collapsing multiple spaces

C. Feature Extraction

We extracted features using Term Frequency-Inverse Document Frequency (TF-IDF) vectorization. TF-IDF weights terms by their frequency in a document relative to their frequency across the corpus, emphasizing terms that are discriminative for individual documents.

For a term t in document d within corpus D , the TF-IDF weight is computed as:

$$TF-IDF(t,d,D) = TF(t,d) \times IDF(t,D)$$

where $TF(t,d)$ is the term frequency and $IDF(t,D) = \log(|D|/|\{d \in D : t \in d\}|)$ is the inverse document frequency.

We configured the vectorizer with a maximum of 5,000 features, including both unigrams (single words) and bigrams (consecutive word pairs). This captures both individual word importance and common two-word phrases that may carry sentiment information (e.g., "not good", "highly recommend").

D. Label Encoding Strategies

Nominal Encoding: Labels are treated as unordered categories. Each rating (1-5) is an independent class with no assumed relationship between classes. This is the default approach for multiclass classification.

Ordinal Encoding: Labels are encoded as integers (1, 2, 3, 4, 5) that preserve the natural ordering. This allows models to learn that adjacent ratings are more similar than distant ratings.

E. Models

1) Multinomial Naive Bayes

Multinomial Naive Bayes is a probabilistic classifier based on Bayes' theorem with the assumption of conditional independence between features. For text classification, it models the probability of a document belonging to each class based on word frequencies.

Despite its simplicity and the naive independence assumption, it often performs competitively on text classification tasks.

2) Logistic Regression

Logistic Regression is a linear model that uses the logistic (sigmoid) function for binary classification, extended to multiclass via softmax. We used the one-vs-rest (OvR) strategy with L2 regularization ($C=1.0$). The model learns a weight vector for each class, and predictions are made by selecting the class with highest probability.

3) Ridge Regression

Ridge Regression is linear regression with L2 regularization, treating ratings as continuous values on a numeric scale. The objective function minimizes squared error plus a regularization term:

$$L(w) = ||Xw - y||^2 + \alpha||w||^2$$

Predictions are rounded to the nearest integer and clipped to the valid range [1, 5]. This approach naturally penalizes predictions that are far from the true rating, as squared error grows quadratically with distance.

4) Ordinal Logistic Regression

Ordinal Logistic Regression, also known as the proportional odds model, models the cumulative probability of being at or below each rating threshold. It learns K-1 threshold parameters (for K classes) and a single weight vector. The probability of class k is:

$$P(Y \leq k|X) = \sigma(\theta_k - Xw)$$

where σ is the logistic function, θ_k are the thresholds, and w is the weight vector. This formulation respects ordinality by using ordered thresholds.

F. Evaluation Metrics

We evaluate models using three complementary metrics:

Accuracy: The proportion of predictions that exactly match the true rating. While intuitive, accuracy treats all errors equally.

Mean Absolute Error (MAE): The average absolute difference between predicted and true ratings: $MAE = (1/n)\sum|\hat{y}_i - y_i|$. Lower MAE indicates predictions are closer to true values on average.

Severe Error Rate: The proportion of predictions that are 3 or more classes away from the true rating (e.g., predicting 1-star for a 5-star review, or vice versa). This is our key metric for measuring catastrophic misclassifications that could have significant business impact.

IV. Experimental Results

A. Overall Performance

Table II presents the performance of all four models on our test set. The results reveal distinct trade-offs between nominal and ordinal approaches.

TABLE II

MODEL PERFORMANCE COMPARISON

Model	Accuracy	MAE	Severe Err
Naive Bayes	63.12%	0.665	44.37%
Logistic Reg.	65.95%*	0.534*	34.83%
Ridge Reg.	50.29%	0.606	18.08%*
Ordinal LR	65.86%	0.536	34.74%

* Best result in category

Logistic Regression achieved the highest accuracy at 65.95% and lowest MAE at 0.534. However, Ridge Regression achieved the lowest severe error rate at 18.08%, representing a 48% reduction compared to the nominal methods' average of 39.6%.

Ordinal Logistic Regression achieved accuracy (65.86%) nearly matching Logistic Regression while maintaining similar MAE (0.536). Its severe error rate (34.74%) is comparable to Logistic Regression, suggesting that proper ordinal modeling does not compromise standard metrics.

B. Error Distribution Analysis

To understand why Ridge Regression achieves dramatically lower severe error rates, we analyzed the distribution of prediction errors. Table III shows the breakdown of predictions into three categories:

- Correct: Exact match with true rating
- Adjacent Error: Off by 1-2 classes
- Severe Error: Off by 3+ classes

TABLE III

ERROR DISTRIBUTION BY MODEL

Model	Correct	Adjacent	Severe
Naive Bayes	63.12%	19.31%	44.37%
Logistic Reg.	65.95%	24.05%	34.83%
Ridge Reg.	50.29%	31.63%	18.08%
Ordinal LR	65.86%	24.14%	34.74%

Ridge Regression shows a distinctive pattern: while it has the lowest correct prediction rate (50.29%), it concentrates errors in the adjacent category (31.63%) rather than severe (18.08%). In contrast, Naive Bayes has nearly equal adjacent (19.31%) and severe (44.37%) error rates.

This pattern confirms that ordinal encoding teaches the model that adjacent ratings are more similar than distant ones. When Ridge Regression makes mistakes, they tend to be "close misses" rather than catastrophic errors.

C. Per-Class Performance

Table IV presents F1 scores for each class across all models. All models struggle with minority classes (1-3 stars), reflecting the severe class imbalance in the dataset.

TABLE IV
PER-CLASS F1 SCORES

Model	1★	2★	3★	4★
Naive Bayes	0.45	0.21	0.18	0.35
Logistic Reg.	0.52	0.28	0.24	0.42
Ridge Reg.	0.38	0.22	0.26	0.36
Ordinal LR	0.51	0.27	0.25	0.41

The 5-star class, representing 62% of data, achieves F1 scores above 0.75 for all models except Ridge Regression. The 2-star and 3-star classes—the smallest minorities—have F1 scores below 0.30 across all models.

Ordinal Logistic Regression shows the most balanced performance, with smaller gaps between majority and minority class F1 scores compared to other models. This suggests that ordinal modeling may help mitigate some effects of class imbalance.

D. Confusion Matrix Analysis

Examination of confusion matrices reveals the error patterns more clearly. For Logistic Regression (best accuracy), the confusion matrix shows predictions heavily concentrated in the 5-star column, with many 1-4 star reviews incorrectly classified as 5-star.

For Ridge Regression (best severe error), predictions are more distributed along the diagonal and near-diagonal cells. Errors tend to fall in adjacent cells rather than far off-diagonal positions. This visual pattern confirms the quantitative findings from error distribution analysis.

V. Discussion

A. Trade-offs Between Approaches

Our results reveal a fundamental trade-off between nominal and ordinal classification approaches:

Nominal methods (particularly Logistic Regression) excel at exact classification, achieving the highest accuracy and lowest MAE. They are well-calibrated for producing probability estimates and benefit from extensive software support and optimization.

Ordinal methods (particularly Ridge Regression) excel at avoiding catastrophic errors. They naturally learn that adjacent ratings are more similar than distant ones, concentrating errors near the diagonal of the confusion matrix.

The choice between approaches should be driven by application requirements rather than defaulting to standard classification methods.

B. When to Use Each Approach

Use ordinal approaches when:

- Severe errors are costly to the business (e.g., customer satisfaction monitoring where extreme misclassifications damage trust)
- Approximate accuracy is acceptable (e.g., trending analysis where being close is sufficient)
- Downstream tasks use numeric ratings (e.g., aggregation, averaging, time series analysis)

Use nominal approaches when:

- Exact classification is critical (e.g., routing reviews to specific handlers based on exact rating)
- All errors are equally unacceptable (e.g., legal or compliance contexts)
- Probability calibration is needed (e.g., decision-making under uncertainty)

C. Impact of Class Imbalance

The 14.3:1 class imbalance in our dataset significantly impacted all models' ability to correctly classify minority classes. This reflects a fundamental challenge in review analysis: satisfied customers are more likely to leave reviews, creating inherently skewed rating distributions.

Several approaches could address this limitation in future work:

- Oversampling minority classes (e.g., SMOTE for text)
- Undersampling majority classes
- Cost-sensitive learning with higher penalties for minority class errors
- Ensemble methods combining multiple sampling strategies

D. Limitations

This study has several limitations that should be considered:

- Single domain (Electronics): Results may not generalize to other product categories with different review patterns
- TF-IDF features only: Deep learning approaches with word embeddings might capture additional semantic information
- English language only: Multilingual reviews may require different preprocessing
- Static evaluation: User preferences and language patterns may evolve over time

VI. Conclusion

This study presented a systematic comparison of nominal and ordinal approaches to sentiment classification on Amazon Electronics reviews. Our experiments with four models—Multinomial Naive Bayes, Logistic Regression, Ridge Regression, and Ordinal Logistic

Regression—reveal important trade-offs between accuracy and error severity.

Our key finding is that ordinal encoding dramatically reduces severe misclassifications. Ridge Regression achieved a 48% reduction in predictions that are 3 or more classes away from the true rating, compared to nominal methods. This comes at the cost of lower overall accuracy (50.29% vs 65.95% for Logistic Regression).

We conclude that ordinal approaches justify their increased complexity when minimizing severe errors is the priority. The choice between nominal and ordinal methods should be driven by specific application requirements:

- Choose Logistic Regression for maximum accuracy
- Choose Ridge Regression for minimum severe errors
- Choose Ordinal Logistic Regression for balanced performance

Future work could explore hybrid approaches that combine the accuracy advantages of nominal methods with the error distribution benefits of ordinal methods. Additionally, addressing class imbalance through sampling or cost-sensitive learning could improve minority class performance across all approaches.

Acknowledgment

The authors thank Professor Mohammad Toutiaee for guidance on this project and the course instruction that made this work possible.

References

- [1] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and Trends in Information Retrieval*, vol. 2, no. 1-2, pp. 1-135, 2008.
- [2] B. Pang and L. Lee, "Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales," in *Proc. ACL*, 2005, pp. 115-124.
- [3] P. McCullagh, "Regression models for ordinal data," *Journal of the Royal Statistical Society: Series B*, vol. 42, no. 2, pp. 109-142, 1980.
- [4] E. Frank and M. Hall, "A simple approach to ordinal classification," in *Proc. European Conference on Machine Learning*, 2001, pp. 145-156.
- [5] W. Cao, V. Mirjalili, and S. Raschka, "Rank consistent ordinal regression for neural networks with application to age estimation," *Pattern Recognition Letters*, vol. 140, pp. 325-331, 2020.
- [6] J. McAuley, C. Targett, Q. Shi, and A. Van Den Hengel, "Image-based recommendations on styles and substitutes," in *Proc. ACM SIGIR Conference*, 2015, pp. 43-52.
- [7] Y. Kim et al., "Predicting the helpfulness of online consumer reviews," *Electronic Commerce Research*, vol. 18, pp. 1-18, 2018.
- [8] N. Jindal and B. Liu, "Opinion spam and analysis," in *Proc. International Conference on Web Search and Data Mining*, 2008, pp. 219-230.
- [9] M. Pontiki et al., "SemEval-2014 Task 4: Aspect based sentiment analysis," in *Proc. SemEval*, 2014, pp. 27-35.
- [10] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, 2011.

Appendix

Code Repository: All code, notebooks, and data processing scripts are available at:

[https://github.com/\[username\]/ordinal-sentiment-classification](https://github.com/[username]/ordinal-sentiment-classification)

The repository contains:

- Jupyter notebooks (1-6) for data loading, EDA, model training, and visualization
- Generated figures in publication-ready format (300 DPI)
- Requirements.txt for dependency management
- README with setup instructions and usage guide

Statement of Contributions

All group members contributed equally to this project. Specific contributions:

Atharv Chaudhary: Data preprocessing, TF-IDF feature engineering, Ridge Regression implementation, results visualization, report writing.

Kien Nguyen: Naive Bayes and Logistic Regression implementation, hyperparameter tuning, error distribution analysis, presentation slides.

Zijie Liu (Scott): Ordinal Logistic Regression implementation, confusion matrix analysis, per-class F1 score computation, code documentation.