# 📊 Notebook 2: EDA & Visualization Final Project - Ordinal vs Nominal Sentiment Analysis
## Atharv Chaudhary

---

**Purpose:** Exploratory Data Analysis and create visualizations for report.

**Input:** `amazon_electronics_cleaned.csv`

**Output:** `class_distribution.png`

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

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call dr
```

```
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    # Plot settings
10   plt.style.use('seaborn-v0_8-whitegrid')
11   sns.set_palette("husl")
12   plt.rcParams['figure.dpi'] = 150
13
14   print("✅ Libraries imported")
```

```
✅ Libraries imported
```

## Step 1: Load Cleaned Data

```
1 # Load cleaned data from Notebook 1
2 df = pd.read_csv(r'/content/drive/MyDrive/Fall 2025/Foundations of Artifici
3
4 print(f"✅ Loaded {len(df):,} reviews")
5 print(f"   Columns: {list(df.columns)}")
6 df.head()
```

✅ Loaded 49,960 reviews
   Columns: ['text', 'rating']

|   | text | rating |
|---|---|---|
| 0 | We got this GPS for my husband who is an (OTR)... | 5 |
| 1 | I'm a professional OTR truck driver, and I bou... | 1 |
| 2 | Well, what can I say. I've had this unit in m... | 3 |
| 3 | Not going to write a long review, even thought... | 2 |
| 4 | I've had mine for a year and here's what we go... | 1 |

Next steps:  ( Generate code with `df` )  ( New interactive sheet )

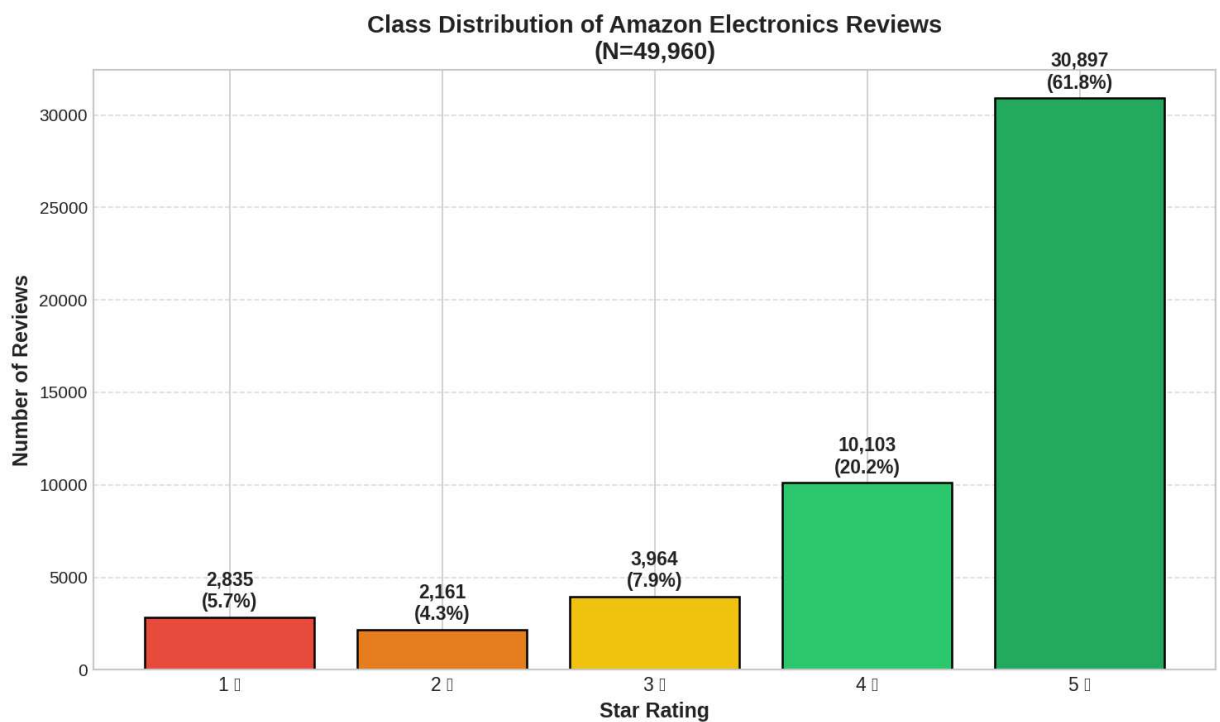## Step 2: Class Distribution Visualization

```
1  # ========================================================================
2  # CLASS DISTRIBUTION PLOT (For Report)
3  # ========================================================================
4
5  fig, ax = plt.subplots(figsize=(10, 6))
6
7  rating_counts = df['rating'].value_counts().sort_index()
8  colors = ['#e74c3c', '#e67e22', '#f1c40f', '#2ecc71', '#27ae60']
9
10 bars = ax.bar(rating_counts.index, rating_counts.values, color=colors,
11               edgecolor='black', linewidth=1.2)
12
13 # Add value labels on bars
14 for bar, count in zip(bars, rating_counts.values):
15     height = bar.get_height()
16     ax.text(bar.get_x() + bar.get_width()/2, height + max(rating_counts)*(
17             f'{count:,}\n({count/len(df)*100:.1f}%)',
18             ha='center', va='bottom', fontsize=11, fontweight='bold')
19
20 ax.set_xlabel('Star Rating', fontsize=12, fontweight='bold')
21 ax.set_ylabel('Number of Reviews', fontsize=12, fontweight='bold')
22 ax.set_title(f'Class Distribution of Amazon Electronics Reviews\n(N={len(
23             fontsize=14, fontweight='bold')
24 ax.set_xticks([1, 2, 3, 4, 5])
25 ax.set_xticklabels(['1 ⭐', '2 ⭐', '3 ⭐', '4 ⭐', '5 ⭐'], fontsize=1
26
27 # Add grid
28 ax.yaxis.grid(True, linestyle='--', alpha=0.7)
29 ax.set_axisbelow(True)
30
31 plt.tight_layout()
32 plt.savefig('/content/drive/MyDrive/Fall 2025/Foundations of Artificial I
```

```
33           facecolor='white', edgecolor='none')
34 plt.show()
35
36 print("\n✅ Saved: /content/drive/MyDrive/Fall 2025/Foundations of Artifi
```

**Class Distribution of Amazon Electronics Reviews**
**(N=49,960)**



✅ Saved: /content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligen

## Step 3: Review Length Analysis

```
1 # Calculate text statistics
2 df['text_length'] = df['text'].str.len()
3 df['word_count'] = df['text'].str.split().str.len()
4
5 print("📏 Review Length Statistics:")
6 print(df[['text_length', 'word_count']].describe())
```
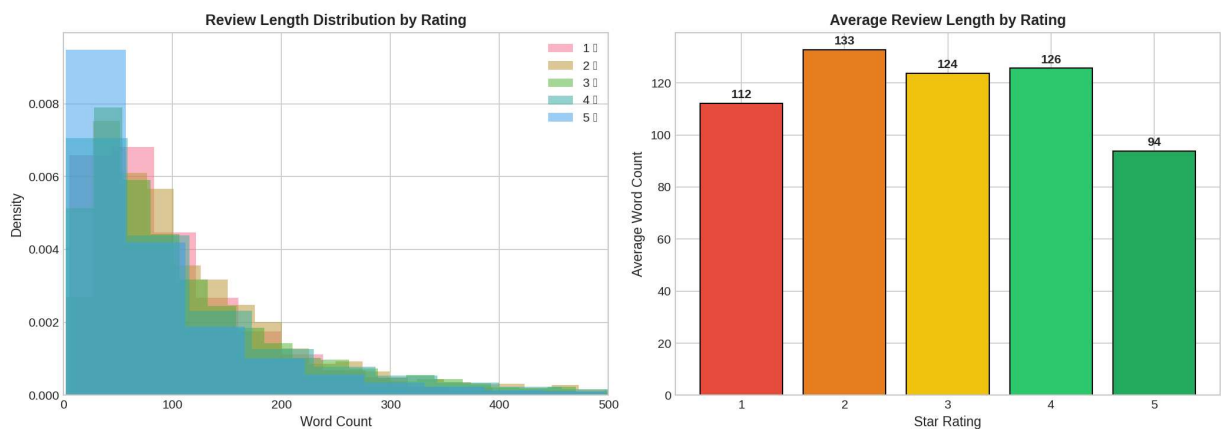
📏 Review Length Statistics:
```
       text_length    word_count
count  49960.000000  49960.000000
mean     577.486930    105.335028
```

```
std       698.808926      124.742919
min        10.000000        2.000000
25%       173.000000       32.000000
50%       337.000000       62.500000
75%       698.000000      128.000000
max     15567.000000     2845.000000
```

```python
 1 # ===========================================================================
 2 # REVIEW LENGTH BY RATING
 3 # ===========================================================================
 4
 5 # Ensure 'word_count' is available in case it was dropped prematurely
 6 if 'word_count' not in df.columns:
 7     df['text_length'] = df['text'].str.len()
 8     df['word_count'] = df['text'].str.split().str.len()
 9
10 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
11
12 # Word count distribution
13 for rating in [1, 2, 3, 4, 5]:
14     subset = df[df['rating'] == rating]['word_count']
15     axes[0].hist(subset, bins=50, alpha=0.5, label=f'{rating} ⭐', density
16
17 axes[0].set_xlabel('Word Count', fontsize=11)
18 axes[0].set_ylabel('Density', fontsize=11)
19 axes[0].set_title('Review Length Distribution by Rating', fontsize=12, font
20 axes[0].legend()
21 axes[0].set_xlim([0, 500])
22
23 # Average word count by rating
24 avg_words = df.groupby('rating')['word_count'].mean()
25 bars = axes[1].bar(avg_words.index, avg_words.values, color=colors, edgecol
26 axes[1].set_xlabel('Star Rating', fontsize=11)
27 axes[1].set_ylabel('Average Word Count', fontsize=11)
28 axes[1].set_title('Average Review Length by Rating', fontsize=12, fontweigh
29 axes[1].set_xticks([1, 2, 3, 4, 5])
30
31 for bar, val in zip(bars, avg_words.values):
32     axes[1].text(bar.get_x() + bar.get_width()/2, val + 2, f'{val:.0f}',
33                  ha='center', fontsize=10, fontweight='bold')
34
35 plt.tight_layout()
36 plt.savefig('/content/drive/MyDrive/Fall 2025/Foundations of Artificial Int
37 plt.show()
38
39 print("\n✅ Saved: /content/drive/MyDrive/Fall 2025/Foundations of Artific
```

✅ Saved: /content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligen

## Step 4: Sample Reviews

```
1 # Show sample reviews for each rating
2 print("=" * 70)
3 print("SAMPLE REVIEWS BY RATING")
4 print("=" * 70)
5
6 for rating in [1, 2, 3, 4, 5]:
7     sample = df[df['rating'] == rating]['text'].iloc[0][:200]
8     print(f"\n{'⭐' * rating} ({rating}-star):")
9     print(f"   {sample}...")
```

```
======================================================================
SAMPLE REVIEWS BY RATING
======================================================================

⭐ (1-star):
   I'm a professional OTR truck driver, and I bought a TND 700 at a truck stop h

⭐⭐ (2-star):
   Not going to write a long review, even thought this unit deserves one. I've d

⭐⭐⭐ (3-star):
   Well, what can I say.  I've had this unit in my truck for about four days now

⭐⭐⭐⭐ (4-star):
   This adapter easily connects my Nook HD 7&#34; to my HDTV through the HDMI ca

⭐⭐⭐⭐⭐ (5-star):
```

> We got this GPS for my husband who is an (OTR) over the road trucker.  Very I

## Step 5: Key Insights

```python
 1 # ============================================================
 2 # KEY INSIGHTS FOR REPORT
 3 # ============================================================
 4
 5 print("=" * 70)
 6 print("📋  KEY INSIGHTS FOR REPORT")
 7 print("=" * 70)
 8
 9 rating_counts = df['rating'].value_counts().sort_index()
10
11 print(f"""
12 1. DATASET SIZE:
13    - Total reviews: {len(df):,}
14    - Source: Amazon Electronics Reviews (McAuley Lab, UCSD)
15
16 2. CLASS IMBALANCE:
17    - 5-star reviews: {rating_counts[5]:,} ({rating_counts[5]/len(df)*100:.1
18    - 1-star reviews: {rating_counts[1]:,} ({rating_counts[1]/len(df)*100:.1
19    - Imbalance ratio: {rating_counts[5]/rating_counts.min():.1f}:1
20
21 3. REVIEW LENGTH:
22    - Average words: {df['word_count'].mean():.0f}
23    - Median words: {df['word_count'].median():.0f}
24    - Negative reviews tend to be longer (more detail about complaints)
25
26 4. CHALLENGE FOR CLASSIFICATION:
27    - Adjacent ratings (4★ vs 5★) use similar vocabulary
28    - Class imbalance affects minority class performance
29    - Ordinal structure: 1 < 2 < 3 < 4 < 5
30 """)
```

```
======================================================================
📋  KEY INSIGHTS FOR REPORT
======================================================================

1. DATASET SIZE:
   - Total reviews: 49,960
   - Source: Amazon Electronics Reviews (McAuley Lab, UCSD)

2. CLASS IMBALANCE:
```

```
      - 5-star reviews: 30,897 (61.8%)
      - 1-star reviews: 2,835 (5.7%)
      - Imbalance ratio: 14.3:1

  3. REVIEW LENGTH:
      - Average words: 105
      - Median words: 62
      - Negative reviews tend to be longer (more detail about complaints)

  4. CHALLENGE FOR CLASSIFICATION:
      - Adjacent ratings (4★ vs 5★) use similar vocabulary
      - Class imbalance affects minority class performance
      - Ordinal structure: 1 < 2 < 3 < 4 < 5
```

```python
1 # Clean up helper columns
2 df = df.drop(columns=['text_length', 'word_count'], errors='ignore')
3
4 # Images are now saved directly to Google Drive, no need for separate downl
5 print("Images have been saved directly to Google Drive.")
```

```
Images have been saved directly to Google Drive.
```

```python
1   # This cell is no longer needed, as images are now saved directly to
    Google Drive by the plt.savefig commands.
2   print("Images have been saved directly to Google Drive, no further move
    operation is required.")
```

```
Images have been saved directly to Google Drive, no further move operation is re
```

---

## ✅ Summary

**Visualizations created:**

- `class_distribution.png` – For Dataset section of report
- `review_length_analysis.png` – Additional analysis

**Key findings:**

- Severe class imbalance (5-star dominant)
- Negative reviews are longer on average
- Ordinal structure should be leveraged

**Next:** Run `3_Models_Nominal.ipynb`