```
1   print("Installing required packages...")
2   import sys
3   import subprocess
4
5   packages = ['datasets', 'scikit-learn', 'pandas', 'numpy']
6   for package in packages:
7       try:
8           __import__(package.replace('-', '_'))
9       except ImportError:
10          print(f"Installing {package}...")
11          subprocess.check_call([sys.executable, "-m", "pip", "install", package, "-q"])
12
13  print("✓ All packages installed!")
14
```

```
Installing required packages...
Installing scikit-learn...
✓ All packages installed!
```

```
1   import pandas as pd
2   import numpy as np
3   from sklearn.model_selection import train_test_split
4   from sklearn.feature_extraction.text import TfidfVectorizer
5   from sklearn.naive_bayes import MultinomialNB
6   from sklearn.linear_model import LogisticRegression
7   from sklearn.metrics import accuracy_score, confusion_matrix, classification_report, mean_absolute_error
8   import re
9   import warnings
10  warnings.filterwarnings('ignore')
11
12  print("✓ Libraries imported successfully")
```

```
✓ Libraries imported successfully
```

```
1 print("=" * 70)
2 print("LOADING AMAZON ELECTRONICS REVIEWS DATASET")
3 print("=" * 70)
```

```
======================================================================
LOADING AMAZON ELECTRONICS REVIEWS DATASET
======================================================================
```

```
1 df = None
2 try:
3     from datasets import load_dataset
4
5     print("\nLoading Electronics reviews from HuggingFace...")
6     print("This may take 2-3 minutes on first download...\n")
7
8     dataset = load_dataset(
9         "McAuley-Lab/Amazon-Reviews-2023",
10        "raw_review_Electronics",
11        split="full[:10000]",
12        trust_remote_code=False  # Updated parameter
13    )
14
15    df = pd.DataFrame(dataset)
16    print(f"✅ Loaded {len(df)} reviews via HuggingFace")
17
18 except Exception as e:
19    print(f"HuggingFace method failed: {e}")
20    print("\nTrying alternative: Old Amazon Dataset (2014)...\n")
21
22    # Method 2: Old UCSD Dataset (backup)
23    try:
24        import gzip
25        import json
26        import urllib.request
27
28        url = "http://snap.stanford.edu/data/amazon/productGraph/categoryFiles/reviews_Electronics_5.json.gz"
29
30        print("Downloading from UCSD...")
31        urllib.request.urlretrieve(url, 'electronics.json.gz')
32
```

```
33        reviews = []
34        with gzip.open('electronics.json.gz', 'rt', encoding='utf-8') as f:
35            for i, line in enumerate(f):
36                if i >= 10000:
37                    break
38                reviews.append(json.loads(line))
39
40        df = pd.DataFrame(reviews)
41        print(f"✅ Loaded {len(df)} reviews from 2014 dataset")
42
43    except Exception as e2:
44        print(f"❌ All methods failed: {e2}")
45        raise
46
47 # Display dataset info
48 print(f"\nDataset Shape: {df.shape}")
49 print(f"Columns: {df.columns.tolist()}")
50 print("\nFirst 3 reviews:")
51 print(df.head(3))
```

```
Loading Electronics reviews from HuggingFace...
This may take 2-3 minutes on first download...

README.md:       30.3k/? [00:00<00:00, 1.33MB/s]

Amazon-Reviews-2023.py:       39.6k/? [00:00<00:00, 1.57MB/s]

HuggingFace method failed: Dataset scripts are no longer supported, but found Amazon-Reviews-2023.py

Trying alternative: Old Amazon Dataset (2014)...

Downloading from UCSD...
✅ Loaded 10000 reviews from 2014 dataset

Dataset Shape: (10000, 9)
Columns: ['reviewerID', 'asin', 'reviewerName', 'helpful', 'reviewText', 'overall', 'summary', 'unixReviewTime', 'reviewTime']

First 3 reviews:
      reviewerID        asin      reviewerName   helpful  \
0   AO94DHGC771SJ  0528881469           amazdnu    [0, 0]
1  AMO214LNFCEI4  0528881469   Amazon Customer  [12, 15]
2  A3N7T0DY83Y4IG  0528881469     C. A. Freeman  [43, 45]

                                         reviewText  overall  \
0  We got this GPS for my husband who is an (OTR)...      5.0
1  I'm a professional OTR truck driver, and I bou...      1.0
2  Well, what can I say.  I've had this unit in m...      3.0

              summary  unixReviewTime   reviewTime
0      Gotta have GPS!      1370131200   06 2, 2013
1  Very Disappointed      1290643200  11 25, 2010
2      1st impression      1283990400   09 9, 2010
```

```
 1 print("\n" + "=" * 70)
 2 print("DATA PREPROCESSING")
 3 print("=" * 70)
 4
 5 # Identify column names
 6 text_col = None
 7 rating_col = None
 8
 9 for col in ['text', 'reviewText', 'review_text']:
10     if col in df.columns:
11         text_col = col
12         break
13
14 for col in ['rating', 'overall', 'stars']:
15     if col in df.columns:
16         rating_col = col
17         break
18
19 print(f"\n√ Text column: '{text_col}'")
20 print(f"√ Rating column: '{rating_col}'")
21
22 # Clean data
23 df_clean = df[[text_col, rating_col]].dropna()
24 print(f"√ Removed {len(df) - len(df_clean)} rows with missing values")
25
26 # Convert ratings to integers
```

```
27 df_clean[rating_col] = pd.to_numeric(df_clean[rating_col], errors='coerce')
28 df_clean = df_clean.dropna()
29 df_clean[rating_col] = df_clean[rating_col].astype(int)
30 df_clean = df_clean[df_clean[rating_col].isin([1, 2, 3, 4, 5])]
31
32 print(f"✓ Final dataset: {len(df_clean)} reviews")
33
34 # Show class distribution
35 print("\n📊 Class Distribution:")
36 dist = df_clean[rating_col].value_counts().sort_index()
37 for rating, count in dist.items():
38     pct = (count / len(df_clean)) * 100
39     print(f"  {rating}-star: {count:5d} ({pct:5.1f}%)")
```

```
======================================================================
DATA PREPROCESSING
======================================================================

✓ Text column: 'reviewText'
✓ Rating column: 'overall'
✓ Removed 0 rows with missing values
✓ Final dataset: 10000 reviews

📊 Class Distribution:
  1-star:   572 (  5.7%)
  2-star:   450 (  4.5%)
  3-star:   822 (  8.2%)
  4-star:  2095 ( 20.9%)
  5-star:  6061 ( 60.6%)
```

```
 1 # %% Cell 5: Text Cleaning
 2 print("\n" + "=" * 70)
 3 print("TEXT PREPROCESSING")
 4 print("=" * 70)
 5
 6 def clean_text(text):
 7     """Clean and normalize text"""
 8     if not isinstance(text, str):
 9         return ""
10     text = text.lower()
11     text = re.sub(r'[^a-z\s]', '', text)
12     text = ' '.join(text.split())
13     return text
14
15 print("\nCleaning text...")
16 df_clean['cleaned_text'] = df_clean[text_col].apply(clean_text)
17
18 # Remove very short reviews
19 df_clean = df_clean[df_clean['cleaned_text'].str.len() >= 10]
20 print(f"✓ Text cleaned. Final reviews: {len(df_clean)}")
21
22 print("\nExample cleaned review:")
23 print(f"Original: {df_clean[text_col].iloc[0][:100]}...")
24 print(f"Cleaned:  {df_clean['cleaned_text'].iloc[0][:100]}...")
25
```

```
======================================================================
TEXT PREPROCESSING
======================================================================

Cleaning text...
✓ Text cleaned. Final reviews: 9991

Example cleaned review:
Original: We got this GPS for my husband who is an (OTR) over the road trucker.  Very Impressed with the shipp...
Cleaned:  we got this gps for my husband who is an otr over the road trucker very impressed with the shipping ...
```

```
 1 print("\n" + "=" * 70)
 2 print("FEATURE EXTRACTION")
 3 print("=" * 70)
 4
 5 X = df_clean['cleaned_text']
 6 y = df_clean[rating_col]
 7
 8 # Train-test split
 9 print("\nSplitting data (80% train, 20% test)...")
10 X_train, X_test, y_train, y_test = train_test_split(
```

```
11      X, y, test_size=0.2, random_state=42, stratify=y
12 )
13
14 # TF-IDF Vectorization
15 print("Extracting TF-IDF features...")
16 vectorizer = TfidfVectorizer(max_features=1000, min_df=5, max_df=0.8)
17 X_train_tfidf = vectorizer.fit_transform(X_train)
18 X_test_tfidf = vectorizer.transform(X_test)
19
20 print(f"✓ Training set: {X_train_tfidf.shape[0]} samples")
21 print(f"✓ Test set: {X_test_tfidf.shape[0]} samples")
22 print(f"✓ Features: {X_train_tfidf.shape[1]} TF-IDF terms")
23
24 print("\nTop 10 features by TF-IDF score:")
25 feature names = vectorizer.get feature names out
```

```
======================================================================
FEATURE EXTRACTION
======================================================================

Splitting data (80% train, 20% test)...
Extracting TF-IDF features...
✓ Training set: 7992 samples
✓ Test set: 1999 samples
✓ Features: 1000 TF-IDF terms

Top 10 features by TF-IDF score:
```

```
 1 print("\n" + "=" * 70)
 2 print("TRAINING MODELS")
 3 print("=" * 70)
 4
 5 # Model 1: Multinomial Naive Bayes
 6 print("\n[1/2] Training Naive Bayes...")
 7 nb_model = MultinomialNB()
 8 nb_model.fit(X_train_tfidf, y_train)
 9 nb_pred = nb_model.predict(X_test_tfidf)
10 print("✓ Naive Bayes trained")
11
12 # Model 2: Logistic Regression
13 print("\n[2/2] Training Logistic Regression...")
14 lr_model = LogisticRegression(max_iter=500, random_state=42, multi_class='multinomial')
15 lr_model.fit(X_train_tfidf, y_train)
16 lr_pred = lr_model.predict(X_test_tfidf)
17 print("✓ Logistic Regression trained")
```

```
======================================================================
TRAINING MODELS
======================================================================

[1/2] Training Naive Bayes...
✓ Naive Bayes trained

[2/2] Training Logistic Regression...
✓ Logistic Regression trained
```

```
 1 # %% Cell 8: Evaluate Results
 2 print("\n" + "=" * 70)
 3 print("EVALUATION RESULTS")
 4 print("=" * 70)
 5
 6 # Calculate metrics
 7 nb_acc = accuracy_score(y_test, nb_pred)
 8 nb_mae = mean_absolute_error(y_test, nb_pred)
 9 nb_cm = confusion_matrix(y_test, nb_pred)
10
11 lr_acc = accuracy_score(y_test, lr_pred)
12 lr_mae = mean_absolute_error(y_test, lr_pred)
13 lr_cm = confusion_matrix(y_test, lr_pred)
14
15 print(f"\n📊 MULTINOMIAL NAIVE BAYES:")
16 print(f"   Accuracy: {nb_acc*100:.1f}% ({nb_acc:.4f})")
17 print(f"   MAE: {nb_mae:.2f}")
18
19 print(f"\n📊 LOGISTIC REGRESSION (Nominal):")
20 print(f"   Accuracy: {lr_acc*100:.1f}% ({lr_acc:.4f})")
```

```
21 print(f"   MAE: {lr_mae:.2f}")
22
23 # Confusion matrices
24 print("\n📈 Confusion Matrix - Naive Bayes:")
25 print(nb_cm)
26
27 print("\n📈 Confusion Matrix - Logistic Regression:")
28 print(lr_cm)
```

```
============================================================
EVALUATION RESULTS
============================================================

📊 MULTINOMIAL NAIVE BAYES:
   Accuracy: 60.5% (0.6053)
   MAE: 0.74

📊 LOGISTIC REGRESSION (Nominal):
   Accuracy: 62.6% (0.6258)
   MAE: 0.62

📈 Confusion Matrix - Naive Bayes:
[[   0    0    0    0  114]
 [   0    0    0    1   89]
 [   0    0    0    1  164]
 [   0    0    0    0  419]
 [   0    0    0    1 1210]]

📈 Confusion Matrix - Logistic Regression:
[[  23    2    5    6   78]
 [   9    3    6   20   52]
 [   6    1    6   47  105]
 [   2    1    3   92  321]
 [   2    0    4   78 1127]]
```

```
 1 print("\n" + "=" * 70)
 2 print("DETAILED ANALYSIS")
 3 print("=" * 70)
 4
 5 # Adjacent rating error analysis
 6 def calc_adjacent_errors(cm):
 7     """Calculate % of errors that are adjacent ratings"""
 8     total_errors = np.sum(cm) - np.trace(cm)
 9     if total_errors == 0:
10         return 0
11     adjacent_errors = 0
12     for i in range(len(cm)):
13         for j in range(len(cm)):
14             if abs(i - j) == 1:
15                 adjacent_errors += cm[i][j]
16     return (adjacent_errors / total_errors) * 100
17
18 nb_adj = calc_adjacent_errors(nb_cm)
19 lr_adj = calc_adjacent_errors(lr_cm)
20
21 print(f"\n📊 Adjacent Rating Errors:")
22 print(f"   Naive Bayes: {nb_adj:.1f}% of errors are adjacent (e.g., 4↔5)")
23 print(f"   Logistic Regression: {lr_adj:.1f}% of errors are adjacent")
24
25 # Per-class performance
26 print("\n📊 Per-Class F1-Scores:")
27 nb_report = classification_report(y_test, nb_pred, output_dict=True, zero_division=0)
28 lr_report = classification_report(y_test, lr_pred, output_dict=True, zero_division=0)
29
30 print("\nNaive Bayes:")
31 for rating in [1, 2, 3, 4, 5]:
32     if str(rating) in nb_report:
33         f1 = nb_report[str(rating)]['f1-score']
34         support = nb_report[str(rating)]['support']
35         print(f"   {rating}-star: F1 = {f1:.3f} (n={int(support)})")
36
37 print("\nLogistic Regression:")
38 for rating in [1, 2, 3, 4, 5]:
39     if str(rating) in lr_report:
40         f1 = lr_report[str(rating)]['f1-score']
41         support = lr_report[str(rating)]['support']
```

```
================================================================
DETAILED ANALYSIS
================================================================

📊 Adjacent Rating Errors:
   Naive Bayes: 53.4% of errors are adjacent (e.g., 4↔5)
   Logistic Regression: 62.4% of errors are adjacent

📊 Per-Class F1-Scores:

Naive Bayes:
   1-star: F1 = 0.000 (n=114)
   2-star: F1 = 0.000 (n=90)
   3-star: F1 = 0.000 (n=165)
   4-star: F1 = 0.000 (n=419)
   5-star: F1 = 0.755 (n=1211)

Logistic Regression:
   1-star: F1 = 0.295 (n=114)
   2-star: F1 = 0.062 (n=90)
   3-star: F1 = 0.063 (n=165)
   4-star: F1 = 0.278 (n=419)
   5-star: F1 = 0.779 (n=1211)
```

```python
 1   # Get key metrics
 2   f1_5star = lr_report.get('5', {}).get('f1-score', 0)
 3   f1_3star = lr_report.get('3', {}).get('f1-score', 0)
 4
 5   paragraph = f"""Initial experiments on {len(df_clean):,} Electronics reviews show promising
 6   directions. Using TF-IDF features (max 1,000 features), Multinomial Naive Bayes
 7   achieved {nb_acc*100:.1f}% accuracy with MAE of {nb_mae:.2f}, while Logistic Regression
 8   (nominal treatment) achieved {lr_acc*100:.1f}% accuracy with MAE of {lr_mae:.2f}.
 9   Confusion matrix analysis reveals that {lr_adj:.0f}% of misclassifications occur
10   between adjacent ratings (particularly 4↔5 stars), supporting our hypothesis that
11   ordinal treatment may improve performance. The class imbalance is evident—models
12   achieve {f1_5star*100:.0f}% F1-score for 5-star reviews but only {f1_3star*100:.0f}%
13   for 3-star reviews. These preliminary findings motivate our investigation into
14   whether ordinal methods can reduce MAE by better modeling rating structure while
15   addressing the adjacent-rating confusion problem."""
16
17   print("\n" + paragraph)
```

```
Initial experiments on 9,991 Electronics reviews show promising
directions. Using TF-IDF features (max 1,000 features), Multinomial Naive Bayes
achieved 60.5% accuracy with MAE of 0.74, while Logistic Regression
(nominal treatment) achieved 62.6% accuracy with MAE of 0.62.
Confusion matrix analysis reveals that 62% of misclassifications occur
between adjacent ratings (particularly 4↔5 stars), supporting our hypothesis that
ordinal treatment may improve performance. The class imbalance is evident—models
achieve 78% F1-score for 5-star reviews but only 6%
for 3-star reviews. These preliminary findings motivate our investigation into
whether ordinal methods can reduce MAE by better modeling rating structure while
addressing the adjacent-rating confusion problem.
```