# Data Sources Documentation

## 1. Introduction:

### a. Brief Overview of the Datasets

There are three main datasets in the data warehouse project: Customers_FlipKart, Products_FlipKart, and Sales_FlipKart. The Customers collection keeps track of information about each customer's demographics and behavior, such as what they buy. The Products dataset has important details about products, like their category, price, and availability. The Sales Transactions dataset is the main fact table. It stores information about every purchase and uses unique identifiers to connect customers and goods.

### How I created my Dataset?

Using the Faker tool in Jupyter Notebook, I made the datasets for this project. I found the most important fields and relationships for a full retail data warehouse by studying e-commerce data structures on Kaggle. I used a star schema design with Sales as the fact table to make three datasets that are all linked to each other: Customers_FlipKart, Products_FlipKart, and Sales_FlipKart. I made realistic fake data using Faker, such as customer information, product details, and transaction records. Each dataset has more than 1000 records to make sure there is enough information for research.

### b. Purpose of the Datasets in the Data Warehouse

These datasets are meant to help people make decisions based on data by making it easy to store, find, and analyze sales and customer contacts. Businesses can learn about the demographics and behavior of their customers with the Customers dataset. The Products dataset, on the other hand, helps them keep track of their supplies, prices, and how well their products are doing. The Sales Transactions dataset is an important part of business intelligence and reporting because it shows trends in income, buying patterns, and operational efficiency.

### c. How Datasets are Connected

The three sets of data are linked together using the primary and foreign keys, creating a star schema in the data warehouse. In this case, the Sales Transactions dataset is the fact table, and the customer_id and product_id in the Customers dataset connect them. This structure lets you run quick and effective analytical questions, which lets you keep an eye on what customers buy, how much demand there is for products, and how much money is being made at different levels of detail.
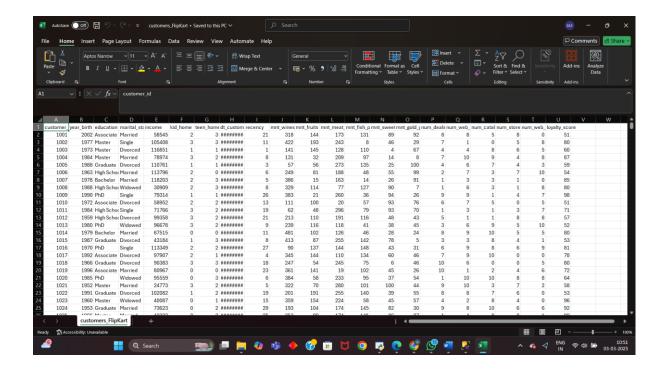
# 2. Dataset Descriptions:

## 2.1 Customers_FlipKart Dataset

The Customers dataset has information about customers' demographics and behaviors, which helps businesses look at buying patterns and divide customers into groups. A primary key (customer_id) is used to identify each customer individually. There are fields in the dataset for things like birth year, level of schooling, marital status, income, number of children, purchase history, and loyalty score. For customer segmentation, personalized marketing, and behavior research, this dataset is a must-have.

**Key Attributes:**

- product_id (Primary Key) – Unique identifier for each product.

- product_name – The name of the product.

- category – The category the product belongs to (e.g., Electronics, Furniture).

- price – The listed price of the product.

- stock_quantity – Available stock levels to monitor inventory.

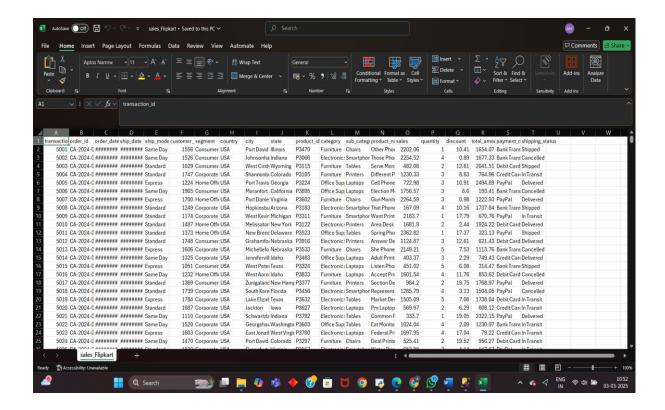- rating – Customer review rating for the product.

## 2.2 Sales_FlipKart Dataset

The Sales_FlipKart dataset is the Fact Table in the data center. It keeps track of all the purchases that customers make. It has information about the transaction, like the order details, buy date, customer ID, product ID, quantity, total amount, and payment method. This dataset is very important for tracking sales success, analyzing revenue, and predicting demand because it has foreign keys (customer_id and product_id) that connect customers and products.

**Key Attributes:**

- transaction_id (Primary Key) – Unique identifier for each transaction.

- customer_id (Foreign Key) – Links to the Customers dataset.

- product_id (Foreign Key) – Links to the Products dataset.

- order_date – Date of the purchase transaction.

- quantity – Number of units purchased in the transaction.

- total_amount – The final purchase amount after discounts.

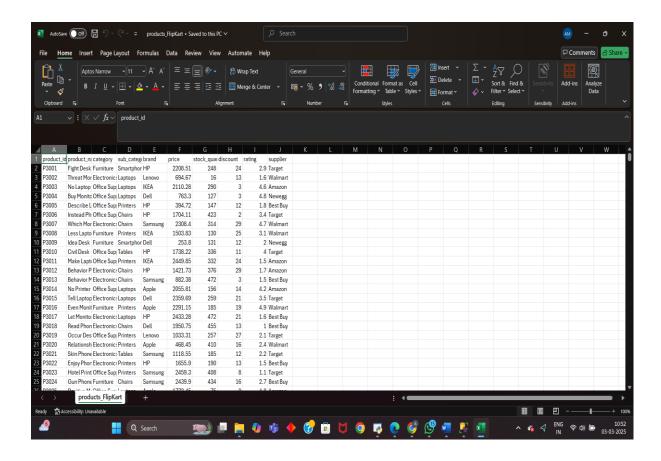- payment_method – Payment type used (Credit Card, PayPal, etc.).

## 2.3 Products_FlipKart Dataset

The Product_FlipKart dataset has information about items that are for sale, such as their categories, brands, prices, availability, and discounts. A primary key (product_id) is used to identify each product individually. This set of data helps companies keep track of their supplies, study how well their products are doing, and find the best ways to set prices.

**Key Attributes:**

- product_id (Primary Key) – Unique identifier for each product.

- product_name – The name of the product.

- category – The category the product belongs to (e.g., Electronics, Furniture).

- price – The listed price of the product.

- stock_quantity – Available stock levels to monitor inventory.

- rating – Customer review rating for the product.

# 3. Relationships between Datasets:

The Customers, Products, and Sales Transactions datasets are interconnected using primary and foreign keys, forming a star schema model that enables efficient analytical querying.

## 3.1 Customers_FlipKart Dataset (customer_id)

- The **primary key** customer_id uniquely identifies each customer.

- The **Sales_FlipKart dataset** contains customer_id as a **foreign key**, linking each transaction to a specific customer.

## 3.2 Products_FlipKart Dataset (product_id)

- The **primary key** product_id uniquely identifies each product.

- The **Sales_FlipKart dataset** contains product_id as a **foreign key**, linking each transaction to the product being purchased.

### 3.3 Sales_FlipKart Dataset (Fact Table)

- Acts as the **Fact Table**, linking **Customers and Products** datasets.

- Contains both customer_id and product_id as **foreign keys**, ensuring that each sale references a valid customer and product.

**Relationship Summary**

- **One-to-Many Relationship** (Customers → Sales Transactions):

  ➢ Each customer can make more than one purchase, but each purchase is only linked to one customer.

- **One-to-Many Relationship** (Products → Sales Transactions):

  ➢ Every product can be part of more than one transaction, but every transaction is only linked to one product.