

Notebook 3: Nominal Models Final Project - Ordinal vs Nominal Sentiment Analysis

Atharv Chaudhary

Purpose: Train and evaluate NOMINAL classification models.

Models:

1. Multinomial Naive Bayes
2. Logistic Regression (Multinomial)

Input: `amazon_electronics_cleaned.csv`

Output: `nominal_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.naive_bayes import MultinomialNB
12 from sklearn.linear_model import LogisticRegression
13 from sklearn.metrics import (
14     accuracy_score,
15     mean_absolute_error,
16     f1_score,
17     classification_report,
18     confusion_matrix
19 )
20
21 # Settings
22 RANDOM_STATE = 42
23 np.random.seed(RANDOM_STATE)
24 plt.style.use('seaborn-v0_8-whitegrid')
25
26 print("✅ Libraries imported")
```

✅ Libraries imported

Step 1: Load Data

```

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 # Show distribution
6 print("\n📊 Rating Distribution:")
7 print(df['rating'].value_counts().sort_index())

```

✅ Loaded 49,960 reviews

📊 Rating Distribution:

rating	count
1	2835
2	2161
3	3964
4	10103
5	30897

Name: count, dtype: int64

Step 2: Feature Extraction (TF-IDF)

```

1 # =====
2 # TF-IDF FEATURE EXTRACTION
3 # =====
4
5 print("=" * 70)
6 print("TF-IDF FEATURE EXTRACTION")
7 print("=" * 70)
8
9 # Configuration
10 MAX_FEATURES = 5000
11 NGRAM_RANGE = (1, 2) # Unigrams and bigrams
12
13 vectorizer = TfidfVectorizer(
14     max_features=MAX_FEATURES,
15     stop_words='english',
16     ngram_range=NGRAM_RANGE,
17     min_df=5,
18     max_df=0.95
19 )
20
21 print(f"\nSettings:")
22 print(f"    Max features: {MAX_FEATURES}")
23 print(f"    N-gram range: {NGRAM_RANGE}")
24

```

```

25 # Transform
26 X = vectorizer.fit_transform(df['text'])
27 y = df['rating'].values
28
29 print(f"\n✅ TF-IDF matrix: {X.shape}")

```

```

=====
TF-IDF FEATURE EXTRACTION
=====

```

Settings:

Max features: 5000
N-gram range: (1, 2)

✅ TF-IDF matrix: (49960, 5000)

Step 3: Train/Test Split

```

1 # =====
2 # TRAIN/TEST SPLIT
3 # =====
4
5 X_train, X_test, y_train, y_test = train_test_split(
6     X, y,
7     test_size=0.2,
8     random_state=RANDOM_STATE,
9     stratify=y
10 )
11
12 print(f"✅ Training set: {X_train.shape[0]:,} samples")
13 print(f"✅ Test set: {X_test.shape[0]:,} samples")

```

✅ Training set: 39,968 samples
✅ Test set: 9,992 samples

Step 4: Helper Functions

```

1 # =====
2 # HELPER FUNCTIONS
3 # =====
4
5 def evaluate_model(y_true, y_pred, model_name):

```

```

6     """Evaluate model and return metrics."""
7     accuracy = accuracy_score(y_true, y_pred)
8     mae = mean_absolute_error(y_true, y_pred)
9     f1_macro = f1_score(y_true, y_pred, average='macro')
10    f1_weighted = f1_score(y_true, y_pred, average='weighted')
11
12    print(f"\n{'='*55}")
13    print(f"🌈 {model_name}")
14    print(f"{'='*55}")
15    print(f"Accuracy:      {accuracy:.4f} ({accuracy*100:.2f}%)")
16    print(f"MAE:           {mae:.4f}")
17    print(f"F1 (macro):     {f1_macro:.4f}")
18    print(f"F1 (weighted): {f1_weighted:.4f}")
19
20    return {
21        'model': model_name,
22        'encoding': 'Nominal',
23        'accuracy': accuracy,
24        'mae': mae,
25        'f1_macro': f1_macro,
26        'f1_weighted': f1_weighted
27    }
28
29
30 def calculate_error_rates(y_true, y_pred):
31     """Calculate adjacent and severe error rates."""
32     errors = y_true != y_pred
33     if errors.sum() == 0:
34         return 0.0, 0.0
35
36     error_distances = np.abs(y_true[errors] - y_pred[errors])
37     adjacent = (error_distances == 1).sum() / errors.sum()
38     severe = (error_distances >= 2).sum() / errors.sum()
39
40     return adjacent, severe
41
42
43 print("✅ Helper functions defined")

```

✅ Helper functions defined

Step 5: Model 1 - Multinomial Naive Bayes

```

1 # =====
2 # MODEL 1: MULTINOMIAL NAIVE BAYES
3 # =====

```

```

4
5 print("\n" + "="*70)
6 print("🔧 MODEL 1: Multinomial Naive Bayes (Nominal)")
7 print("="*70)
8 print("\nTreats classes as UNORDERED categories.")
9 print("Formula:  $P(Y=k|x) \propto P(Y=k) \times \prod P(x_j|Y=k)$ ")
10
11 # Train
12 nb_model = MultinomialNB(alpha=1.0) # Laplace smoothing
13 nb_model.fit(X_train, y_train)
14
15 # Predict
16 nb_pred = nb_model.predict(X_test)
17
18 # Evaluate
19 nb_results = evaluate_model(y_test, nb_pred, "Naive Bayes")
20
21 # Error analysis
22 nb_adjacent, nb_severe = calculate_error_rates(y_test, nb_pred)
23 nb_results['adjacent_error'] = nb_adjacent
24 nb_results['severe_error'] = nb_severe
25
26 print(f"\nError Analysis:")
27 print(f"    Adjacent Error Rate ( $\pm 1$ ): {nb_adjacent:.2%}")

```

```

=====
🔧 MODEL 1: Multinomial Naive Bayes (Nominal)
=====

```

Treats classes as UNORDERED categories.
 Formula: $P(Y=k|x) \propto P(Y=k) \times \prod P(x_j|Y=k)$

```
=====
```

📊 Naive Bayes

```
=====
```

Accuracy: 0.6312 (63.12%)
 MAE: 0.6651
 F1 (macro): 0.2321
 F1 (weighted): 0.5129


Error Analysis:

Adjacent Error Rate (± 1): 55.63%
 Severe Error Rate ($\pm 2+$): 44.37%

```

1 # Classification report
2 print("\n📄 Classification Report - Naive Bayes:")
3 print(classification_report(y_test, nb_pred, digits=4))

```

 Classification Report - Naive Bayes:

	precision	recall	f1-score	support
1	0.6236	0.1958	0.2980	567
2	0.0000	0.0000	0.0000	432
3	0.3750	0.0038	0.0075	793
4	0.3797	0.0445	0.0797	2021
5	0.6378	0.9877	0.7751	6179
accuracy			0.6312	9992
macro avg	0.4032	0.2464	0.2321	9992
weighted avg	0.5364	0.6312	0.5129	9992

Step 6: Model 2 - Logistic Regression (Multinomial)

```

1 # =====
2 # MODEL 2: LOGISTIC REGRESSION (MULTINOMIAL)
3 # =====
4
5 print("\n" + "="*70)
6 print("🔧 MODEL 2: Logistic Regression (Nominal - Multinomial)")
7 print("="*70)
8 print("\nUses softmax, treats classes as UNORDERED.")
9 print("Formula:  $P(Y=k|x) = \frac{\exp(w_k^T x + b_k)}{\sum \exp(w_j^T x + b_j)}$ ")
10
11 # Train
12 lr_model = LogisticRegression(
13     multi_class='multinomial',
14     solver='lbfgs',
15     max_iter=1000,
16     random_state=RANDOM_STATE,
17     n_jobs=-1
18 )
19 lr_model.fit(X_train, y_train)
20
21 # Predict
22 lr_pred = lr_model.predict(X_test)
23
24 # Evaluate
25 lr_results = evaluate_model(y_test, lr_pred, "Logistic Regression")
26
27 # Error analysis
28 lr_adjacent, lr_severe = calculate_error_rates(y_test, lr_pred)
29 lr_results['adjacent_error'] = lr_adjacent
30 lr_results['severe_error'] = lr_severe
31
32 print(f"\nError Analysis:")
33 print(f"    Adjacent Error Rate ( $\pm 1$ ): {lr_adjacent:.2%}")
34 print(f"    Severe Error Rate ( $\pm 2$ ): {lr_severe:.2%}")

```

```

=====
🔧 MODEL 2: Logistic Regression (Nominal - Multinomial)
=====

```

Uses softmax, treats classes as UNORDERED.
Formula: $P(Y=k|x) = \frac{\exp(w_k^T x + b_k)}{\sum \exp(w_j^T x + b_j)}$

```

=====
📊 Logistic Regression
=====
Accuracy:      0.6595 (65.95%)
MAE:           0.5337

```

F1 (macro): 0.3793
F1 (weighted): 0.6064

Error Analysis:

Adjacent Error Rate (± 1): 65.17%
Severe Error Rate ($\pm 2+$): 34.83%

```
1 # Classification report
2 print("\n 📄 Classification Report - Logistic Regression:")
3 print(classification_report(y_test, lr_pred, digits=4))
```

📄 Classification Report - Logistic Regression:

	precision	recall	f1-score	support
1	0.6108	0.4374	0.5098	567
2	0.3059	0.0602	0.1006	432
3	0.3404	0.1223	0.1800	793
4	0.4197	0.2340	0.3005	2021
5	0.7103	0.9299	0.8054	6179
accuracy			0.6595	9992
macro avg	0.4774	0.3568	0.3793	9992
weighted avg	0.5991	0.6595	0.6064	9992

Step 7: Confusion Matrices

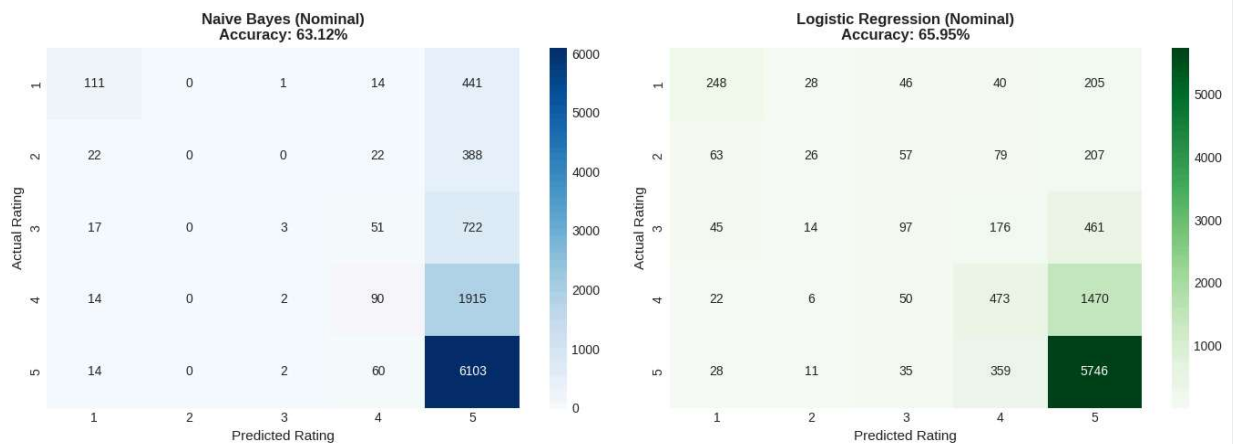
```
1 # =====
2 # CONFUSION MATRICES
3 # =====
4
5 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
6
7 # Naive Bayes
8 cm_nb = confusion_matrix(y_test, nb_pred)
9 sns.heatmap(cm_nb, annot=True, fmt='d', cmap='Blues', ax=axes[0],
10             xticklabels=[1, 2, 3, 4, 5], yticklabels=[1, 2, 3, 4, 5])
11 axes[0].set_xlabel('Predicted Rating', fontsize=11)
12 axes[0].set_ylabel('Actual Rating', fontsize=11)
13 axes[0].set_title(f'Naive Bayes (Nominal)\nAccuracy: {nb_results["accuracy"]}',
14                  fontsize=12, fontweight='bold')
15
16 # Logistic Regression
17 cm_lr = confusion_matrix(y_test, lr_pred)
18 sns.heatmap(cm_lr, annot=True, fmt='d', cmap='Greens', 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'Logistic Regression (Nominal)\nAccuracy: {lr_results["accuracy"]}',
23                  fontsize=12, fontweight='bold')
```



```

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

```




✅ Saved: confusion_matrices_nominal.png

Step 8: Save Results


```


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


```

```
=====
 NOMINAL MODELS SUMMARY
=====
```

	model	encoding	accuracy	mae	f1_macro	f1_weighted	adjacent
	Naive Bayes	Nominal	0.631205	0.665132	0.232055	0.512934	0.
	Logistic Regression	Nominal	0.659528	0.533727	0.379251	0.606418	0.

 Saved: nominal_results.csv

```
1 # Save predictions for later analysis
2 predictions_df = pd.DataFrame({
3     'actual': y_test,
4     'nb_pred': nb_pred,
5     'lr_pred': lr_pred
6 })
7 predictions_df.to_csv('nominal_predictions.csv', index=False)
8 print( Saved: nominal_predictions.csv)
```

 Saved: nominal_predictions.csv

```
1 # Download files
2 try:
3     from google.colab import files
4     files.download('nominal_results.csv')
5     files.download('confusion_matrices_nominal.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 # Create the directory if it doesn't exist
7 os.makedirs(drive_path, exist_ok=True)
8
9 files_to_save = [
10     'nominal_results.csv',
11     'confusion_matrices_nominal.png',
12     'nominal_predictions.csv'
13 ]
14
15 for file_name in files_to_save:
16     source_path = os.path.join('/content', file_name)
17     destination_path = os.path.join(drive_path, file_name)
18     try:
```

```
19         shutil.copy(source_path, destination_path)
20         print(f"✅ Saved '{file_name}' to Drive: {destination_path}")
21     except FileNotFoundError:
22         print(f"⚠️ Error: '{file_name}' not found. Skipping.")
23     except Exception as e:
24         print(f"❌ Error saving '{file_name}' to Drive: {e}")
```

✅ Saved 'nominal_results.csv' to Drive: /content/drive/MyDrive/Fall 2025/Founda
✅ Saved 'confusion_matrices_nominal.png' to Drive: /content/drive/MyDrive/Fall
✅ Saved 'nominal_predictions.csv' to Drive: /content/drive/MyDrive/Fall 2025/Fa

📌 Summary