

# Amazon India Headphones Market Data Analysis

Ishita Dingare

June 2025

## Executive Summary

This project delivers a comprehensive analysis of the Amazon India headphones market using real-world data. Through a combination of SQL, Python and Power BI, I cleaned and structured the dataset, performed advanced statistical and machine learning analysis, and built interactive dashboards for business users. My findings reveal a highly competitive market with clear segmentation, strong brand dynamics, and actionable opportunities for both volume and premium strategies. The Power BI dashboard enables business stakeholders to explore these insights, supporting data-driven decisions in product positioning, pricing and brand strategy.

## 1. Business Problem and Objective

### 1.1 Objective:

To analyze the Amazon headphones market, uncover brand and product performance trends, understand pricing and value dynamics, and provide actionable recommendations for business growth.

### 1.2 Key Business Questions:

- Which brands and products dominate the market?
- How do pricing and customer ratings vary across segments?
- Where are the opportunities for value, premium, or volume strategies?
- How can data-driven insights inform product, pricing, and marketing decisions?

## 2. Data Collection and Preparation

### 2.1 Source:

Data was web-scraped from Amazon India, capturing product details, pricing, ratings, reviews, and brand information for headphones.

### 2.2 Cleaning:

- SQL and Python were used to standardize brand names, remove duplicates, fix web scraping errors, and handle the missing values.
- Final dataset: 723 products, 72 unique brands, with key columns: Title, Brand, Price, MRP, Discount, Rating, Review Count, Prime, ASIN, URL.

### 2.3 Feature Engineering:

- Added columns for Saving, Price Segment (Budget, Mid-range, Premium, Luxury), Value Score (rating per ₹1000), and product clusters.

## 3. Analytical Approach

### 3.1 SQL:

Initial data cleaning, deduplication, and summary statistics.

### 3.2 Python:

Advanced EDA, statistical correlation, clustering (KMeans), regression (Random Forest), and value analysis.

### 3.3 Power BI:

Interactive dashboards for executive overview, brand/product performance, price segmentation, and advanced insights.

## 4. Key Findings & Insights

### 4.1 Market Overview:

- **Market Size:** 723 unique products, 72 brands.
- **Price Range:** ₹99 - ₹29,990 (Avg: ₹1,892).
- **Average Rating:** 3.81/5.
- **Top Brands by Product Count:** boAt (159), Noise (60), Boult (60), ZEBRONICS (47), JBL (38).
- **Market Share:** boAt leads with 22% of products; top 10 brands control 70% of the market.

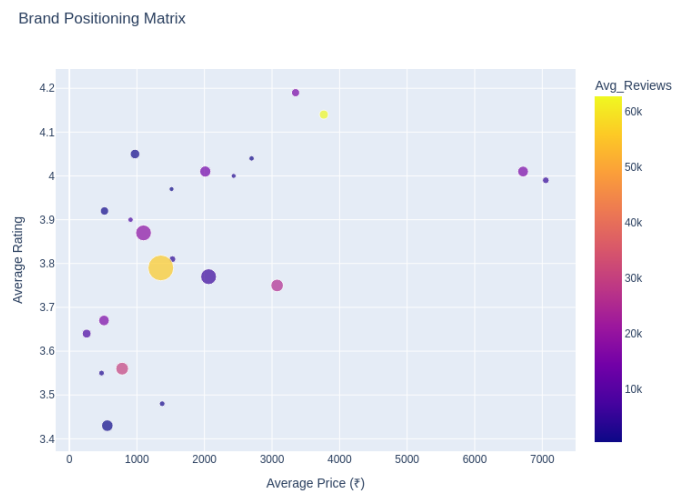


Figure 1: Overview

### 4.2 Brand & Product Performance:

- **Volume Leaders:** boAt dominates by product count and reviews engagement.
- **Premium Brands:** OnePlus (₹3,767, 4.14 rating), Sony (₹6,716, 4.01), soundcore (₹3,349, 4.19).
- **Value Brands:** Amazon Basics (₹562, 3.43), pTron (₹514, 3.67), Ambrane (₹256, 3.64).

- **Customer Engagement:** boAt and Noise have the highest review volumes, indicating strong brand presence.

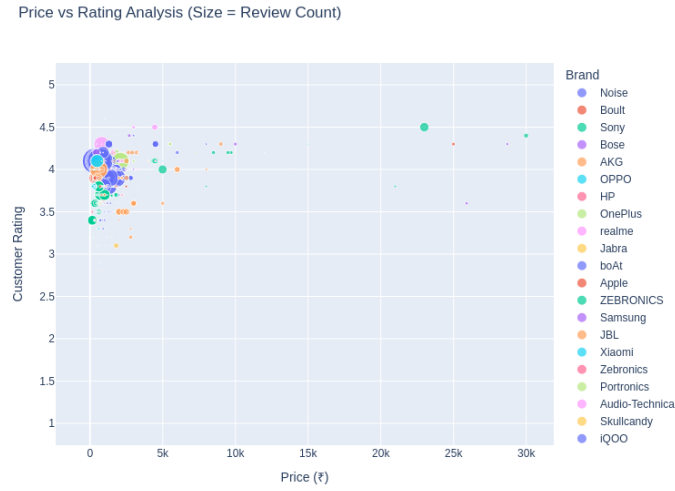


Figure 2: Price vs. Rating

### 4.3 Price Segmentation & Value:

- **Product Distribution:**
  - Mid-range (₹500-₹1,500): 345 products (48% of the market).
  - Premium (₹1,500-₹5,000): 197 products (27%).
  - Budget (<₹500): 138 products (19%).
  - Luxury (>₹5,000): 43 products (6%).
- **Quality-Price Correlation:**
  - Average rating increases with price: Budget (3.69), Mid-range (3.75), Premium (3.94), Luxury (4.10).
  - Value Score drops as price increases: Budget (11.36), Luxury (0.47).
- **Customer Engagement:**
  - Budget and mid-range products have the highest average review counts, luxury products have lower engagement but higher satisfaction.

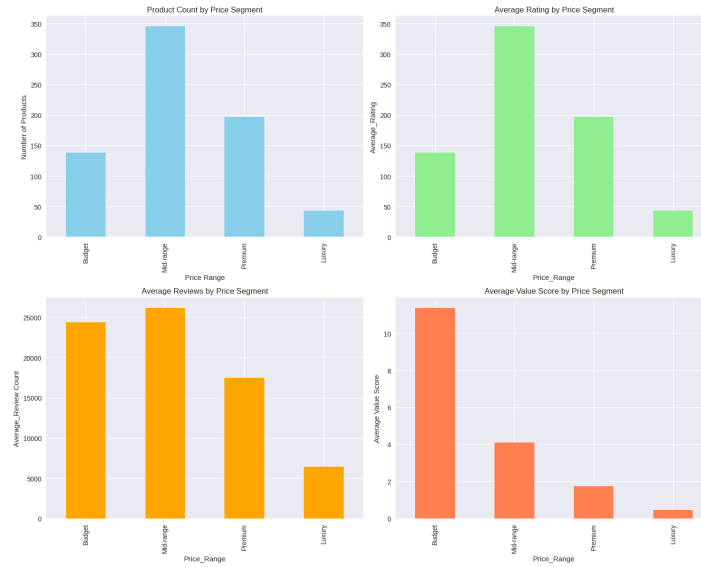


Figure 3: Market Segmentation Analysis

#### 4.4 Advanced Insights:

- **Correlation Analysis:**

- Moderate positive correlation between price and rating ( $r=0.22$ ,  $p<0.001$ ).
- Weak negative correlation between price and review count ( $r = -0.07$ ).
- Value Score is strongly negatively correlated with price ( $r = -0.39$ ): Cheaper products offer better rating-per-rupee.

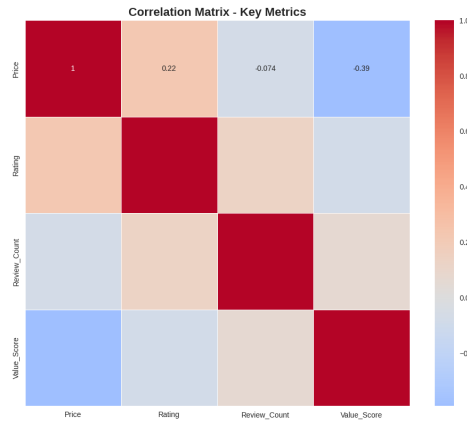


Figure 4: Correlation Analysis

- **Clustering:** Four clusters identified:

- Cluster 0: Mid-priced, high rating, very high reviews (popular, good value).
- Cluster 1: Upper mid-range, moderate rating/reviews.
- Cluster 2: Premium/luxury, highest price/rating, low value score.
- Cluster 3: Budget, moderate rating, high reviews, highest value score.

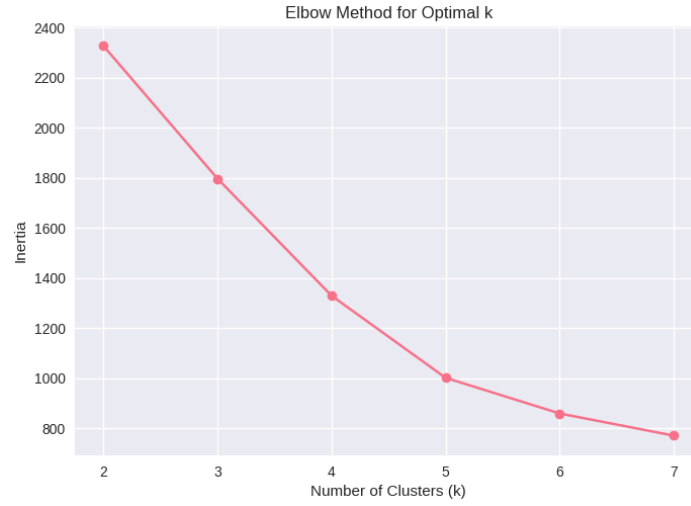


Figure 5: Elbow Method Plot

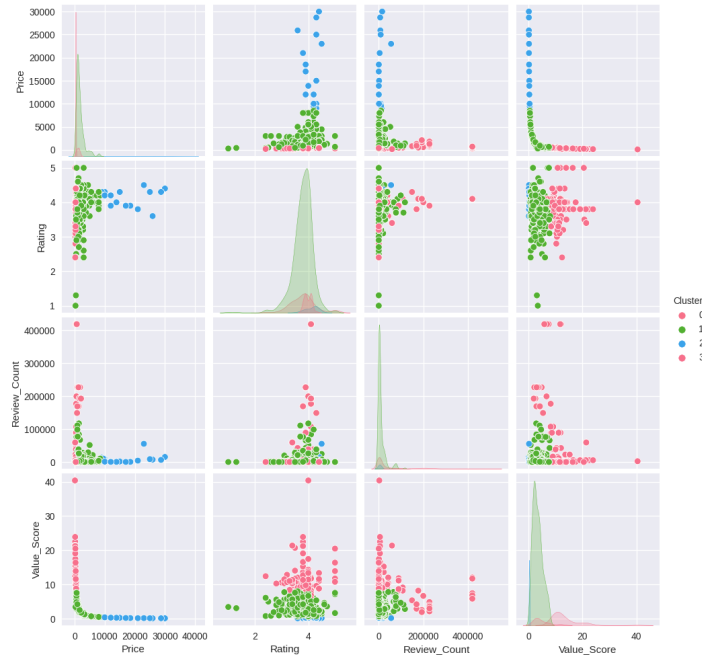


Figure 6: Cluster Analysis

- **Brand vs. Price Segment:**
  - Premium brands (OnePlus, Sony, soundcore) lead in ratings and price.
  - Value brands (truke, Portonics) deliver high ratings at low prices.
- **Prime vs. Non-Prime:**
  - Prime products (89% of market) command higher prices and ratings.

#### 4.5 Price Prediction Modeling:

##### *Random Forest Regression:*

- Test  $R^2$  up to 0.51 with brand, rating, and review count as features.

- Rating and brand are the most important predictors of price, but much price variation remains unexplained due to product features not captured in the dataset.

## 5. Recommendations

### 5.1 For New Entrants:

- **Target the mid-range segment (₹500 - ₹1,500):** Highest product count and customer engagement.
- **Focus on achieving 4.0+ ratings** to compete in premium segments.
- **Consider budget entry for volume strategies**—high review counts indicate strong demand.

### 5.2 For Existing Brands:

- **Premium migration:** Brands in the mid-range should improve quality to move into premium segments.
- **Exploit quality gaps:** Bridge the rating gap between segments and differentiation.
- **Leverage customer reviews:** Use feedback to improve products and boost ratings.

### 5.3 For Pricing Strategy:

- **₹1,000 - ₹1,500:** Highly competitive, requires strong differentiation.
- **₹2,000 - ₹3,000:** Premium positioning quality focus.
- **Under ₹500:** Volume play; maintain acceptable quality standards.

### 5.4 For Value Optimization:

- **Ultra-budget products (<₹200:)** Can achieve exceptional value scores and mass-market appeal.
- **Premium products:** Must justify higher prices with superior ratings and features.

## 6. Dashboard Support

### 6.1 Power BI dashboards allow business users to:

- Explore market share and brand performance interactively.
- Analyze product distribution, ratings, reviews, and value by price segment.
- Identify top-performing products and clusters.
- Drill down by brand, price segment, or Prime status for targeted insights.

### 6.2 Visuals include:

KPI cards, bar/column charts, scatter plots, heatmaps, and cluster analysis, all designed for executive decision-making.

## 7. Limitations & Next Steps

### 7.1 Data limitations:

- Web scraping errors and missing values were addressed, but some product features (e.g., technical specs, launch date) are not captured.
- Price prediction is limited by available features; including more product attributes could improve accuracy.

### 7.2 Future work:

- Add time-series analysis for price and rating trends.
- Integrate sales volume or revenue data for deeper business insights.
- Apply NLP on product reviews for sentiment analysis.

## 8. Appendix

### 8.1 Dashboard screenshots:

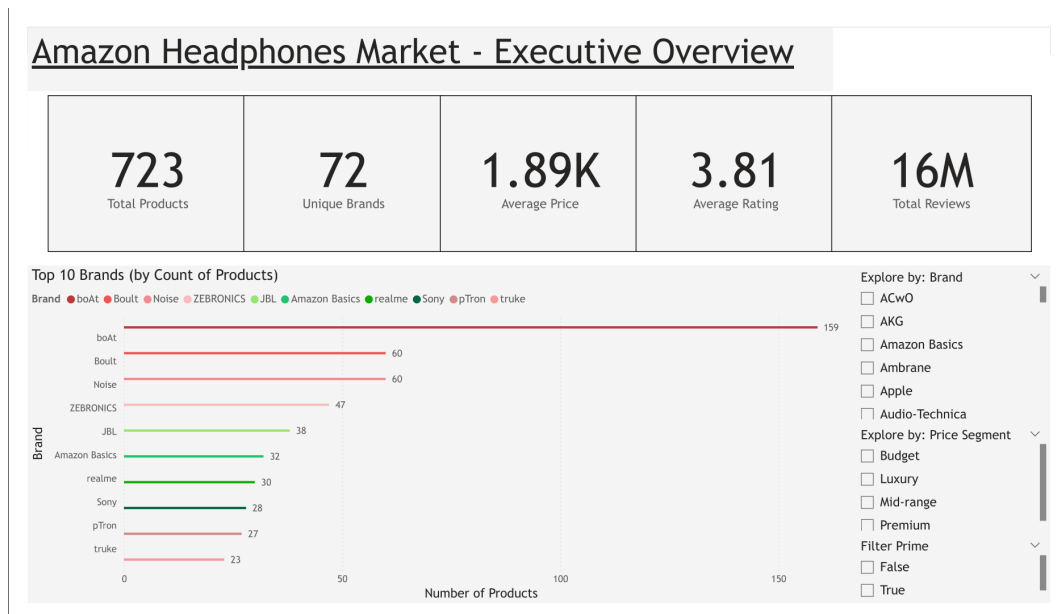


Figure 7: Executive Overview

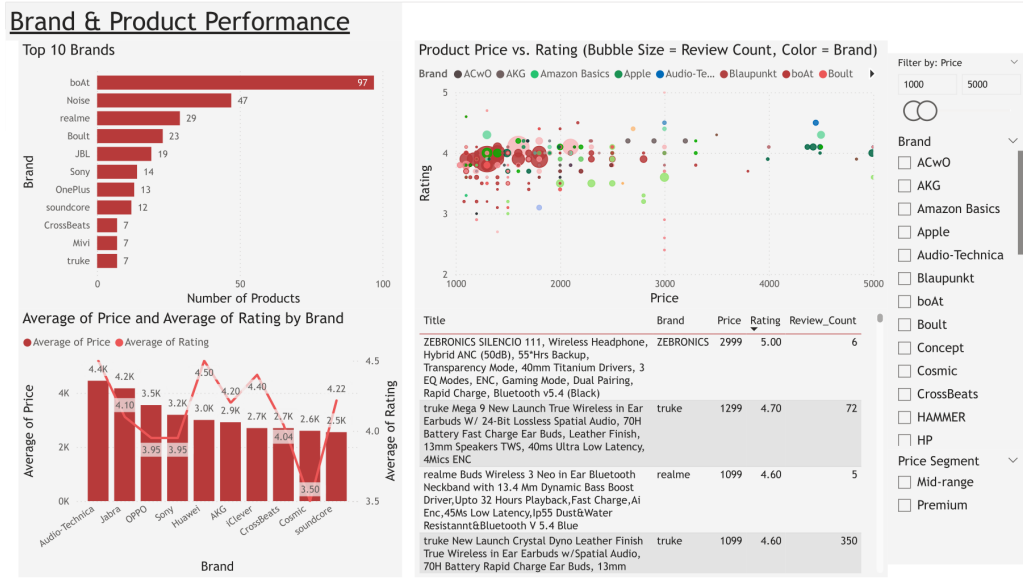


Figure 8: Brand & Product Performance

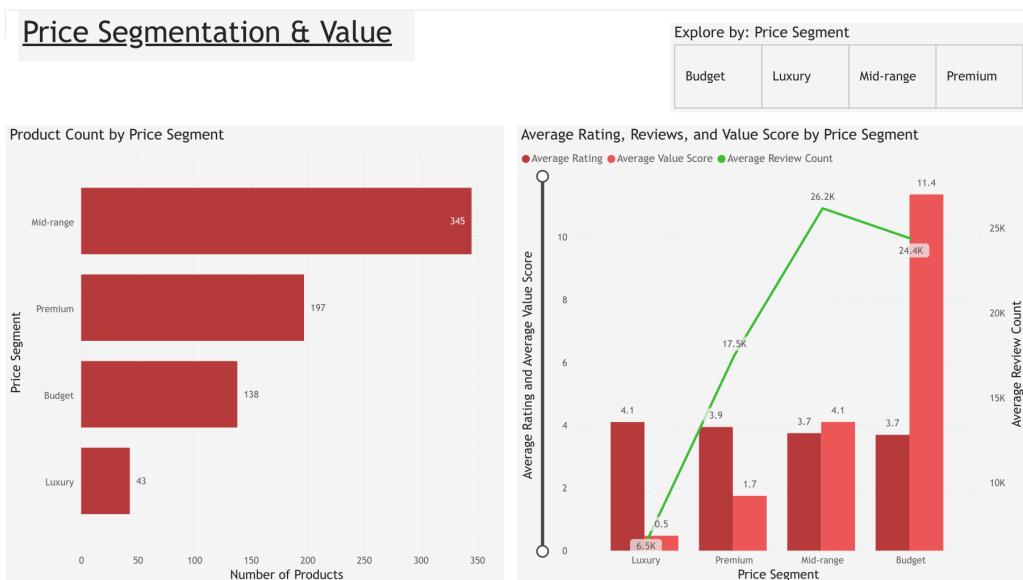


Figure 9: Price Segmentation & Value

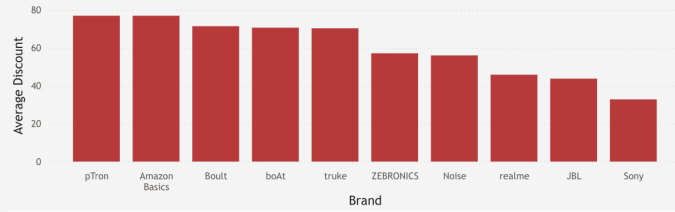


## Advanced Insights

Brand vs Price Segment Rating

Brand	Budget	Luxury	Mid-range	Premium	Total
Total	3.68	4.10	3.75	3.94	3.81
soundcore		4.27	4.08	4.23	4.19
OnePlus	3.80	4.25		4.12	4.14
CrossBeats			4.10	4.03	4.04
truke			4.06	3.80	4.04
Apple		3.96		4.10	4.02
Sony		4.19	3.83	4.00	4.01
realme			3.96	4.04	4.01
Redmi				4.00	4.00
Samsung	4.20	4.02	3.80	4.00	3.99
KZ			3.93	4.00	3.97
Portronics	4.03		3.83		3.92
Blaupunkt	3.90		3.93	3.80	3.90
Boult	3.90		3.84	3.99	3.87
Concept			3.80	4.10	3.86
Mivi			3.73	3.90	3.81
boAt	3.84		3.70	3.96	3.79
Noise		3.90	3.79	3.73	3.77
JBL	3.90	4.02	3.65	3.60	3.75
pTron	3.95		3.76		3.67
Ambrane	3.64				3.64
ZEBRONICS	3.51		3.53	4.35	3.56
HAMMER	3.25		3.30	3.68	3.48
Amazon Basics	3.43		3.43		3.43

Average Discount by Brand



Discount vs Price Segment



Figure 10: Advanced Insights

## 8.2 Sample codes:

### 1. Sample SQL code for Data Cleaning

```
1  -- Standardizing brand names
2  UPDATE amazon_data
3  SET "Brand" = CASE
4      WHEN LOWER("Brand") = 'boat' THEN 'boAt'
5      WHEN LOWER("Brand") = 'oneplus' THEN 'OnePlus'
6      WHEN LOWER("Brand") = 'amazon' THEN 'Amazon Basics'
7      WHEN LOWER("Brand") = 'zebronics' THEN 'ZEBRONICS'
8      WHEN "Brand" IN ('Q', 'W20', 'Mustang') THEN 'Boult'
9      ELSE "Brand"
10 END;
11
12 -- Deduplication using ASIN (unique for each product)
13 CREATE TABLE cleaned_amazon_data AS
14 SELECT
15     "Title", "Brand", "Price", "MRP", "Discount",
16     "Rating", "Review_Count", "Prime", "ASIN", "URL"
17 FROM (
18     SELECT *, ROW_NUMBER() OVER (PARTITION BY "ASIN" ORDER BY "Title") as rn
19     FROM amazon_data
20     WHERE "ASIN" IS NOT NULL AND "ASIN" != ''
21 ) ranked
22 WHERE rn = 1;
23
24 -- Clean Price column
25 UPDATE cleaned_amazon_data
26 SET "Price" = REGEXP_REPLACE("Price", '[^0-9]', '', 'g')::INTEGER
27 WHERE "Price" IS NOT NULL AND "Price" ~ '[0-9]';
28
29 -- Clean Rating column
30 UPDATE cleaned_amazon_data
31 SET "Rating" = SPLIT_PART("Rating", ' ', 1)::FLOAT
32 WHERE "Rating" IS NOT NULL AND "Rating" ~ '[0-9\.]+';
```

### 2. Sample Python code for Feature Engineering and Analysis

```
1  # Feature engineering
2  df['Price_Range'] = pd.cut(df['Price'],
3      bins=[0, 500, 1500, 5000, float('inf')],
4      labels=['Budget', 'Mid-range', 'Premium', 'Luxury'])
5
6  df['Value_Score'] = (df['Rating'] / (df['Price'] / 1000)).round(3)
7
8  # Correlation analysis
9  numeric_cols = ['Price', 'Rating', 'Review_Count', 'Value_Score']
10 correlation_matrix = df[numeric_cols].corr()
11 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0)
12
13 # Random Forest regression
14 from sklearn.ensemble import RandomForestRegressor
15 from sklearn.model_selection import train_test_split
16 from sklearn.preprocessing import OneHotEncoder
17 from sklearn.compose import ColumnTransformer
18
19 preprocessor = ColumnTransformer([
20     ('brand', OneHotEncoder(handle_unknown='ignore'), ['Brand'])
21 ], remainder='passthrough')
22
23 X = df[['Brand', 'Rating', 'Review_Count']]
24 y = df['Price']
25 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
26
27 X_train_processed = preprocessor.fit_transform(X_train)
28 X_test_processed = preprocessor.transform(X_test)
29
30 rf = RandomForestRegressor(n_estimators=100, max_depth=20, min_samples_leaf=2,
31     max_features=0.8, random_state=42)
32 rf.fit(X_train_processed, y_train)
```

### 8.3 Data dictionary:

Column	Description
Title	Product title/name as listed on Amazon
Brand	Brand name (standardized, e.g., "boAt", "OnePlus")
Price	Current selling price (in ₹)
MRP	Maximum Retail Price (in ₹)
Discount	Discount percentage from MRP (integer)
Rating	Customer rating (out of 5, float)
Review_Count	Number of customer reviews (integer)
Prime	Whether the product is Amazon Prime eligible (TRUE/FALSE)
ASIN	Amazon Standard Identification Number (unique product identifier)
URL	Product page URL on Amazon
Savings	Amount saved from MRP (MRP - Price, in ₹)
Price_Segment	Price category: Budget, Mid-range, Premium, or Luxury (engineered feature)
Value_Score	Rating per ₹1000 spent (engineered feature: Rating / (Price/1000))

## 9. Conclusion

This project demonstrates a full-cycle, real-world data analytics workflow:

- Data collection, cleaning, and feature engineering
- Advanced analysis (statistical, clustering, regression)
- Business-focused visualization and reporting

The insights and dashboards produced here enable better strategic decisions in product positioning, pricing, and brand management for the Amazon headphones market.