

Notebook 1: Data Loading & Preprocessing

Purpose: Load Amazon Electronics Reviews, clean, and save for other notebooks.

Output: `amazon_electronics_cleaned.csv`

```
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 gzip
5 import json
6 import urllib.request
7 import warnings
8 warnings.filterwarnings('ignore')
9
10 print("✅ Libraries imported")
```

✅ Libraries imported

Step 1: Load Dataset

```
1 # =====
2 # CONFIGURATION
3 # =====
4
5 SAMPLE_SIZE = 50000 # Number of reviews to load (increase for final run)
6 RANDOM_STATE = 42
7
8 print(f"Configuration:")
9 print(f"  Sample size: {SAMPLE_SIZE},")
10 print(f"  Random state: {RANDOM_STATE}")
```

Configuration:
Sample size: 50,000
Random state: 42

```
1 # =====
2 # LOAD AMAZON ELECTRONICS REVIEWS (2014 Dataset - Reliable)
3 # =====
4
5 print("=" * 70)
6 print("LOADING AMAZON ELECTRONICS REVIEWS DATASET")
7 print("Source: UCSD McAuley Lab (Stanford SNAP)")
8 print("=" * 70)
9
10 url = "http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/reviews_Electronics_5.json.gz"
11
12 print(f"\nDownloading {SAMPLE_SIZE:,} reviews...")
13 print("⏳ This may take 1-2 minutes...\n")
14
15 try:
16     # Download file
17     urllib.request.urlretrieve(url, 'electronics.json.gz')
18
19     # Load JSONL.gz file
20     reviews = []
21     with gzip.open('electronics.json.gz', 'rt', encoding='utf-8') as f:
22         for i, line in enumerate(f):
23             if i >= SAMPLE_SIZE:
24                 break
25             if i % 10000 == 0 and i > 0:
26                 print(f"  Loaded {i:,} reviews...")
27             reviews.append(json.loads(line))
28
29     df_raw = pd.DataFrame(reviews)
30     print(f"\n✅ SUCCESS! Loaded {len(df_raw):,} reviews")
```

```

31     print(f"  Columns: {df_raw.columns.tolist()}")
32
33 except Exception as e:
34     print(f"❌ Error: {e}")
35     df_raw = None

```

=====
LOADING AMAZON ELECTRONICS REVIEWS DATASET
Source: UCSD McAuley Lab (Stanford SNAP)
=====

Downloading 50,000 reviews...
⌚ This may take 1-2 minutes...

Loaded 10,000 reviews...
Loaded 20,000 reviews...
Loaded 30,000 reviews...
Loaded 40,000 reviews...

SUCCESS! Loaded 50,000 reviews

Columns: ['reviewerID', 'asin', 'reviewerName', 'helpful', 'reviewText', 'overall', 'summary', 'unixReviewTime', 'reviewTime']

```

1 # Preview raw data
2 print("\n📋 Raw Data Preview:")
3 df_raw.head(3)

```

📋 Raw Data Preview:

	reviewerID	asin	reviewerName	helpful	reviewText	overall	summary	unixReviewTime	reviewTime	grid icon	more icon
0	AO94DHGC771SJ	0528881469	amazdnu	[0, 0]	We got this GPS for my husband who is an (OTR)...	5.0	Gotta have GPS!	1370131200	06 2, 2013		
1	AMO214LNFC EI4	0528881469	Amazon Customer	[12, 15]	I'm a professional OTR truck driver, and I bou...	1.0	Very Disappointed	1290643200	11 25, 2010		
					...						

Next steps: [Generate code with df_raw](#) [New interactive sheet](#)

Step 2: Data Cleaning

```

1 # =====
2 # DATA CLEANING
3 # =====
4
5 print("=" * 70)
6 print("DATA CLEANING")
7 print("=" * 70)
8
9 # Rename columns to standard names
10 df = df_raw.rename(columns={
11     'overall': 'rating',
12     'reviewText': 'text'
13 }).copy()
14
15 print(f"\nOriginal size: {len(df):,} reviews")
16
17 # Keep only relevant columns
18 df = df[['text', 'rating']].copy()
19
20 # Step 1: Remove missing values
21 missing_before = df.isna().sum()
22 print(f"\nMissing values: text={missing_before['text']}, rating={missing_before['rating']}")
23 df = df.dropna(subset=['text', 'rating'])
24 print(f"After removing nulls: {len(df):,} reviews")
25
26 # Step 2: Remove very short reviews
27 df = df[df['text'].str.len() >= 10]
28 print(f"After removing short reviews (<10 chars): {len(df):,} reviews")
29
30 # Step 3: Convert rating to integer
31 df['rating'] = df['rating'].astype(int)
32
33 # Step 4: Verify ratings are 1-5

```

```

34 df = df[df['rating'].between(1, 5)]
35 print(f"After rating validation: {len(df)} reviews")
36
37 # Reset index
38 df = df.reset_index(drop=True)
39

```

```
=====
DATA CLEANING
=====
```

Original size: 50,000 reviews

Missing values: text=0, rating=0
After removing nulls: 50,000 reviews
After removing short reviews (<10 chars): 49,960 reviews
After rating validation: 49,960 reviews

Final cleaned dataset: 49,960 reviews

```

1 # Preview cleaned data
2 print("\n📝 Cleaned Data Preview:")
3 df.head()

```

📝 Cleaned Data Preview:

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

Next steps: [Generate code with df](#) [New interactive sheet](#)

Step 3: Class Distribution

```

1 # =====
2 # CLASS DISTRIBUTION
3 # =====
4
5 print("=" * 70)
6 print("CLASS DISTRIBUTION")
7 print("=" * 70)
8
9 rating_counts = df['rating'].value_counts().sort_index()
10
11 print("\n📊 Rating Distribution:")
12 for rating, count in rating_counts.items():
13     pct = count / len(df) * 100
14     bar = '█' * int(pct / 2)
15     print(f"  {rating} ⭐: {count:,>6,} ({pct:,>5.1f}%) {bar}")
16
17 print(f"\n  Total: {len(df),} reviews")

```

```
=====
CLASS DISTRIBUTION
=====
```

📊 Rating Distribution:

1	⭐:	2,835	(5.7%)	
2	⭐:	2,161	(4.3%)	
3	⭐:	3,964	(7.9%)	
4	⭐:	10,103	(20.2%)	
5	⭐:	30,897	(61.8%)	

Total: 49,960 reviews

```

1 # Check class imbalance
2 print("\n⚠️ Class Imbalance Analysis:")
3 majority_class = rating_counts.max()
4 minority_class = rating_counts.min()

```

```

5 imbalance_ratio = majority_class / minority_class
6
7 print(f"  Majority class (5-star): {majority_class:,}")
8 print(f"  Minority class: {minority_class:,}")
9

```

⚠ Class Imbalance Analysis:
 Majority class (5-star): 30,897
 Minority class: 2,161
 Imbalance ratio: 14.3:1

Step 4: Text Statistics

```

1 # =====
2 # TEXT STATISTICS
3 # =====
4
5 print("=" * 70)
6 print("TEXT STATISTICS")
7 print("=" * 70)
8
9 # Calculate text length
10 df['text_length'] = df['text'].str.len()
11 df['word_count'] = df['text'].str.split().str.len()
12
13 print("\n📝 Review Length (characters):")
14 print(f"  Min: {df['text_length'].min():,.0f}")
15 print(f"  Max: {df['text_length'].max():,.0f}")
16 print(f"  Mean: {df['text_length'].mean():,.0f}")
17 print(f"  Median: {df['text_length'].median():,.0f}")
18
19 print("\n📌 Word Count:")
20 print(f"  Min: {df['word_count'].min():,.0f}")
21 print(f"  Max: {df['word_count'].max():,.0f}")
22 print(f"  Mean: {df['word_count'].mean():,.0f}")
23 print(f"  Median: {df['word_count'].median():,.0f}")
24
25 # Drop helper columns before saving
26 df = df.drop(columns=['text_length', 'word_count'])

```

```
=====
TEXT STATISTICS
=====

📝 Review Length (characters):
Min: 10
Max: 15,567
Mean: 577
Median: 337

📌 Word Count:
Min: 2
Max: 2,845
Mean: 105
Median: 62
```

Step 5: Save Cleaned Data

```

1 # =====
2 # SAVE CLEANED DATA
3 # =====
4
5 print("=" * 70)
6 print("SAVING CLEANED DATA")
7 print("=" * 70)
8
9 import os
10
11 # Define output path
12 drive_path = '/content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligence/Final Project/data'
13 output_file = os.path.join(drive_path, 'amazon_electronics_cleaned.csv')
14
15 # Create directory if it doesn't exist
16 os.makedirs(drive_path, exist_ok=True)
17
18 # Save to CSV

```

```

19 df.to_csv(output_file, index=False)
20
21 print(f"\n\n✅ Successfully saved cleaned data to: {output_file}")
22

=====
SAVING CLEANED DATA
=====

✅ Successfully saved cleaned data to: /content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligence/Final Project/data/amazon_electronics_cleaned.csv

```

```

1 # Download for Google Colab
2 try:
3     from google.colab import files
4     files.download(output_file)
5     print("⬇️ Download started...")
6 except:
7     print("Not in Colab - file saved locally")

```

⬇️ Download started...

✓ Summary

Data loaded and cleaned!

Metric	Value
Total reviews	See output above
Columns	text, rating
Rating range	1–5
Output file	amazon_electronics_cleaned.csv

Next: Run 2_EDA_Visualization.ipynb

```

1 # Final summary
2 print("\n" + "=" * 70)
3 print("📝 NOTEBOOK 1 COMPLETE")
4 print("=" * 70)
5 print(f"\nDataset: Amazon Electronics Reviews")
6 print(f"Source: UCSD McAuley Lab")
7 print(f"Reviews: {len(df)}")
8 print(f"Output: {output_file}")
9 print("\n→ Next: Run Notebook 2 (EDA & Visualization)")

```

📝 NOTEBOOK 1 COMPLETE

Dataset: Amazon Electronics Reviews

Source: UCSD McAuley Lab

Reviews: 49,960

Output: /content/drive/MyDrive/Fall 2025/Foundations of Artificial Intelligence/Final Project/data/amazon_electronics_cleaned.csv

→ Next: Run Notebook 2 (EDA & Visualization)

