

Notebook 4: Ordinal Models Final Project - Ordinal vs Nominal Sentiment Analysis

Atharv Chaudhary

Purpose: Train and evaluate ORDINAL classification models.

Models:

1. Ridge Regression (treats ratings as continuous)
2. Ordinal Logistic Regression (threshold-based)

Input: `amazon_electronics_cleaned.csv`

Output: `ordinal_results.csv`, confusion matrices

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount('/content/drive', force_remount=True)`

```
1 # Import libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 from sklearn.model_selection import train_test_split
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 from sklearn.linear_model import Ridge, LogisticRegression
12 from sklearn.metrics import (
13     accuracy_score,
14     mean_absolute_error,
15     f1_score,
16     classification_report,
17     confusion_matrix
18 )
19
20 # Settings
21 RANDOM_STATE = 42
22 np.random.seed(RANDOM_STATE)
23 plt.style.use('seaborn-v0_8-whitegrid')
24
25 print("✅ Libraries imported")
```

✅ Libraries imported

✓ Step 1: Load Data & Prepare Features

```

1 # Load cleaned data
2 df = pd.read_csv('/content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligence/reviews.csv')
3 print(f"✓ Loaded {len(df):,} reviews")
4
5 # TF-IDF
6 vectorizer = TfidfVectorizer(
7     max_features=5000,
8     stop_words='english',
9     ngram_range=(1, 2),
10    min_df=5,
11    max_df=0.95
12 )
13
14 X = vectorizer.fit_transform(df['text'])
15 y = df['rating'].values
16
17 # Train/Test Split
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.2, random_state=RANDOM_STATE, stratify=y
20 )
21
22 print(f"✓ Training: {X_train.shape[0]:,} | Test: {X_test.shape[0]:,}")

```

✓ Loaded 49,960 reviews
 ✓ Training: 39,968 | Test: 9,992

✓ Step 2: Helper Functions

```

1 def evaluate_model(y_true, y_pred, model_name):
2     """Evaluate model and return metrics."""
3     accuracy = accuracy_score(y_true, y_pred)
4     mae = mean_absolute_error(y_true, y_pred)
5     f1_macro = f1_score(y_true, y_pred, average='macro')
6     f1_weighted = f1_score(y_true, y_pred, average='weighted')
7
8     print(f"\n{'='*55}")
9     print(f"🌈 {model_name}")
10    print(f"{'='*55}")
11    print(f"Accuracy:      {accuracy:.4f} ({accuracy*100:.2f}%)")
12    print(f"MAE:           {mae:.4f}")
13    print(f"F1 (macro):     {f1_macro:.4f}")
14    print(f"F1 (weighted): {f1_weighted:.4f}")
15
16    return {

```

```

17         'model': model_name,
18         'encoding': 'Ordinal',
19         'accuracy': accuracy,
20         'mae': mae,
21         'f1_macro': f1_macro,
22         'f1_weighted': f1_weighted
23     }
24
25
26 def calculate_error_rates(y_true, y_pred):
27     """Calculate adjacent and severe error rates."""
28     errors = y_true != y_pred
29     if errors.sum() == 0:
30         return 0.0, 0.0
31
32     error_distances = np.abs(y_true[errors] - y_pred[errors])
33     adjacent = (error_distances == 1).sum() / errors.sum()
34     severe = (error_distances >= 2).sum() / errors.sum()
35
36     return adjacent, severe
37
38

```

✓ Helper functions defined

✓ Step 3: Model 1 - Ridge Regression (Ordinal)

```

1 # =====
2 # MODEL 1: RIDGE REGRESSION (ORDINAL)
3 # =====
4
5 print("\n" + "="*70)
6 print("🔧 MODEL 1: Ridge Regression (Ordinal)")
7 print("="*70)
8 print("\nTreats ratings as CONTINUOUS ordinal values: 1 < 2 < 3 < 4 < 5")
9 print("Formula:  $\hat{y} = w^T x + b$ , then round to nearest integer")
10 print("Loss:  $\sum (y_i - \hat{y}_i)^2 + \lambda ||w||^2$  (naturally penalizes large errors more)")
11
12 # Train
13 ridge_model = Ridge(alpha=1.0, random_state=RANDOM_STATE)
14 ridge_model.fit(X_train, y_train)
15
16 # Predict (continuous) then round and clip
17 ridge_pred_continuous = ridge_model.predict(X_test)
18 ridge_pred = np.clip(np.round(ridge_pred_continuous), 1, 5).astype(int)
19
20 # Evaluate
21 ridge_results = evaluate_model(y_test, ridge_pred, "Ridge Regression")
22

```

```

23 # Error analysis
24 ridge_adjacent, ridge_severe = calculate_error_rates(y_test, ridge_pred)
25 ridge_results['adjacent_error'] = ridge_adjacent
26 ridge_results['severe_error'] = ridge_severe
27
28 print(f"\nError Analysis:")
29 print(f"    Adjacent Error Rate (±1): {ridge_adjacent:.2%}")
30 print(f"    Severe Error Rate (±2+): {ridge_severe:.2%}")

```

```

=====
🔧 MODEL 1: Ridge Regression (Ordinal)
=====

```

Treats ratings as CONTINUOUS ordinal values: $1 < 2 < 3 < 4 < 5$
 Formula: $\hat{y} = w^T x + b$, then round to nearest integer
 Loss: $\sum (y_i - \hat{y}_i)^2 + \lambda ||w||^2$ (naturally penalizes large errors more)

```

=====
📊 Ridge Regression
=====

```

```

Accuracy:      0.5029 (50.29%)
MAE:           0.6055
F1 (macro):    0.3063
F1 (weighted): 0.5244

```

Error Analysis:
 Adjacent Error Rate (±1): 81.92%
 Severe Error Rate (±2+): 18.08%

```

1 # Show continuous prediction distribution
2 print("\n📊 Ridge Continuous Predictions:")
3 print(f"    Min: {ridge_pred_continuous.min():.2f}")
4 print(f"    Max: {ridge_pred_continuous.max():.2f}")
5 print(f"    Mean: {ridge_pred_continuous.mean():.2f}")
6
7 # Classification report
8 print("\n📄 Classification Report - Ridge:")
9 print(classification_report(y_test, ridge_pred, digits=4))

```

```

📊 Ridge Continuous Predictions:
Min: -0.19
Max: 6.23
Mean: 4.29

```

```

📄 Classification Report - Ridge:

```

	precision	recall	f1-score	support
1	0.7222	0.0688	0.1256	567
2	0.1756	0.1065	0.1326	432
3	0.2135	0.2762	0.2408	793
4	0.2606	0.5359	0.3507	2021
5	0.8093	0.5888	0.6817	6179

accuracy			0.5029	9992
macro avg	0.4362	0.3152	0.3063	9992
weighted avg	0.6187	0.5029	0.5244	9992

✓ Step 4: Model 2 - Ordinal Logistic Regression

```

1 # =====
2 # ORDINAL LOGISTIC REGRESSION CLASS
3 # =====
4
5 class OrdinalLogisticRegression:
6     """
7     Ordinal Logistic Regression using threshold (cumulative) approach.
8
9     For K classes, trains K-1 binary classifiers.
10    Each classifier k models:  $P(Y \leq k \mid x) = \sigma(\theta_k - w^T x)$ 
11
12    This EXPLICITLY preserves ordinal structure:  $1 < 2 < 3 < 4 < 5$ 
13
14    Key insight: Same feature weights  $w$  for all thresholds,
15    only the threshold  $\theta_k$  changes.
16    """
17
18    def __init__(self, max_iter=1000, C=1.0):
19        self.max_iter = max_iter
20        self.C = C
21        self.classifiers = []
22        self.classes_ = None
23
24    def fit(self, X, y):
25        self.classes_ = np.sort(np.unique(y))
26        n_classes = len(self.classes_)
27
28        print(f"    Training {n_classes - 1} binary classifiers...")
29        print(f"    Classes: {self.classes_}")
30
31        # Train K-1 binary classifiers for cumulative probabilities
32        for k in range(n_classes - 1):
33            threshold = self.classes_[k]
34            # Binary target: 1 if y <= threshold, 0 otherwise
35            y_binary = (y <= threshold).astype(int)
36
37            clf = LogisticRegression(
38                max_iter=self.max_iter,
39                C=self.C,
40                random_state=42,
41                solver='lbfgs'

```

```

42         )
43         clf.fit(X, y_binary)
44         self.classifiers.append(clf)
45
46         pos_rate = y_binary.mean()
47         print(f"      Classifier {k+1}: P(Y ≤ {threshold}) - {pos_rate}")
48
49     return self
50
51     def predict_proba(self, X):
52         """Predict class probabilities using cumulative approach."""
53         n_samples = X.shape[0]
54         n_classes = len(self.classes_)
55
56         # Get cumulative probabilities P(Y ≤ k)
57         cumulative_probs = np.zeros((n_samples, n_classes))
58         cumulative_probs[:, -1] = 1.0 # P(Y ≤ max_class) = 1
59
60         for k, clf in enumerate(self.classifiers):
61             cumulative_probs[:, k] = clf.predict_proba(X)[:, 1]
62
63         # Convert cumulative to class probabilities: P(Y = k) = P(Y ≤ k) - P(Y ≤ k-1)
64         class_probs = np.zeros((n_samples, n_classes))
65         class_probs[:, 0] = cumulative_probs[:, 0] # P(Y = 1) = P(Y ≤ 1)
66
67         for k in range(1, n_classes):
68             class_probs[:, k] = cumulative_probs[:, k] - cumulative_probs[:, k-1]
69
70         # Ensure non-negative (numerical stability)
71         class_probs = np.maximum(class_probs, 0)
72
73         # Normalize rows to sum to 1
74         row_sums = class_probs.sum(axis=1, keepdims=True)
75         row_sums[row_sums == 0] = 1 # Avoid division by zero
76         class_probs = class_probs / row_sums
77
78         return class_probs
79
80     def predict(self, X):
81         """Predict class labels."""
82         probs = self.predict_proba(X)
83         return self.classes_[np.argmax(probs, axis=1)]
84
85

```

✓ OrdinalLogisticRegression class defined

```

1 # =====
2 # MODEL 2: ORDINAL LOGISTIC REGRESSION
3 # =====

```

```

4
5 print("\n" + "="*70)
6 print("🔧 MODEL 2: Ordinal Logistic Regression (Threshold-Based)")
7 print("="*70)
8 print("\nUses K-1 binary classifiers to model cumulative probabilities.")
9 print("Formula:  $P(Y \leq k \mid x) = \sigma(\theta_k - w^T x)$ ")
10 print("Class prob:  $P(Y = k) = P(Y \leq k) - P(Y \leq k-1)$ ")
11 print("\nExplicitly encodes ordinal constraint:  $\theta_1 < \theta_2 < \theta_3 < \theta_4$ \n")
12
13 # Train
14 olr_model = OrdinalLogisticRegression(max_iter=1000, C=1.0)
15 olr_model.fit(X_train, y_train)
16
17 # Predict
18 print("\n  Making predictions...")
19 olr_pred = olr_model.predict(X_test)
20
21 # Evaluate
22 olr_results = evaluate_model(y_test, olr_pred, "Ordinal Logistic Regression")
23
24 # Error analysis
25 olr_adjacent, olr_severe = calculate_error_rates(y_test, olr_pred)
26 olr_results['adjacent_error'] = olr_adjacent
27 olr_results['severe_error'] = olr_severe
28
29 print(f"\nError Analysis:")
30 print(f"  Adjacent Error Rate ( $\pm 1$ ): {olr_adjacent:.2%}")

```

```

=====
🔧 MODEL 2: Ordinal Logistic Regression (Threshold-Based)
=====

```

Uses K-1 binary classifiers to model cumulative probabilities.

Formula: $P(Y \leq k \mid x) = \sigma(\theta_k - w^T x)$

Class prob: $P(Y = k) = P(Y \leq k) - P(Y \leq k-1)$

Explicitly encodes ordinal constraint: $\theta_1 < \theta_2 < \theta_3 < \theta_4$

Training 4 binary classifiers...

Classes: [1 2 3 4 5]

Classifier 1: $P(Y \leq 1)$ - 5.7% positive

Classifier 2: $P(Y \leq 2)$ - 10.0% positive

Classifier 3: $P(Y \leq 3)$ - 17.9% positive

Classifier 4: $P(Y \leq 4)$ - 38.2% positive

Making predictions...

```

=====
📊 Ordinal Logistic Regression
=====
Accuracy:      0.6586 (65.86%)
MAE:           0.5360

```

```
F1 (macro):    0.3702
F1 (weighted): 0.6049
```

Error Analysis:

```
Adjacent Error Rate ( $\pm 1$ ): 65.26%
Severe Error Rate ( $\pm 2+$ ): 34.74%
```

```
1 # Classification report
2 print("\n 📄 Classification Report - Ordinal LR:")
3 print(classification_report(y_test, olr_pred, digits=4))
```

```
📄 Classification Report - Ordinal LR:
```

	precision	recall	f1-score	support
1	0.6246	0.3316	0.4332	567
2	0.2481	0.0741	0.1141	432
3	0.3248	0.1438	0.1993	793
4	0.4318	0.2256	0.2964	2021
5	0.7101	0.9372	0.8080	6179
accuracy			0.6586	9992
macro avg	0.4679	0.3424	0.3702	9992
weighted avg	0.5984	0.6586	0.6049	9992

✓ Step 5: Confusion Matrices

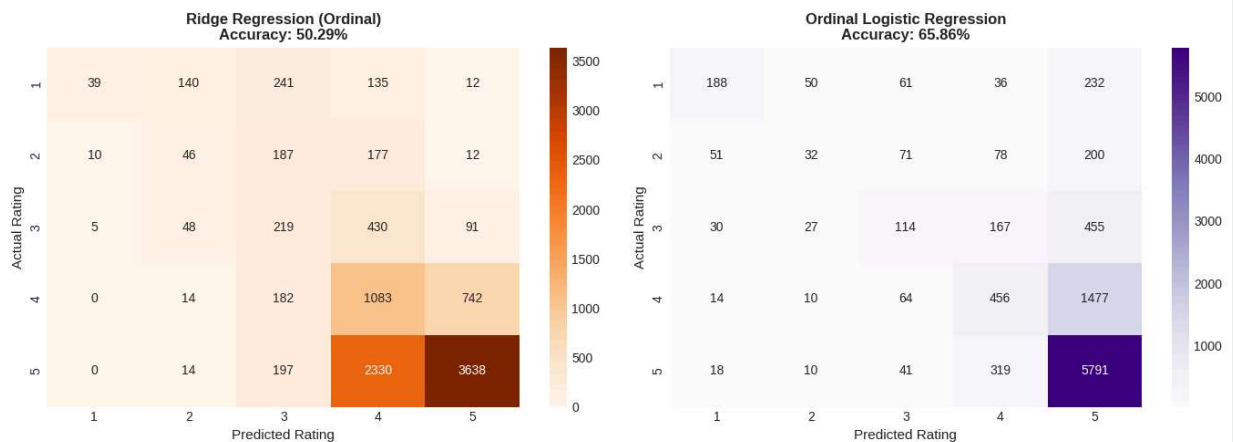
```
1 # =====
2 # CONFUSION MATRICES
3 # =====
4
5 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
6
7 # Ridge Regression
8 cm_ridge = confusion_matrix(y_test, ridge_pred)
9 sns.heatmap(cm_ridge, annot=True, fmt='d', cmap='Oranges', ax=axes[0],
10             xticklabels=[1, 2, 3, 4, 5], yticklabels=[1, 2, 3, 4, 5])
11 axes[0].set_xlabel('Predicted Rating', fontsize=11)
12 axes[0].set_ylabel('Actual Rating', fontsize=11)
13 axes[0].set_title(f'Ridge Regression (Ordinal)\nAccuracy: {ridge_results[']
14                  fontsize=12, fontweight='bold')
15
16 # Ordinal Logistic Regression
17 cm_olr = confusion_matrix(y_test, olr_pred)
18 sns.heatmap(cm_olr, annot=True, fmt='d', cmap='Purples', ax=axes[1],
19             xticklabels=[1, 2, 3, 4, 5], yticklabels=[1, 2, 3, 4, 5])
20 axes[1].set_xlabel('Predicted Rating', fontsize=11)
21 axes[1].set_ylabel('Actual Rating', fontsize=11)
22 axes[1].set_title(f'Ordinal Logistic Regression\nAccuracy: {olr_results[']
```



```

23             fontsize=12, fontweight='bold')
24
25 plt.tight_layout()
26 plt.savefig('confusion_matrices_ordinal.png', dpi=150, bbox_inches='tight')
27 plt.show()
28
29 print("\n✅ Saved: confusion_matrices_ordinal.png")

```



✅ Saved: confusion_matrices_ordinal.png

✓ Step 6: Save Results

```

1 # =====
2 # SAVE RESULTS
3 # =====
4
5 # Combine results
6 ordinal_results = pd.DataFrame([ridge_results, olr_results])
7
8 print("\n" + "="*70)
9 print("📊 ORDINAL MODELS SUMMARY")
10 print("="*70)
11 print(ordinal_results.to_string(index=False))
12
13 # Save to CSV
14 ordinal_results.to_csv('ordinal_results.csv', index=False)
15 print("\n✅ Saved: ordinal_results.csv")

```

=====

📊 ORDINAL MODELS SUMMARY

=====

model	encoding	accuracy	mae	f1_macro	f1_weighted
-------	----------	----------	-----	----------	-------------

Ridge Regression	Ordinal	0.502902	0.605484	0.306266	0.524437
Ordinal Logistic Regression	Ordinal	0.658627	0.536029	0.370193	0.604948

✓ Saved: ordinal_results.csv

```

1 # Save predictions
2 predictions_df = pd.DataFrame({
3     'actual': y_test,
4     'ridge_pred': ridge_pred,
5     'olr_pred': olr_pred
6 })
7 predictions_df.to_csv('ordinal_predictions.csv', index=False)
8 print("✓ Saved: ordinal_predictions.csv")

```

✓ Saved: ordinal_predictions.csv

```

1 # Download files
2 try:
3     from google.colab import files
4     files.download('ordinal_results.csv')
5     files.download('confusion_matrices_ordinal.png')
6 except:
7     print("Files saved locally")

```

```

1 import shutil
2 import os
3
4 drive_path = '/content/drive/MyDrive/Fall 2025/Foundations of Artificial
Intelligence/Final Project/data'
5
6 # Ensure the directory exists
7 os.makedirs(drive_path, exist_ok=True)
8
9 files_to_save = [
10     'ordinal_results.csv',
11     'ordinal_predictions.csv',
12     'confusion_matrices_ordinal.png'
13 ]
14
15 for file_name in files_to_save:
16     try:
17         shutil.copy(file_name, os.path.join(drive_path, file_name))
18         print(f"✓ Saved '{file_name}' to Google Drive")
19     except FileNotFoundError:
20         print(f"✗ Error: '{file_name}' not found. Skipping.")
21     except Exception as e:
22         print(f"✗ Error saving '{file_name}' to Google Drive: {e}")

```

- ✓ Saved 'ordinal_results.csv' to Google Drive
- ✓ Saved 'ordinal_predictions.csv' to Google Drive
- ✓ Saved 'confusion_matrices_ordinal.png' to Google Drive

✓ Summary

Ordinal Models Trained:

Model	Accuracy	MAE	Adjacent Error	Severe Error
Ridge Regression	See above	See above	See above	See above
Ordinal LR	See above	See above	See above	See above

Key Insight: Ordinal models should have LOWER severe error rates because they understand that $1 \rightarrow 5$ is worse than $4 \rightarrow 5$.

Next: Run `5_Results_Analysis.ipynb` to compare all models.