

# Data Selection & Exploration

## 1. Dataset Description

### 1. Data Sources & Structure

- **Sources:** Amazon Reviews datasets for **Books, Watches, and Electronics**.
- **File Format:** Originally in TSV (tab-separated values) format; merged into a single dataset.
- **Schema Overview (common columns):**
  - marketplace (VARCHAR)
  - review\_id (BIGINT or VARCHAR)
  - customer\_id (BIGINT) – *If available for cross-category analysis*
  - product\_id (VARCHAR)
  - star\_rating (FLOAT or INT)
  - helpful\_votes (INT)
  - total\_votes (INT)
  - verified\_purchase (VARCHAR/BOOLEAN)
  - review\_headline (VARCHAR)
  - review\_body (TEXT)
  - review\_date (DATE)
  - category (VARCHAR) – *A custom field labeling each record as “Books,” “Watches,” or “Electronics.”*

### 2. Dataset Size & Volume

- The merged dataset exceeds **100,000 total records** across all three categories.
- Each category contributes thousands (or tens of thousands) of reviews, ensuring substantial coverage for analysis.

### 3. Data Types & Missing Values

- **Data Types:**
  - Numeric columns (e.g., star\_rating, helpful\_votes, total\_votes).
  - Text columns (e.g., review\_body, review\_headline).
  - Date column (review\_date).
  - Categorical columns (category, verified\_purchase).
- **Missing Values:**
  - Some reviews may lack certain fields (e.g., missing review\_headline or review\_body).
  - Inconsistent or null entries in helpful\_votes or total\_votes if reviewers did not receive or cast votes.
- **Initial Handling:**
  - Identified these missing values during a preliminary scan in **SQL Server Management Studio (SSMS)**.
  - Will address these in the next phase (ETL) to ensure data consistency.

#### 4. Key Attributes

- **review\_id** or **customer\_id** can serve as a unique identifier (depending on availability).
- **product\_id** links to the item being reviewed, which is crucial for cross-category analysis.
- **category** helps distinguish among books, watches, and electronics.

---

## 2. Business Problem

### Focus: Cross-Selling & Customer Overlap

- **Primary Question:**  
*Are there customers who purchase across multiple categories (books + watches + electronics)?*

- **Supporting Questions:**

1. *Which purchasing patterns exist across categories?*

- Do customers who buy certain electronics also tend to buy certain books or watches?

2. *Can we identify bundle opportunities or product recommendations?*

- Is there a correlation between specific electronics and watch purchases that might indicate a potential bundle?
- Are there popular book genres frequently bought alongside certain tech items?

- **Goal:**

Leverage the reviews dataset to uncover insights into cross-category behavior, ultimately guiding **marketing strategies, personalized recommendations, and inventory decisions**. By pinpointing overlaps, the business can target customers with more relevant promotions and optimize stock for high-demand bundle opportunities.