| **Student Numbers:** Please list numbers of all group members | **5587248,5550544,2289197,5587165,5576953** |
|---|---|
| **Module Code:** | **IB9HP0** |
| **Module Title:** | **Data Management** |
| **Submission Deadline:** | **12:00 PM (UK Time) Monday 20th March 2024** |
| **Date Submitted:** | **19th March 2024** |
| **Word Count:** | |
| **Number of Pages:** | **72** |
| **Question Attempted:** (question number/title, or description of assignment) | **5** |
| **Have you used Artificial Intelligence (AI) in any part of this assignment?** | **Section 2 and 3** |

# Group 28 IB9HP0

## Introduction

In today's constantly expanding digital commerce landscape, having a deep understanding of data is essential for making well-informed decisions and strategizing for business success. This report offers an examination of an e-commerce data environment, from the initial creation of a database to the sophisticated analysis of its contents. Beginning with the E-R diagram as a foundation, it serves as the basis for constructing the entire SQL schema. For generating data, Python was utilized to produce artificial information. Subsequently, a data pipeline was established using GitHub to ensure smooth version control and transparent workflow management - keeping our dataset up-to-date and aligned with business expansion needs. Lastly, leveraging R along with tools such as dplyr and ggplot2 enabled us to conduct advanced analysis on the dataset revealing valuable trends and insights.

The link to our team's repository: [https://github.com/Isaacdexi/dm28]

# 1. Database Design and Implementation

## 1.1 Entity Relationship Diagram



1. *Customers Entity*: This represents the customers who use the platform. Attributes include personal information such as first name, last name, email, phone number, and customer ID. Customers also have an associated address that includes street, country, and ZIP code, as well as additional information like gender, date of birth, and platform (the platform through which the customer was acquired).

2. *Order Relationship*: It is linked to the Customers and Products entity. This shows that customers can place orders on products. Each order has an order number, date, payment method, quantity of items ordered, customer rating, and a review of the purchase.

3. *Shipment Entity*: Connected to the Order entity. This represents the shipping details of an order, including a shipment ID, delay in days, cost, and whether a refund was issued.

4. *Products Entity*: This represents the items that are for sale on the platform. Product attributes include a unique product ID, name, description, the number of views, price, and weight. Products are also linked to categories and have an inventory level.

5. *Product_Category Entity*: A junction table that associates products with categories, indicating that products can belong to one or more categories. Categories have a unique ID, name, and description.

6. *Sellers Entity*: Represents the sellers or suppliers on the platform. This includes a seller ID, address details, email, phone, and company name..

7. *Discount Entity*: Connected to the Products entity. This represents discounts that can be applied to products, with attributes like discount ID, start and end date, and the percentage of the discount.

Each rectangle represents an entity, which is a table in the database, while the ovals represent the attributes or fields within those tables. Diamonds represent relationships between entities, and the lines connecting the entities indicate how they are related. The 'M' and 'N' notation specifies the nature of the relationship

**Relationship Sets**



Some instances in a "Orders" relationship set of M:N relationship, between customer and products.

Some instances in a "Sells" relationship set of 1:N, between seller and product.

Customers, Orders, Products

In a many-to-many association involving customers, orders, and products, each customer has the ability to make multiple orders while each order can encompass multiple products. This arrangement provides flexibility in recording the interactions between customers and the items they purchase via orders. It also facilitates monitoring of customer preferences, buying history, and product popularity.

Sellers, Sells, Products

A one-to-many relationship exists among sellers, sales activities (sells), and products where each seller can vend numerous products; however, each product is vended by only one seller. Sellers are connected to the products they sell through their unique identifiers with the products being linked to sellers via the identifier belonging to the respective seller.



Some instances in a "Has" relationship set of N:1 relationship, between product and product category.

Some instances in a "Has" relationship set of 1:N relationship, between product and discount.

Products, Has, Product Category

In a many-to-one association between products and product categories, several items can be part of the same product category, while each item is associated with only one category. This structure helps organize items into specific categories, making it easier to classify products.

Product, Has, Discount

When it comes to the connection between products and discounts, there is a one-to-many correlation where each product can have multiple related discounts. However, each discount is only applicable to a single product. This arrangement aids in managing the various discounts offered within a product catalog.

Some instances in a "Order" relationship set of
M:1 relationship, between customers and
shipment.

Customer, Order, Shipment

In a one-to-many connection linking customers, orders, and shipments, a single customer can have multiple associated orders while individual shipments can fulfill various orders. This arrangement facilitates the effective supervision and monitoring of order processing for each customer.

**Challenges:**

The complexities primarily revolve around excessive entities at the outset, difficulties in normalization, ambiguity in identifying entities and relationships, scope creep, data integrity concerns, user requirements changes, lack of expertise, and conflicting requirements. These multifaceted challenges impede progress, leading to prolonged revisions and delays in schema and data generation. Effective communication, collaboration, and a thorough understanding of database design principles are essential to address these obstacles and ensure the creation of an accurate and efficient Entity-Relationship diagram.

## 1.2 SQL Database Schema Creation

### 1.2.1 Logical Schema

Customer (customer_id(primary key), first_name, last_name, gender, date_of_birth, email, phone, customer_street, customer_country, customer_zip_code, platform)

Product (product_id(primary key), product_name, product_description, price, weight, inventory, category_id (foreign key), seller_id (foreign key), product_views)

Category (category_id(primary key), p_category_id, cat_name, cat_description)

Discount (discount_id(primary key), discount_percentage, discount_start_date, discount_end_date, product_id (foreign key))

Order (order_number(primary key), payment_method, order_date, quantity, review, customer_id (foreign key), product_id (primary key), shipment_id (foreign key), customer_rating)

Shipment (shipment_id(primary key), shipment_delay_days, shipment_cost, order_number (foreign key), refund)

Sellers (seller_id(primary key), company_name, supplier_phone, supplier_email, seller_street, seller_country, seller_zip_code)

### 1.2.2 Physical Schema

```
sqlite3 Ecommerce.db
```

```
# Establish a connection to the SQLite database
database <- dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
```

The implementation of a robust physical schema essentially converts the conceptual design of the database into a tangible structure, comprising tables, columns, and constraints. It deals with the detailed conversion of logical entities, attributes, and relationships into concrete database tables, indices, and storage mechanisms. Accurate mapping of logical entities and attributes to their physical counterparts ensures efficient storage, retrieval, and integrity of data. This forms the backbone of the platform's data management infrastructure, enabling seamless operations and facilitating informed decision-making in the competitive landscape of online commerce

Customer Table

```
# If Customer table exist, drop it
if(dbExistsTable(database, "Customer")){
  dbExecute(database, "DROP TABLE Customer")
}

dbExecute(database, "CREATE TABLE 'Customer' (
    'customer_id' TEXT PRIMARY KEY,
```

```
8      'first_name' VARCHAR(50) NOT NULL,
9      'last_name' VARCHAR(50) NOT NULL,
10     'gender' VARCHAR(50) NOT NULL,
11     'date_of_birth' DATETIME NOT NULL,
12     'email' VARCHAR(50) NOT NULL,
13     'phone' VARCHAR(20) NOT NULL,
14     'customer_street' VARCHAR(50) NOT NULL,
15     'customer_country' VARCHAR(50) NOT NULL,
16     'customer_zip_code' VARCHAR(10) NOT NULL,
17     'platform'TEXT NOT NULL)")
```

The Customer entity, representing registered users, is translated into a database table named "Customer". Each attribute, such as customer_id, first_name, last_name, etc., corresponds to a column within this table. Data types are assigned to each column based on their nature (e.g., VARCHAR for textual data, DATETIME for date of birth).

Product Table

```
1  # If Product table exist, drop it
2  if(dbExistsTable(database, "Product")){
3    dbExecute(database, "DROP TABLE Product")
4  }
5
6  dbExecute(database,"CREATE TABLE  'Product' (
7      'product_id' TEXT PRIMARY KEY,
8      'product_name' VARCHAR(255) NOT NULL,
9      'price' DECIMAL(10,2) NOT NULL,
10     'product_description' TEXT NOT NULL,
11     'inventory' INT NOT NULL,
12     'weight' DECIMAL(10,2) NOT NULL,
13     'category_id' TEXT NOT NULL,
14     'seller_id' TEXT NOT NULL,
15     'product_views' INT NOT NULL,
16     FOREIGN KEY ('category_id') REFERENCES Category ('category_id'),
17     FOREIGN KEY ('seller_id') REFERENCES Sellers ('seller_id'))")
```

The Product entity, representing items available for sale, is mapped to the "Product" table in the database. Attributes like product_id, product_name, etc., are represented as columns within this table. Foreign key constraints are implemented for attributes like category_id and seller_id, ensuring referential integrity with the Category and Seller tables.

Category Table

```
1  # If Category table exist, drop it
2  if(dbExistsTable(database, "Category")){
3    dbExecute(database, "DROP TABLE Category")
4  }
5  dbExecute(database,"CREATE TABLE 'Category' (
6      'category_id' TEXT PRIMARY KEY,
7      'p_category_id' TEXT,
8      'cat_name' VARCHAR(255) NOT NULL,
9      'cat_description' TEXT NOT NULL)")
```

The Category entity, defining product categories, is instantiated as the "Category" table.
Attributes like category_id, cat_name, etc., are mapped to corresponding columns in this
table.

Discount Table

```
1  # If Discount table exist, drop it
2  if(dbExistsTable(database, "Discount")){
3    dbExecute(database, "DROP TABLE Discount")
4  }
5  dbExecute(database,"CREATE TABLE  'Discount' (
6      'discount_id' INT PRIMARY KEY,
7      'discount_percentage' DECIMAL(10,2) NOT NULL,
8      'discount_start_date' DATETIME NOT NULL,
9      'discount_end_date' DATETIME NOT NULL,
10     'product_id' INT NOT NULL,
11     FOREIGN KEY ('product_id') REFERENCES Product ('product_id')
12  )")
```

The Discount entity, representing discounts applicable to products, is materialized as the "Dis-
count" table. Attributes like discount_id, discount_percentage, etc., are mapped to columns
within this table. A foreign key constraint is implemented for the product_id attribute, en-
suring integrity with the Product table.

Order Table

```
1  # If Order table exists, drop it
2  if (dbExistsTable(database, "Order")) {
3    dbExecute(database, "DROP TABLE 'Order'")
4  }
5  dbExecute(database,"CREATE TABLE 'Order' (
6      'order_number' TEXT NOT NULL,
```

```
7      'payment_method' TEXT NOT NULL ,
8      'order_date' DATETIME NOT NULL ,
9      'quantity' INTEGER NOT NULL ,
10     'review' TEXT,
11     'customer_id' TEXT NOT NULL ,
12     'product_id' TEXT NOT NULL ,
13     'shipment_id' TEXT NOT NULL ,
14     'customer_rating' INT NOT NULL,
15     PRIMARY KEY ('order_number', 'product_id'),
16     FOREIGN KEY ('product_id') REFERENCES Product ('product_id'),
17     FOREIGN KEY ('customer_id') REFERENCES Customer ('customer_id'),
18     FOREIGN KEY ('shipment_id') REFERENCES Shipment ('shipment_id')
19  )")
```

The Order entity, capturing customer orders, is realized as the "Order" table. Each attribute, such as order_number, payment_method, etc., is represented as a column in this table.

Foreign key constraints are established for attributes like customer_id, product_id, and shipment_id, maintaining referential integrity with the Customer, Product, and Shipment tables, respectively.

Shipment Table

```
1   # If shipment table exist, drop it
2   if(dbExistsTable(database, "Shipment")){
3     dbExecute(database, "DROP TABLE 'Shipment'")
4   }
5
6   dbExecute(database,"CREATE TABLE 'Shipment' (
7     'shipment_id' TEXT PRIMARY KEY,
8     'shipment_delay_days' INT NOT NULL,
9     'shipment_cost' DECIMAL(10,2) NOT NULL,
10    'order_number' TEXT NOT NULL,
11    'refund' TEXT NOT NULL,
12    FOREIGN KEY ('order_number') REFERENCES `Order` ('order_number')
13  )")
```

The Shipment entity, representing shipment details, is translated into the "Shipment" table. Attributes like shipment_id, shipment_delay_days, etc., are mapped to columns in this table. A foreign key constraint is implemented for the order_number attribute, ensuring integrity with the Order table.

Seller Table

```
1  # If Sellers table exist, drop it
2  if(dbExistsTable(database, "Sellers")){
3    dbExecute(database, "DROP TABLE 'Sellers'")
4  }
5
6  dbExecute(database,"CREATE TABLE 'Sellers' (
7      'seller_Id' TEXT PRIMARY KEY,
8      'company_name' VARCHAR(100) NOT NULL ,
9      'supplier_phone' VARCHAR(20) NOT NULL,
10     'supplier_email' VARCHAR(100) NOT NULL UNIQUE,
11     'seller_Street' VARCHAR(255) NOT NULL,
12     'seller_country' VARCHAR(255) NOT NULL,
13     'seller_zip_code' VARCHAR(10) NOT NULL)")
```

The Sellers entity, representing platform vendors, is instantiated as the "Sellers" table. Attributes like seller_Id, company_name, etc., are represented as columns within this table.

**Challenges:**

While designing and implementing the schema for our e-commerce platform database, we encountered several challenges that required careful consideration and problem-solving.

*Complex Relationships*:

• Managing the many links between entities was one of the main difficulties. For instance, managing many-to-many relationships was necessary to build relationships between customers, orders, and items.

• It was difficult to conclude what can be an Entity, Attribute, or associative relationship. For instance, Order could be an Entity, Attribute as well as associative relationship depending on the scenario. However, we concluded that order is an associative relationship because it acts as a bridge between the customer and the product entities. In a typical e-commerce scenario, a customer places an order for one or more products. This establishes a many-to-many relationship between customers and products since a single customer can place multiple orders, and each order can contain multiple products.

*Normalization*: Achieving normalization while avoiding excessive data redundancy was another challenge. We had to brainstorm about various attributes aligned to different entities. For instance, we made two entities Product and Product Category to ensure the minimization of data redundancy.

*Performance Optimization*: Optimizing database performance, especially for complex queries was an important task. We ensured that the query was dynamic and ran efficiently in a loop without any halt. For instance, in our query, we added a drop at the start of every

creation query, and instead of the 'overwrite' function, we used 'append' to maintain the data consistency of the database.

*Data Integrity*: Data Integrity is important because it measures the accuracy, consistency, and reliability of data stored in a database. We faced challenges while deciding on the primary key of Order Relationship because we were particularly concerned about the need to uniquely identify each order and its associated products. We needed a key that would uniquely identify each order while also because an order could consist of multiple products. After careful consideration, we concluded that a single attribute alone, such as 'order_number', may not be sufficient to uniquely identify orders, especially in scenarios where multiple orders with the same number could exist if placed by different customers or at different times. To address this challenge, we opted for a composite primary key consisting of 'order_number' and 'product_id'.

## 2. Data Generation and Management

### 2.1 Synthetic Data Generation

During this stage, the team employed Python V3.8.1 and Spyder IDE 5.4.3 to create synthetic data in a CSV file using ChatGPT. The task was aided by the Faker and random libraries to produce a highly authentic eCommerce dataset. Additionally, ChatGPT was used for generating textual data lists, function enhancements, and debugging assistance. The code can be seen in Appendix 2



The above diagram depicts the comprehensive approach used to generate data, involving an ongoing cycle between schema design and data generation. This iterative process ensures the creation of highly realistic datasets that accurately reflect an eCommerce database. A specific list is needed to store various textual data, such as reviews, which will be randomized based on

purchase orders. Meanwhile, details like product descriptions will be linked to their respective products. ChatGPT supports this process by automating the arduous task of generating large amounts of textual information. Additionally, it aids in organizing codes coherently and enables real-time improvisation, thereby enhancing code interpretability and streamlining data generation efficiency.



### Sequencing Logic between Entities

The logic of sequencing as illustrated can be segmented into 3 tranches, in which, dependencies is considered. The first would be to generate customer and seller information, so category can be created followed by products. Order information is dependent on product and customer, so it makes logical sense to generate the table after the first and second tranches. Shipment is dependent on order, which can be created as the last table alongside discount. It is important to generate this systematically, as certain attributes such as customer_id would need to be first created as it is included in the order table.

### Attribute Data Generation Logic

Faker library is able to do at least 70% of the data generation work, which leaves the remaining to the use of ChatGPT, random library and logic binding of variables (focus of this section). All tables have such operations except for product and discount.

| Seller | |
|---|---|
| Attribute that requires binding | Attributes |
| email | company_name |
| seller_country<br>supplier_phone | Country code function and inter-bind<br>between seller_country<br>and supplier_phone |

Since Faker is unable to match email domais, data generated for email variables would be based on a [random point-of-contact's name + '@' + extraction of company first name + '-' + '.com', to reflect true email domains in real business setting. As for country, it extracts from

another function and list of country code, then based on that, phone number randomise 10 digits attached with the "(+country code)".

| Customer | |
|---|---|
| Attribute that requires binding | Attributes |
| email | first_name |
| | last_name |
| customer_country | Country code function and inter-bind |
| phone | between customer_country and phone |

Since Faker is unable to match email names, data generated for email variables would be based on extracting [first_name + '_' + last_name + '@gmail.com'. As for country, it extracts from another function and list of country code, then based on that, phone number randomise 10 digits attached with the "(+country code)".

| Category | |
|---|---|
| Attribute that requires binding | Attributes |
| p_category_id | cat_name |

Parent category id is based on the mapping of category names since, it is subsumed to a main category, for instance, 'Sportswear' is under 'Sports and Outdoors'.

| Order | |
|---|---|
| Attribute that requires binding | Attributes |
| payment_method | order_number |

As a single order number can include multiple products, it is essential to bind the payment method to ensure that it makes logical sense that one order can only have the same payment method. Also, to reflect the M:N relationship between customers and products, the order table would have records that reflect purchases with multiple products.

Overall, certain simplifications are also made, such as a single customer only has one delivery address, an order date would be the same as shipment data, where it is immediately dispatched. Also, the number of records should be tallied in the shipment table, where it references the latest order number in the order table for the creation of shipping_id.

Prompt Strategy

The team defined the ChatGPT prompt strategy into two main fronts, the contextual phase, and the output management stage. To ensure precise output, providing contextual information is crucial, for instance: "This is a database project, which requires the generation of datasets that reflect the eCommerce database environment". Following which, specifying the role is important for accuracy too, "I am a data engineer working on this project". Input definition would constitute the granular details and parameters, "The dataset should be in a CSV file, using the Faker library in python, I would like X1, X2, X3, X3....Xn variables, such that X1 should take on the format of a unique identifier, "c00001", and so on sequentially. In total, there should be 500 rows of data."

Generating multiple outputs also provide us more flexibility in selecting the most applicable codes to the project, and combined with interpretation prompts, "Please explain the codes", we can conduct testing of the data generation on our local environment more intuitively. We can create review statements prompts for improvisation, "X2 variable is not supposed to be an integer, please make amendments", and if errors occur, simply pasting the codes and prompting it to fix the errors would ensure that our data generation is executed smoothly.

**Challenges:**

| A | B | C | |
|---|---|---|---|
| product_id | prod_name | price | description |
| p00001 | Digital Air Fryer | 365.6296 | Front yourself large fall discuss. Material situation prepare car. |
| p00002 | Foldable Electric Scoote | 76.39194 | Quality support picture once. Better left control raise gas. |
| p00003 | Portable Bluetooth Spea | 487.4114 | Personal star guy recently morning decide father. Control professional assume truth. |

Although Faker library is considerably powerful in creating realistic datasets, there are certain limitations, particularly, the inability to match records correctly. This is especially striking when generating review and product descriptions, and as shown above, the description is totally not relevant to the product at all.

```
product_list = {
    "Electronics": [
        {"name": "SmartHome Hub", "description": "A central control hub for all your smart home devices."},
        {"name": "Wireless Bluetooth Earbuds", "description": "Enjoy wireless freedom with high-quality audio."},
        {"name": "4K Ultra HD Smart TV", "description": "Bring the cinema experience home with our 4K Ultra HD Smart TV, d
    ],
    "Home and Kitchen": [
        {"name": "Smart Coffee Maker", "description": "Brew your favorite coffee with smart features."},
        {"name": "Non-Stick Cookware Set", "description": "Premium cookware for your kitchen."},
        {"name": "Smart Mini Refrigerator", "description": "Keep your food fresh and organized with our Smart Refrigerator
    ],
    "Sports and Outdoors": [
        {"name": "Fitness Tracker", "description": "Track your fitness activities with this advanced device."},
        {"name": "Camping Tent", "description": "Explore the outdoors with a durable camping tent."},
        {"name": "Waterproof Hiking Boots", "description": "Conquer any trail with our Waterproof Hiking Boots, keeping you
    ],
```

The team overcame the challenge through the use of ChatGPT, to create a list of products and relevant descriptions generated by ChatGPT, then store and integrate it in the Python codes with Faker. This requires multiple rounds of as well as multiple trial and error before arriving at a satisfactory solution.

## 2.2 Data Import and Quality Assurance

### 2.2.1 Data Loading

```
1  setwd("/cloud/project")
2  Customer <- readr::read_csv("customer.csv")
3  Category <- readr::read_csv("category.csv")
4  Sellers <- readr::read_csv("seller.csv")
5  Product <- readr::read_csv("product.csv")
6  Discount <- readr::read_csv("discount.csv")
7  Shipment <- readr::read_csv("shipment.csv")
8  Order <- readr::read_csv("order.csv")
```

### 2.2.2 Data Validation

Customer Data Validation

```
1  ## Validation for customer data
2
3  # Function to check if a datetime is in the desired format ("%Y-%m-%d %H:%M:%S")
4  is_datetime_format <- function(datetime_string) {
5    tryCatch({
6      as.POSIXlt(datetime_string, format = "%Y-%m-%d %H:%M:%S")
7      TRUE
8    }, error = function(e) {
```

```
 9       FALSE
10     })
11   }
12
13   # Check if dates are in the desired format, if not, convert them
14   for (i in 1:nrow(Customer)) {
15     if (!is_datetime_format(Customer$date_of_birth[i])) {
16       Customer$date_of_birth[i] <- as.POSIXct(Customer$date_of_birth[i], format = "%Y-%m-%d %H
17     }
18   }
19
20
21   # Perform the rest of the validation checks
22   missing_values <- apply(is.na(Customer), 2, sum)
23
24   # Check unique customer IDs
25   if (length(unique(Customer$customer_id)) != nrow(Customer)) {
26     print("Customer ID is not unique.")
27   }
28
29   # Check data types for first_name and last_name
30   if (!all(sapply(Customer$first_name, is.character)) || !all(sapply(Customer$last_name, is.cha
31     print("First name and last name should be character.")
32   }
33
34   # Check valid gender values
35   valid_genders <- c("Male", "Female", "Other")
36   if (any(!Customer$gender %in% valid_genders)) {
37     print("Gender should be Male, Female, or Other.")
38   }
39
40   # Check email format
41   if (any(!grepl("^\\S+@\\S+\\.\\S+$", Customer$email))) {
42     print("Invalid email format")
43   }
44
45   # Regular expressions for phone number formats of Belgium, China, France, United Kingdom, Uni
46   phone_regex <- "^\\(\\+\\d+\\)\\d{10}$"
47
48   # Check phone number format for specific countries
49   invalid_phone_indices <- which(!grepl(phone_regex, Customer$phone))
50   if (length(invalid_phone_indices) > 0) {
```

```
51    print("Invalid phone numbers:")
52    print(Customer[invalid_phone_indices, ])
53  }
54
55  # Regular expressions for zip code formats of Belgium, China, France, United Kingdom, United
56  zip_regex <- c(
57    "^[0-9]{4}$",  # Belgium
58    "^[0-9]{6}$",  # China
59    "^[0-9]{5}$",  # France
60    "^[A-Z]{2}[0-9]{1,2}[A-Z]? [0-9][A-Z]{2}$",  # United Kingdom
61    "^[0-9]{5}-[0-9]{4}$"    # United States
62  )
63
64  # Check zip code format for specific countries
65  invalid_zip_indices <- which(!grepl(paste(zip_regex, collapse = "|"), Customer$customer_zip_
66  if (length(invalid_zip_indices) > 0) {
67    print("Invalid zip codes:")
68    print(Customer[invalid_zip_indices, ])
69  }
70
71  # Check platform values
72  valid_platforms <- c("Referral", "Instagram", "Facebook", "Others")
73  if (any(!Customer$platform %in% valid_platforms)) {
74    print("Invalid platform values.")
75  }
76
77  # If no errors are found, print a message indicating that the data is valid
78  if (!any(is.na(missing_values)) &&
79      length(unique(Customer$customer_id)) == nrow(Customer) &&
80      all(sapply(Customer$first_name, is.character)) &&
81      all(sapply(Customer$last_name, is.character)) &&
82      all(Customer$gender %in% valid_genders) &&
83      all(grepl("^\\S+@\\S+\\.\\S+$", Customer$email)) &&
84      length(invalid_phone_indices) == 0 &&
85      length(invalid_zip_indices) == 0 &&
86      all(Customer$platform %in% valid_platforms)) {
87    print("Customer Data is valid. Loading data into the database...")
88    RSQLite::dbWriteTable(database, "Customer", Customer, append = TRUE)
89    # Load the data into the database
90  } else {
91    print("Data is not valid. Please correct the errors.")
92  }
```

In the data validation process for the "Customer" dataset, It starts by making sure that the birthdates of customers are in the correct format, so they can be properly analyzed. Then, it checks for common issues like duplicate customer IDs, ensuring that each customer is unique in the dataset. It also verifies that names are stored as text, and gender is recorded as either "Male", "Female", or "Other". Email addresses are checked to ensure they follow a valid format. The script also examines phone numbers and zip codes, making sure they match the expected patterns for different countries. If everything checks out, it gives the green light to load the data into the database; otherwise, it flags any errors that need attention.

Discount Data Validation

```r
# Function to check if date is in the desired format
is_datetime_format <- function(x) {
  tryCatch({
    as.POSIXlt(x, format = "%Y-%m-%d %H:%M:%S")
    TRUE
  }, error = function(e) {
    FALSE
  })
}

# Convert discount_start_date and discount_end_date to desired format if not already in that
if (!all(sapply(Discount$discount_start_date, is_datetime_format))) {
  Discount$discount_start_date <- as.POSIXlt(Discount$discount_start_date, format = "%Y-%m-%
}

if (!all(sapply(Discount$discount_end_date, is_datetime_format))) {
  Discount$discount_end_date <- as.POSIXlt(Discount$discount_end_date, format = "%Y-%m-%d %H
}

# Check for missing values in Discount dataframe
na_disc <- apply(is.na(Discount), 2, sum)

# Validate discount_percentage, discount_start_date, and discount_end_date data types
valid_decimal <- function(x) {
  !is.na(as.numeric(x))
}

valid_datetime <- function(x) {
  !is.na(as.POSIXlt(x))
}

# Check discount percentage range (assuming it's between 0 and 100)
```

```r
33   if (any(Discount$discount_percentage < 0 | Discount$discount_percentage > 100) ||
34       !all(sapply(Discount$discount_percentage, valid_decimal))) {
35     print("Invalid discount percentage.")
36   }
37
38   # Check discount dates
39   if (any(Discount$discount_start_date >= Discount$discount_end_date) ||
40       !all(sapply(Discount$discount_start_date, valid_datetime)) ||
41       !all(sapply(Discount$discount_end_date, valid_datetime))) {
42     print("Discount start date should be before the end date.")
43   }
44
45   # Check if discount_id is unique
46   if (any(duplicated(Discount$discount_id))) {
47     print("Duplicate discount IDs found.")
48   }
49
50   # Check if product_id exists in Product table
51   if (any(!Discount$product_id %in% Product$product_id)) {
52     print("Invalid product IDs. Some product IDs do not exist in the Product table.")
53   }
54
55   # If no errors are found, print a message indicating that the data is valid
56   if (!any(is.na(na_disc)) &&
57       all(Discount$discount_percentage >= 0 & Discount$discount_percentage <= 100) &&
58       all(Discount$discount_start_date < Discount$discount_end_date) &&
59       !any(duplicated(Discount$discount_id)) &&
60       all(Discount$product_id %in% Product$product_id) &&
61       all(sapply(Discount$discount_percentage, valid_decimal)) &&
62       all(sapply(Discount$discount_start_date, valid_datetime)) &&
63       all(sapply(Discount$discount_end_date, valid_datetime))) {
64     print("Discount data is valid. Loading data into the database...")
65     RSQLite::dbWriteTable(database, "Discount", Discount, append = TRUE)
66     # Load the data into the database
67   } else {
68     print("Data is not valid. Please correct the errors.")
69   }
```

The validation process for the "Discount" dataset begins with checking if the format of discount start and end dates conforms to the specified datetime format ("%Y-%m-%d %H:%M:%S"). If needed, it converts them accordingly. Subsequently, it detects any missing values in the discount dataset. Then, it verifies the data types of discount percentage, start date, and end date.

Furthermore, it ensures that the discount percentage is within the expected range (0 to 100) and validates that the start date precedes the end date. The process also includes confirming unique discount IDs and validating product IDs referenced in the discount table against those existing in a separate product table. If all checks pass successfully, then script affirms validity of data for loading into database; otherwise prompts correction before proceeding further.

Order Data Validation

```r
na_order <- apply(is.na(Order), 2, sum)

# Check quantity (assuming it should be a positive integer)
if (any(Order$quantity <= 0)) {
  print("Invalid quantity.")
}

# Check customer rating (assuming it should be between 1 and 5)
if (any(Order$customer_rating < 1 | Order$customer_rating > 5)) {
  print("Invalid customer rating.")
}

# Check if product_id exists in Product table
if (any(!Order$product_id %in% Product$product_id)) {
  print("Invalid product IDs. Some product IDs do not exist in the Product table.")
}

# Check if customer_id exists in Customer table
if (any(!Order$customer_id %in% Customer$customer_id)) {
  print("Invalid customer IDs. Some customer IDs do not exist in the Customer table.")
}

# Check if shipment_id exists in Shipment table
if (any(!Order$shipment_id %in% Shipment$shipment_id)) {
  print("Invalid shipment IDs. Some shipment IDs do not exist in the Shipment table.")
}

# Check uniqueness based on primary key (order_number, customer_id, product_id)
if (any(duplicated(Order[c("order_number", "customer_id", "product_id")]))) {
  print("Duplicate records found based on order_number, customer_id, and product_id.")
}

# Check order date format and range
if (any(!is_datetime_format(Order$order_date))) {
  # Convert order date to the desired format if not already
```

```
36    Order$order_date <- as.POSIXct(Order$order_date, format = "%Y-%m-%d %H:%M:%S", tz = "UTC")
37  }
38
39  # If no errors are found, print a message indicating that the data is valid
40  if (!any(is.na(na_order)) &&
41      all(Order$quantity > 0) &&
42      all(Order$customer_rating >= 1 & Order$customer_rating <= 5)&&
43      all(Order$product_id %in% Product$product_id) &&
44      all(Order$customer_id %in% Customer$customer_id) &&
45      all(Order$shipment_id %in% Shipment$shipment_id) &&
46      !any(duplicated(Order[c("order_number", "customer_id", "product_id")])) &&
47      all(is_datetime_format(Order$order_date))) {
48    print("Order data is valid. Loading data into the database...")
49    RSQLite::dbWriteTable(database, "Order", Order, append = TRUE)
50    # Load the data into the database
51  } else {
52    print("Order data is not valid. Please correct the errors.")
53  }
```

The validation procedure for the "Order" dataset begins by checking for missing values in the order dataset. It then ensures that the quantity of items ordered is a positive integer and that customer ratings fall within the range of 1 to 5. Additionally, it verifies if the product IDs referenced in the order table exist in the product table, and similarly for customer IDs in the customer table and shipment IDs in the shipment table. The script also checks for uniqueness based on the primary key consisting of order number, customer ID, and product ID. If all validation checks pass without errors, the script confirms the validity of the order data and suggests loading it into the database; otherwise, it prompts to correct the errors before proceeding.

Product Category Data Validation

```
1   na_prod_cat <- apply(is.na(Category), 2, sum)
2
3   # Ensure "category_id" values are unique
4   if (length(unique(Category$category_id)) != nrow(Category)) {
5     print("category_id values are not unique.")
6   }
7
8   # Check length of "cat_name"
9   if (any(nchar(Category$cat_name) > 255)) {
10    print("cat_name exceeds 255 characters.")
11  }
```

```
12
13   # Check data type of each column
14   if (!all(sapply(Category$category_id, is.character)) ||
15       !all(sapply(Category$cat_name, is.character)) ||
16       !all(sapply(Category$cat_description, is.character))) {
17     print("Invalid data type for one or more columns.")
18   }
19
20   # If no errors are found, print a message indicating that the data is valid
21   if (!any(is.na(na_prod_cat)) &&
22       length(unique(Category$category_id)) == nrow(Category) &&
23       !any(nchar(Category$cat_name) > 255) &&
24       all(sapply(Category$category_id, is.character)) &&
25       all(sapply(Category$cat_name, is.character)) &&
26       all(sapply(Category$cat_description, is.character))) {
27     print("product_category data is valid. Loading data into the database...")
28     RSQLite::dbWriteTable(database, "Category", Category, append = TRUE)
29     # Load the data into the database
30   } else {
31     print("product_category data is not valid. Please correct the errors.")
32   }
```

First, it checks if there are any missing values in the dataset. Then, it ensures that each category has a unique identifier. Next, it examines the length of category names to make sure they're not too long, as this could cause issues with storing the data. After that, it verifies the data type of each column to ensure they are all in the expected format. If everything checks out, it prints a message confirming the validity of the product category data and suggests loading it into the database. However, if any issues are detected during the validation process, it prompts the user to correct the errors before proceeding with data loading.

Products Data Validation

```
1    # Function to check if a value is decimal
2    valid_decimal <- function(x) {
3      !is.na(as.numeric(x))
4    }
5
6    # Function to check if a value is an integer
7    valid_integer <- function(x) {
8      !is.na(as.integer(x))
9    }
10
11   na_Product <- apply(is.na(Product), 2, sum)
```

```r
12
13  # Ensure "product_id" values are unique
14  if (length(unique(Product$product_id)) != nrow(Product)) {
15    print("product_id values are not unique.")
16  }
17
18  # Check length of "product_name"
19  if (any(nchar(Product$product_name) > 255)) {
20    print("product_name exceeds 255 characters.")
21  }
22
23  if (any(!Product$category_id %in% Category$category_id)) {
24    print("Invalid category IDs. Some category IDs do not exist in the product_category table.'
25  }
26
27  if (any(!Product$seller_id %in% Sellers$seller_id)) {
28    print("Invalid seller IDs. Some seller IDs do not exist in the Sellers table.")
29  }
30
31  # Check if inventory and product views are integers
32  if (any(!sapply(Product$inventory, valid_integer)) || any(!sapply(Product$product_views, val:
33    print("Inventory and product views should be integers.")
34  }
35
36  # Check if price and weight are decimal
37  if (any(!sapply(Product$price, valid_decimal)) || any(!sapply(Product$weight, valid_decimal)
38    print("Price and weight should be decimal values.")
39  }
40
41  # If no errors are found, print a message indicating that the data is valid
42  if (!any(is.na(na_Product)) &&
43      length(unique(Product$product_id)) == nrow(Product) &&
44      !any(nchar(Product$product_name) > 255) &&
45      all(Product$category_id %in% Category$category_id) &&
46      all(Product$seller_id %in% Sellers$seller_id) &&
47      all(sapply(Product$inventory, valid_integer)) &&
48      all(sapply(Product$product_views, valid_integer)) &&
49      all(sapply(Product$price, valid_decimal)) &&
50      all(sapply(Product$weight, valid_decimal))) {
51    print("Product data is valid. Loading data into the database...")
52    RSQLite::dbWriteTable(database, "Product", Product, append = TRUE)
53    # Load the data into the database
```

```
54  } else {
55    print("Product data is not valid. Please correct the errors.")
56  }
```

This script begins by defining two functions to check if a value is a decimal or an integer. Then, it checks for missing values in the product dataset. It ensures that each product has a unique identifier and checks the length of product names to ensure they don't exceed 255 characters. Additionally, it verifies that the category IDs referenced in the product table exist in the product category table and that the seller IDs exist in the sellers table. Furthermore, it checks if inventory and product views are integers and if price and weight are decimal values. If all validation checks pass without errors, the script confirms the validity of the product data and suggests loading it into the database. Otherwise, it prompts to correct the errors before proceeding with data loading.

Seller Data Validation

```
1   library(stringr)
2   na_sellers <- apply(is.na(Sellers), 2, sum)
3
4   # Ensure "seller_Id" values are unique
5   if (length(unique(Sellers$seller_id)) != nrow(Sellers)) {
6     print("seller_Id values are not unique.")
7   }
8
9   # Check length of "company_name"
10  if (any(nchar(Sellers$company_name) > 100)) {
11    print("company_name exceeds 100 characters.")
12  }
13
14  # Check email format
15  invalid_emails <- which(!str_detect(Sellers$supplier_email, "\\b[A-Za-z0-9._%+-]+@[A-Za-z0-9
16  if (length(invalid_emails) > 0) {
17    print("Invalid email addresses:")
18    print(Sellers[invalid_emails, ])
19  }
20
21  # If no errors are found, print a message indicating that the data is valid
22  if (!any(is.na(na_sellers)) &&
23      length(unique(Sellers$seller_id)) == nrow(Sellers) &&
24      !any(nchar(Sellers$company_name) > 100) &&
25      length(invalid_emails) == 0) {
26    print("Sellers data is valid. Loading data into the database...")
```

24

```
27    RSQLite::dbWriteTable(database, "Sellers", Sellers, append = TRUE)
28    # Load the data into the database
29  } else {
30    print("Sellers data is not valid. Please correct the errors.")
31  }
```

Initially, missing values are identified within the dataset. Then, the uniqueness of "seller_Id" values is confirmed to ensure each identifier is distinct. Moreover, the length of "company_name" is examined to ensure it does not exceed 100 characters. Additionally, email addresses are validated for correct formatting, with any invalid entries flagged for review. If no errors are detected and all validation criteria are met, the "Sellers" data is deemed valid and ready for database loading. However, should any discrepancies arise, corrective measures must be taken before proceeding with data integration to maintain data integrity.

Shipment Data Validation

```
1   na_shipment <- sapply(Shipment, function(x) sum(is.na(x)))
2
3   # Ensure "shipment_id" values are unique
4   if (length(unique(Shipment$shipment_id)) != nrow(Shipment)) {
5     print("shipment_id values are not unique.")
6   }
7
8   # Validate "refund" column
9   valid_refunds <- c("Yes", "No")
10  if (!all(Shipment$refund %in% valid_refunds)) {
11    print("Invalid values in the 'refund' column.")
12  }
13
14  # Validate "shipment_delay_days" and "shipment_cost" columns
15  if (any(Shipment$shipment_delay_days <= 0) || any(Shipment$shipment_cost <= 0)) {
16    print("shipment_delay_days and shipment_cost should be positive numbers.")
17  }
18
19  # Ensure that "shipment_delay_days" is an integer
20  if (any(!as.integer(Shipment$shipment_delay_days) == Shipment$shipment_delay_days)) {
21    print("shipment_delay_days should be integers.")
22  }
23
24  # Ensure that all "order_number" values exist in the "Order" table
25  order_numbers <- unique(Shipment$order_number)
26  if (!all(order_numbers %in% Order$order_number)) {
```

```
27    print("Some order numbers do not exist in the 'Order' table.")
28  }
29
30  # If no errors are found, print a message indicating that the data is valid
31  if (all(na_shipment == 0) &&
32      length(unique(Shipment$shipment_id)) == nrow(Shipment) &&
33      all(Shipment$refund %in% valid_refunds) &&
34      all(Shipment$shipment_delay_days > 0) &&
35      all(Shipment$shipment_cost > 0) &&
36      all(as.integer(Shipment$shipment_delay_days) == Shipment$shipment_delay_days) &&
37      all(order_numbers %in% Order$order_number)) {
38    print("Shipment data is valid. Loading data to database ...")
39    RSQLite::dbWriteTable(database, "Shipment",Shipment, append = TRUE)
40    # Load the data into the database
41  } else {
42    print("Shipment data is not valid. Please correct the errors.")
43  }
```

The script is designed to check if our shipment dataset is in good shape. It first counts how many missing values we have in each column, then ensures that every shipment has a unique ID. Next, it looks at the "refund" column to make sure it only has "Yes" or "No" values. It also checks that our shipment delay days and costs are positive numbers, with the delay days being whole numbers. Lastly, it confirms that all the order numbers mentioned in the shipment dataset actually exist in our orders. If everything checks out, it gives the green light to load the shipment data into our database. But if there are issues, it asks us to fix them before moving forward with the data loading process.

Validation Test

```
+     Customer$date_of_birth[i] <- as.POSIXct(Customer$date_of_birth[i], format = "%Y-%m-%d %H:%M:%S")
+   }
+ }
There were 50 or more warnings (use warnings() to see the first 50)
>
>
> # Perform the rest of the validation checks
> missing_values <- apply(is.na(Customer), 2, sum)
>
> # Check unique customer IDs
> if (length(unique(Customer$customer_id)) != nrow(Customer)) {
+   print("Customer ID is not unique.")
+ }
[1] "Customer ID is not unique."
Warning message:
Unknown or uninitialised column: `customer_id`.
>
> # Check data types for first_name and last_name
> if (!all(sapply(Customer$first_name, is.character)) || !all(sapply(Customer$last_name, is.character))) {
+   print("First name and last name should be character.")
+ }
Warning messages:
1: Unknown or uninitialised column: `first_name`.
2: Unknown or uninitialised column: `last_name`.
>
> # Check valid gender values
> valid_genders <- c("Male", "Female", "Other")
> if (any(!Customer$gender %in% valid_genders)) {
+   print("Gender should be Male, Female, or Other.")
+ }
Warning message:
Unknown or uninitialised column: `gender`.
>
> # Check email format
> if (any(!grepl("^\\S+@\\S+\\.\\S+$", Customer$email))) {
+   print("Invalid email format")
```

From the validation results, it is evident that attempts were made to input duplicate values, and the algorithm was able to identify this.

```
<0 rows> (or 0-length row.names)
> cust_result <- RSQLite::dbGetQuery(database, "SELECT * FROM Customer")
> print(cust_result[c("customer_id", "first_name")])
[1] customer_id first_name
<0 rows> (or 0-length row.names)
>
```
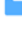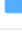
Afterward, we attempted to verify if the data had been stored in the database. Since it was not stored, the code correctly prevents invalid data from being sent to the database.

**Challenges:**

One of the primary challenges in crafting data validation rules lies in achieving the delicate balance between specificity and generality. These rules must be finely tuned to detect errors effectively without overly restricting legitimate data variations, demanding a profound grasp of the business context to identify critical integrity constraints while accommodating genuine data fluctuations. Moreover, determining how to address invalid data introduces further complexity, requiring careful consideration of factors like data significance, downstream implications, and regulatory compliance. Furthermore, ensuring that validation rules remain dynamic and adaptable to evolving business needs and data trends adds another layer of intricacy, underscoring the importance of ongoing collaboration among teams and iterative refinement of schema and data generation processes.

# 3. Data Pipeline Generation

## 3.1 GitHub Repository and Workflow Setup

| | | |
|---|---|---|
| 📁 .github/workflows | commit all db cust new data | 13 hours ago |
| 📁 Dataset | add commit dataset | 20 hours ago |
| 📁 DatasetTest | commit db add all test | 12 hours ago |
| 📁 R | db commit changes product new test | 11 hours ago |
| 📁 figures | Create placeholder | 20 hours ago |
| 📄 .gitignore | Initial commit | yesterday |
| 📄 Ecommerce.db | commit db add all test | 12 hours ago |
| 📄 README.md | Initial commit | yesterday |
| 📄 dm28.Rproj | demo | 18 hours ago |

For version control and data pipeline management purposes, the team created a Github repository, which consists of multiple folders, Dataset, stores our initial database CSV, and DatasetTest, for testing if there are new updates made to the existing database. In R, the team stores 2 Rscripts, one for the schema, validation and analysis, while another one for functions of loading new datasets.

```
1    name: ETL workflow for group 28
2
3    on:
4      schedule:
5        - cron: '0 */3 * * *' # Run every 3 hours
6      push:
7        branches: [ main ]
8
9    jobs:
10     build:
11       runs-on: ubuntu-latest
12
13       steps:
14         - name: Checkout code
15           uses: actions/checkout@v2
16
17         - name: Setup R environment
18           uses: r-lib/actions/setup-r@v2
19           with:
20             r-version: '4.2.0'
21
22         - name: Cache R packages
23           uses: actions/cache@v2
24           with:
25             path: ${{ env.R_LIBS_USER }}
26             key: ${{ runner.os }}-r-${{ hashFiles('**/lockfile') }}
27             restore-keys: |
28               ${{ runner.os }}-r-
29
30         - name: Install packages
```

28

```
31          if: steps.cache.outputs.cache-hit != 'true'
32          run: |
33            Rscript -e 'install.packages(c("ggplot2","dplyr","readr","RSQLite", "DBI","stringr","lubridate"))'
34
35        - name: Execute R script
36          run: |
37            Rscript R/schema_validation_plots.R
38
39        - name: Execute R script New Data Load
40          run: |
41            Rscript R/new_data_load.R
42
43        - name: Add files
44          run: |
45            git config --local --unset-all "http.https://github.com/.extraheader"
46            git config --global user.email "isaacchuadx@gmail.com"
47            git config --global user.name "Isaacdexi"
48            git add --all figures/
49
50        - name: Check for changes in Ecommerce.db
51          id: check_changes
52          run: |
53            git diff --quiet --exit-code -- Ecommerce.db || echo "Ecommerce.db has changed"
54
55        - name: Commit files
56          if: steps.check_changes.outputs.return == 'Ecommerce.db has changed'
57          run: |
58            git add Ecommerce.db
59            git commit -m "Update Ecommerce.db"
60
61        - name: Push changes
62          uses: ad-m/github-push-action@v0.6.0
63          with:
64              github_token: ${{ secrets.DM_28 }}
65              branch: main
```

Workflow file, "etl.yaml", is also created, and this provides the main anchor to automate the
process of updating the database. It can be seen that the workflow is set by the team to
execute the Rscripts, any related files, Ecommerce.db, and pushed accordingly if there are any
changes made.

## 3.2 Github Actions for Continuous Integration

```
82   New records added to the Customer table.
83      customer_id first_name
84   1       c00501      Randy
85   [1] "Customer data is valid. Data loaded into the database..."
86   New records added to the Seller table.
87      seller_id company_name
88   1      s00501 Ruden-Forson
89   [1] "Sellers data is valid. Data loaded into the database..."
90   [1] "No differences found between the old and new data in Category Table. No updates needed.\n"
91   NULL
92   [1] "No differences found between the old and new data in Product Table. No updates needed.\n"
93   NULL
94   [1] "No differences found between the old and new data in Discount Table. No updates needed.\n"
95   NULL
96   [1] "No differences found between the old and new data in Order Table. No updates needed.\n"
97   NULL
98   [1] "No differences found between the old and new data in Shipment Table. No updates needed.\n"
99   NULL

✓  Add files

✓  Check for changes in Ecommerce.db
```

Testing our new datasets involves assuming daily uploads of csv files into the environment. We have developed multiple test datasets, modifying customer and supplier tables while keeping the rest unchanged. The outcomes above exhibit the addition of a new customer (c00501) and seller (s00501). This automated database update confirms the successful functioning of our data pipeline.

| | | | |
|---|---|---|---|
| ✓ **ETL workflow for group 28**<br>ETL workflow for group 28 #37: Scheduled | main | 3 hours ago<br>3m 57s | ... |
| ✓ **ETL workflow for group 28**<br>ETL workflow for group 28 #36: Scheduled | main | 6 hours ago<br>4m 20s | ... |
| ✓ **ETL workflow for group 28**<br>ETL workflow for group 28 #35: Scheduled | main | 9 hours ago<br>3m 11s | ... |
| ✓ **ETL workflow for group 28**<br>ETL workflow for group 28 #34: Scheduled | main | 12 hours ago<br>3m 16s | ... |

As shown, there are continuous and seamless integration of our workflow in intervals of 3 hours. This ensure that our overall database ecosystem stays updated automatically, whenever there are changes make to the datasets or scripts.

30

As for the testing of our new datasets, we make the assumption that there would be daily uploads of csv pushed into environment. We created various test datasets, with changes made to customer and supplier tables, while no changes are made to the rest. We also included data validation of the new dataset to ensure that in the event, if there are any errors, the database will not load the new data accordingly (refer to Appendix 1 for Rscript Codes). The results are reflected above, reflecting new customer c00501 and seller, s00501. Therefore, the database is automatically updated, indicating the success of our data pipeline.

**Challenges:**



Certainly, the team faces insurmountable challenges in the Github phase of the project. There were multiple errors in the initial stage, and it is due to a variety of reasons such as failure to locate file directory, wrong libraries, syntax errors and authentication failure.



The team overcome these barriers by taking a step-by-step approach to resolve the different erros in the different stages, of which, the most problematic issue was authentication as shown above. The team tried setting up a new repository, and the challenge persisted, until we decide to push the files and change the repository token in our local desktop instead of Posit Cloud. As a result, the workflow was successfully run without any errors and commit any changes accordingly, as shown below.

31

# 4. Data Analysis and Reporting with Quarto in R

## 4.1 Advanced Data Analysis in R

1. <u>Top Locations by Purchasing Power</u>

```r
# Query Order and Customer tables, joining them on customer_id
order_customer <- dbGetQuery(database, "
  SELECT O.order_number, O.customer_id, O.product_id, O.quantity, C.customer_country
  FROM `Order` AS O
  JOIN Customer AS C ON O.customer_id = C.customer_id
")

# Join Product table to get product_price
order_customer_product <- dbGetQuery(database, "
  SELECT OC.*, P.price
  FROM (SELECT O.order_number, O.customer_id, O.product_id, O.quantity, C.customer_country
        FROM `Order` AS O
        JOIN Customer AS C ON O.customer_id = C.customer_id) AS OC
  JOIN Product AS P ON OC.product_id = P.product_id
")

# Calculate total sales amount by multiplying quantity and price for each order
order_customer_product <- mutate(order_customer_product, total_sales = quantity * price)
```

```
20  # Group by country and sum the total sales amount
21  country_sales <- order_customer_product %>%
22    group_by(customer_country) %>%
23    summarize(total_sales = sum(total_sales))
24
25  # Sort the countries by total sales amount in descending order
26  country_sales <- arrange(country_sales, desc(total_sales))
27
28  # Visualize top locations by purchasing power (total sales amount)
29  ggplot(country_sales[1:5, ], aes(x = reorder(customer_country, -total_sales), y = total_sales
30    geom_bar(stat = "identity") +
31    labs(x = "Location", y = "Total Sales Amount", title = "Top Locations by Purchasing Power")
32    theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
33    scale_fill_brewer(palette = "Paired")
```



The analysis begins with identifying the top selling locations by purchasing power, revealing which countries generate the highest total sales. This insight is crucial for understanding market dynamics and targeting regions with higher spending capabilities. The total sales is calculated by joining the 'Order' and 'Customer' tables and incorporating product prices for each country, identifying the markets with significant economic impact on the e-commerce operations. The figure shows that China leads in purchasing power, followed by France and Belgium, suggesting a strategic opportunity to further tailor marketing and expansion efforts in these countries.

2. Top 5 Products by Number of Purchases

```
1  # Query Order table and join with Product table to get product names
2  top_products <- dbGetQuery(database, "
3    SELECT P.product_name, COUNT(*) AS purchase_count
4    FROM `Order` AS O
5    JOIN Product AS P ON O.product_id = P.product_id
6    GROUP BY P.product_name
7    ORDER BY purchase_count DESC
8    LIMIT 5
9  ")
10
11
12 # Visualize the top 5 products by purchase count
13 ggplot(top_products, aes(x = reorder(product_name, -purchase_count), y = purchase_count)) +
14   geom_bar(stat = "identity", fill = "skyblue") +
15   labs(x = "Product Name", y = "Purchase Count", title = "Top 5 Products by Purchase Count")
16   theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Utilise the sales volume of each product, the product popularity figure shows the top 5 products: Meditation Pillow Set, Hydrating Facial Moisturizer, Moisture-Wicking Athletic Shirt, Smart Coffee Maker, Wireless Bluetooth Earbuds, offering a clear view of consumer preferences. This analysis is invaluable for inventory management and marketing, highlighting the need to focus on these popular products.

Furthermore, we calculate the sales conversion rate for the top 5 products, which shows the direct correlation between product views and purchase behavior.

3. Top 5 Products by Conversion Rate

```
1   # Query top 5 products by views
2   top_products <- dbGetQuery(database, "
3     SELECT product_id, product_name, product_views
4     FROM Product
5     ORDER BY product_views DESC
6     LIMIT 5
7   ")
8
9   # Query total number of purchases for each of the top 5 products
10  product_purchases <- dbGetQuery(database, "
11    SELECT O.product_id, COUNT(*) AS purchases
12    FROM `Order` AS O
13    WHERE O.product_id IN (SELECT product_id FROM Product ORDER BY product_views DESC LIMIT 5)
14    GROUP BY O.product_id
15  ")
16
17  # Join product views and purchases
18  product_conversion <- left_join(top_products, product_purchases, by = "product_id")
19
20  # Calculate sales conversion rate (purchases / views) and handle NA values
21  product_conversion <- mutate(product_conversion, conversion_rate = ifelse(is.na(purchases) |
22
23  # Remove NA values
24  product_conversion <- na.omit(product_conversion)
25
26  # Visualize the sales conversion rates for top 5 products
27  ggplot(product_conversion, aes(x = reorder(product_name, -conversion_rate), y = conversion_r
28    geom_bar(stat = "identity") +
29    labs(x = "Product Name", y = "Sales Conversion Rate (%)", title = "Sales Conversion Rate by
30    scale_fill_gradient(low = "skyblue", high = "darkblue") +
31    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

## Sales Conversion Rate by Views for Top 5 Products



The figure shows that the sales conversion rate for the 'Moisture-Wicking Athletic Shirt' is notably high, suggesting that this product has an effective market fit and presentation. Emphasizing such products in promotional materials could enhance conversion rates.

Then, an extended analysis on the top 10 products offered deeper insights into how visibility influences sales.

4. Sales Conversion Rate by Purchased Products for Top 10 Products

```
1   # Query top 10 products by views
2   top_products <- dbGetQuery(database, "
3     SELECT product_id, product_name, product_views
4     FROM Product
5     ORDER BY product_views DESC
6     LIMIT 10
7   ")
8
9   # Query total number of purchases for each of the top 10 products
10  product_purchases <- dbGetQuery(database, "
11    SELECT O.product_id, COUNT(*) AS purchases
12    FROM `Order` AS O
13    WHERE O.product_id IN (SELECT product_id FROM Product ORDER BY product_views DESC LIMIT 10)
14    GROUP BY O.product_id
15  ")
16
```

```
17  # Join product views and purchases
18  product_conversion <- left_join(top_products, product_purchases, by = "product_id")
19
20  # Calculate sales conversion rate (purchases / views) and handle NA values
21  product_conversion <- mutate(product_conversion, conversion_rate = ifelse(is.na(purchases) |
22
23  # Remove NA values
24  product_conversion <- na.omit(product_conversion)
25
26  # Visualize the sales conversion rates for top 10 products
27  ggplot(product_conversion, aes(x = reorder(product_name, -conversion_rate), y = conversion_ra
28    geom_bar(stat = "identity") +
29    labs(x = "Product Name", y = "Sales Conversion Rate (%)", title = "Sales Conversion Rate by
30    scale_fill_gradient(low = "skyblue", high = "darkblue") +
31    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



The figure shows that across the top 10 products, the 'Jing Outdoor Lounge Chair' has the highest conversion rate, indicating strong consumer interest and the potential for outdoor furniture as a category for business growth.

5. Average orders and revenue by Country

```
1   # Query the total revenue and count of orders by country
2   country_orders_revenue <- dbGetQuery(database, "
3     SELECT C.customer_country,
4            COUNT(O.order_number) AS order_count,
5            SUM(O.quantity * P.price) AS total_revenue
6     FROM `Order` AS O
7     INNER JOIN Customer AS C ON O.customer_id = C.customer_id
8     INNER JOIN Product AS P ON O.product_id = P.product_id
9     GROUP BY C.customer_country
10  ")
11
12  # Calculate the average order value for each country
13  country_orders_revenue <- mutate(country_orders_revenue, average_order_value = total_revenue
14
15  # Visualize the average order value and total revenue by countries
16  ggplot(country_orders_revenue, aes(x = reorder(customer_country, -total_revenue), y = total_r
17    geom_bar(stat = "identity") +
18    labs(x = "Country", y = "Total Revenue", title = "Total Revenue by Country", fill = "Countr
19    scale_fill_brewer(palette = "Paired") +
20    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Total Revenue by Country

```
1  ggplot(country_orders_revenue, aes(x = reorder(customer_country, -average_order_value), y = a
2    geom_bar(stat = "identity") +
3    labs(x = "Country", y = "Average Order Value", title = "Average Order Value by Country", f:
4    scale_fill_brewer(palette = "Paired") +
5    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Average Order Value by Country

It can be seen in these two figures that China dominates total revenue by country but doesn't lead in average order value, where France takes the lead. Tailoring strategies to increase the average order value in high-revenue countries could balance revenue streams.

6. Rate of returning customers

```
1  # Query the number of orders for each customer
2  customer_order_count <- dbGetQuery(database, "
3    SELECT customer_id, COUNT(DISTINCT order_number) AS order_count
4    FROM `Order`
5    GROUP BY customer_id
6  ")
7
8  # Calculate the number of returning customers
9  returning_customers <- sum(customer_order_count$order_count > 1)
10
11  # Calculate the total number of customers
```

```r
12  total_customers <- nrow(customer_order_count)
13
14  # Calculate the returning customer rate
15  returning_customer_rate <- (returning_customers / total_customers) * 100
16
17  # Create a data frame for visualization
18  data <- data.frame(
19    Customer_Status = c("Returning Customers", "New Customers"),
20    Count = c(returning_customers, total_customers - returning_customers)
21  )
22
23  # Visualize the returning customer rate using a pie chart
24  ggplot(data, aes(x = "", y = Count, fill = Customer_Status)) +
25    geom_bar(stat = "identity", width = 1) +
26    coord_polar(theta = "y") +
27    labs(fill = "Customer Status", title = "Returning Customer Rate") +
28    theme_void() +
29    theme(legend.position = "bottom")
```

## Returning Customer Rate



Customer Status  ▮ New Customers  ▮ Returning Customers

In order to understand customer loyalty, the rate of returning customers was calculated, revealing the proportion of the customer base that makes repeated purchases. This metric is essential for evaluating the success of retention strategies and the overall satisfaction of the

customer base. A high rate of returning customers is indicative of a healthy e-commerce ecosystem with strong customer loyalty.

The significant segment of returning customers in this figure implies a strong base of customer loyalty. Retention strategies should continue to be prioritized to maintain this valuable consumer segment.

## 7. Demographics and Platform Analysis

```r
# Query the platforms used by customers
customer_platforms <- dbGetQuery(database, "
  SELECT platform, COUNT(*) AS customer_count
  FROM Customer
  WHERE platform IS NOT NULL
  GROUP BY platform
")

# Visualize the analysis
ggplot(customer_platforms, aes(x = platform, y = customer_count)) +
  geom_segment(aes(xend = platform, yend = 0), color = "skyblue") +  # Lollipop stems
  geom_point(color = "blue", size = 3) +  # Lollipop heads
  labs(x = "Platform", y = "Number of Customers", title = "Number of Customers by Platform")
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  theme_classic()
```



Number of Customers by Platform

The number of customers by platform reveals the most popular channels among the customer base, guiding efforts to optimise presence across platforms. This analysis aids in understanding where to focus digital marketing efforts for maximum engagement and conversion. Shown in this figure, the distribution of customers across four different platforms is relatively even, with 'Facebook' having a slight lead. This suggests an opportunity to optimise and diversify platform-specific marketing strategies for better engagement.

8. The effectiveness of discounts on product sales.

```
# Query discount effectiveness
discount_analysis <- dbGetQuery(database, "
  SELECT D.discount_percentage, COUNT(*) AS order_count
  FROM Discount AS D
  INNER JOIN `Order` AS O ON D.product_id = O.product_id
  GROUP BY D.discount_percentage
")

# Visualize the results
ggplot(discount_analysis, aes(x = discount_percentage, y = order_count)) +
  geom_line(group=1, color = "blue") +
  geom_point(color = "blue") +
  labs(x = "Discount Percentage", y = "Number of Orders", title = "Effectiveness of Discount
  theme_classic()
```

Lastly, the effectiveness of discounts on each product sales is analysed, uncovering how different discount percentages influence the number of orders. This analysis provides a deep understanding of pricing strategies' impact on sales volume, essential for optimising pricing for increased sales and customer acquisition.

Discount effectiveness seems to peak at a moderate discount percentage, with diminishing returns as the discount increases. This suggests an optimal discount rate can be found that balances attractiveness to customers with profitability.

9. Number of orders refunds by month

```r
# Query refunded orders with order dates and extract month directly in SQL
refund_by_month <- dbGetQuery(database, "
  SELECT strftime('%Y-%m', O.order_date) AS month,
         COUNT(*) AS refund_count
  FROM `Order` AS O
  INNER JOIN Shipment AS S ON O.order_number = S.order_number
  WHERE S.refund = 'Yes'
  GROUP BY month
")

# Convert month to Date type
refund_by_month$month <- ymd(paste(refund_by_month$month, "-01", sep = "-"))

# Convert month labels to abbreviated month names
refund_by_month$month <- format(refund_by_month$month, "%b")

# Set the order of months
months_order <- c("Jan", "Feb", "Mar", "Apr", "May", "Jun", "Jul", "Aug", "Sep", "Oct", "Nov"
refund_by_month$month <- factor(refund_by_month$month, levels = months_order)

# Visualize the number of refunds by month
ggplot(refund_by_month, aes(x = month, y = refund_count)) +
  geom_bar(stat = "identity", fill = "skyblue") +
  labs(x = "Month", y = "Number of Orders Refunds", title = "Refunds by Month") +
  theme_classic()
```

The bar graph displays monthly order refunds, highlighting noticeable trends. February and December exhibit the highest refund counts at four, while March, June, October, and November have the lowest, with only one refund each.

**Challenges:**

*Choosing the right data*: Given the large number of data encompassing customer behavior, order, and product details, it becomes a challenging endeavor to distill information that is not only aligned with our analytical objectives but is also precise and readily comprehensible.

*Incorrect Join Operations*: Using joins to link different tables is quite tricky. It's crucial for us to determine which key to use for the connection. For example, if the 'JOIN' conditions between Order and Shipment tables are specified incorrectly, it could lead to an inaccurate dataset, such as joining on the wrong key or missing the correct join logic.

*Visualization Complexity*: Designing clear and informative visualizations to present the data effectively can be challenging. Selecting appropriate visualization techniques, formatting charts and graphs, and ensuring visual clarity while conveying complex information require careful attention to detail.

## 5. Conclusion

In summary, this report delves deeply into the intricacies of managing and analyzing e-commerce data. It starts by laying the groundwork with an Entity-Relationship Diagram (ERD), providing a clear blueprint of how various elements like customers, orders, products,

and discounts are interrelated within the database structure. Through meticulous validation processes, the integrity and quality of data are ensured before integration into the database.

Utilizing tools like Python for data generation and R for advanced analysis, the report uncovers valuable insights into customer behavior, product popularity, revenue distribution by country, and the impact of discounts on sales. These insights are crucial for strategic decision-making in areas such as inventory management, marketing optimization, and customer retention strategies. However, the report also acknowledges the challenges inherent in managing e-commerce data, such as database complexity, uncertainty in identifying relationship between elements and lack of expertise issues. Despite these challenges, the meticulous validation and insightful analysis provided in this report offer a solid foundation for informed decision-making and strategic planning in the ever-evolving landscape of digital commerce.

# 6. Appendix

## Appendix 1

```r
## Customer

compare_and_update_database_cust <- function(old_csv, new_csv, table_name, primary_key) {
  # Load old and new data
  old_data <- read.csv(old_csv)
  new_data <- read.csv(new_csv)

  # Check for differences
  added_rows <- anti_join(new_data, old_data, by = primary_key)

  if (nrow(added_rows) == 0) {
    print("No differences found between the old and new data in Customer Table. No updates ne
    return(NULL)
  }
  # Function to check if a datetime is in the desired format ("%Y-%m-%d %H:%M:%S")
  is_datetime_format <- function(datetime_string) {
    tryCatch({
      as.POSIXlt(datetime_string, format = "%Y-%m-%d %H:%M:%S")
      TRUE
    }, error = function(e) {
      FALSE
    })
  }
```

```r
     # Check if dates are in the desired format, if not, convert them
     for (i in 1:nrow(Customer)) {
       if (!is_datetime_format(Customer$date_of_birth[i])) {
         Customer$date_of_birth[i] <- as.POSIXct(Customer$date_of_birth[i], format = "%Y-%m-%d :
       }
     }

     # Data validation for new customer data
     validate_customer_data <- function(data) {
       # Check unique customer IDs
       if (length(unique(data$customer_id)) != nrow(data)) {
         stop("Customer ID is not unique.")
       }

       # Check data types for first_name and last_name
       if (!all(sapply(data$first_name, is.character)) || !all(sapply(data$last_name, is.charac
         stop("First name and last name should be character.")
       }

       # Check valid gender values
       valid_genders <- c("Male", "Female", "Other")
       if (any(!data$gender %in% valid_genders)) {
         stop("Gender should be Male, Female, or Other.")
       }

       # Check email format
       if (any(!grepl("^\\S+@\\S+\\.\\S+$", data$email))) {
         stop("Invalid email format")
       }

       # Regular expressions for phone number formats of Belgium, China, France, United Kingdom
       phone_regex <- "^\\(\\+\\d+\\)\\d{10}$"

       # Check phone number format for specific countries
       invalid_phone_indices <- which(!grepl(phone_regex, data$phone))
       if (length(invalid_phone_indices) > 0) {
         stop("Invalid phone numbers.")
       }

       # Regular expressions for zip code formats of Belgium, China, France, United Kingdom, Un
       zip_regex <- c(
         "^[0-9]{4}$",  # Belgium
```

```r
        "^[0-9]{6}$",  # China
        "^[0-9]{5}$",  # France
        "^[A-Z]{2}[0-9]{1,2}[A-Z]? [0-9][A-Z]{2}$",  # United Kingdom
        "^[0-9]{5}-[0-9]{4}$"    # United States
    )

    # Check zip code format for specific countries
    invalid_zip_indices <- which(!grepl(paste(zip_regex, collapse = "|"), data$customer_zip_
    if (length(invalid_zip_indices) > 0) {
      stop("Invalid zip codes.")
    }

    # Check platform values
    valid_platforms <- c("Referral", "Instagram", "Facebook", "Others")
    if (any(!data$platform %in% valid_platforms)) {
      stop("Invalid platform values.")
    }

    return(TRUE)
  }

  # Validate new customer data
  if (!validate_customer_data(added_rows)) {
    return(NULL)
  }

  # Create a database connection
  database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')

  # Drop existing Customer table if it exists
  if (dbExistsTable(database, table_name)) {
    dbExecute(database, paste0("DROP TABLE ", table_name))
  }

  # Write the new data to the database with append
  RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
  cat("New records added to the Customer table.\n")

  # Print information about the last added record
  if (nrow(added_rows) > 0) {
    cust_result <- RSQLite::dbGetQuery(database, "SELECT * FROM Customer ORDER BY ROWID DESC
    print(cust_result[c("customer_id", "first_name")])
```

```
109      print("Customer data is valid. Data loaded into the database...")
110    }
111
112    # Close database connection
113    dbDisconnect(database)
114 }
115
116 # Usage example
117 compare_and_update_database_cust("Dataset/customer.csv", "DatasetTest/customer_data_test_new
118
119 ## Seller
120
121 compare_and_update_database_seller <- function(old_csv, new_csv, table_name, primary_key) {
122    # Load old and new data
123    old_data <- read.csv(old_csv)
124    new_data <- read.csv(new_csv)
125
126    # Check for differences
127    added_rows <- anti_join(new_data, old_data, by = primary_key)
128
129    if (nrow(added_rows) == 0) {
130      print("No differences found between the old and new data in Seller Table. No updates nee
131      return(NULL)
132    }
133
134    # Create a database connection
135    database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
136
137    # Drop existing Seller table if it exists
138    if (dbExistsTable(database, table_name)) {
139      dbExecute(database, paste0("DROP TABLE ", table_name))
140    }
141
142    # Write the new data to the database with append
143    RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
144    cat("New records added to the Seller table.\n")
145
146    # Print information about the last added record
147    if (nrow(added_rows) > 0) {
148      seller_result <- RSQLite::dbGetQuery(database, paste0("SELECT * FROM ", table_name, " ORI
149      print(seller_result[c("seller_id", "company_name")])
150    }
```

```
151
152     ## Validation of Seller Data
153     library(stringr)
154     na_sellers <- apply(is.na(added_rows), 2, sum)
155
156     # Ensure "seller_Id" values are unique
157     if (length(unique(added_rows$seller_id)) != nrow(added_rows)) {
158       print("seller_Id values are not unique.")
159     }
160
161     # Check length of "company_name"
162     if (any(nchar(added_rows$company_name) > 100)) {
163       print("company_name exceeds 100 characters.")
164     }
165
166     # Check email format
167     invalid_emails <- which(!str_detect(added_rows$supplier_email, "\\b[A-Za-z0-9._%+-]+@[A-Za-
168     if (length(invalid_emails) > 0) {
169       print("Invalid email addresses:")
170       print(added_rows[invalid_emails, ])
171     }
172
173     # If no errors are found, print a message indicating that the data is valid
174     if (!any(is.na(na_sellers)) &&
175         length(unique(added_rows$seller_id)) == nrow(added_rows) &&
176         !any(nchar(added_rows$company_name) > 100) &&
177         length(invalid_emails) == 0) {
178       print("Sellers data is valid. Data loaded into the database...")
179       RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE
180       # Load the data into the database
181     } else {
182       print("Sellers data is not valid. Please correct the errors.")
183     }
184
185     # Close database connection
186     RSQLite::dbDisconnect(database)
187 }
188
189 # Usage example
190 compare_and_update_database_seller("Dataset/seller.csv", "DatasetTest/seller_data_test_new_r
191
192 ## Category
```

```r
compare_and_update_database_category <- function(old_csv, new_csv, table_name, primary_key)
    # Load old and new data
    old_data <- read.csv(old_csv)
    new_data <- read.csv(new_csv)

    # Check for differences
    added_rows <- anti_join(new_data, old_data, by = primary_key)

    if (nrow(added_rows) == 0) {
      print(paste("No differences found between the old and new data in", table_name, "Table. ]
      return(NULL)
    }

    # Validation for category data
    na_prod_cat <- apply(is.na(added_rows), 2, sum)

    # Ensure "category_id" values are unique
    if (length(unique(added_rows$category_id)) != nrow(added_rows)) {
      print("category_id values are not unique.")
    }

    # Check length of "cat_name"
    if (any(nchar(added_rows$cat_name) > 255)) {
      print("cat_name exceeds 255 characters.")
    }

    # Check data type of each column
    if (!all(sapply(added_rows$category_id, is.character)) ||
        !all(sapply(added_rows$cat_name, is.character)) ||
        !all(sapply(added_rows$cat_description, is.character))) {
      print("Invalid data type for one or more columns.")
    }

    # If any errors are found, do not proceed with writing to the database
    if (any(is.na(na_prod_cat)) ||
        length(unique(added_rows$category_id)) != nrow(added_rows) ||
        any(nchar(added_rows$cat_name) > 255) ||
        !all(sapply(added_rows$category_id, is.character)) ||
        !all(sapply(added_rows$cat_name, is.character)) ||
        !all(sapply(added_rows$cat_description, is.character))) {
      print(paste(table_name, "data is not valid. Please correct the errors."))
```

```r
235        return(NULL)
236      }
237
238      # Create a database connection
239      database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
240
241      # Drop existing Category table if it exists
242      if (dbExistsTable(database, table_name)) {
243        dbExecute(database, paste0("DROP TABLE ", table_name))
244      }
245
246      # Write the new data to the database with append
247      RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
248      cat(paste("New records added to the", table_name, "table.\n"))
249
250      # Print information about the last added record
251      if (nrow(added_rows) > 0) {
252        category_result <- RSQLite::dbGetQuery(database, paste0("SELECT * FROM ", table_name, " (
253        print(category_result[c("category_id", "cat_name")])
254        print("Category data is valid. Loaded data into the database...")
255      }
256
257      # Close database connection
258      dbDisconnect(database)
259    }
260
261    compare_and_update_database_category("Dataset/category.csv", "DatasetTest/category_data_no_n
262
263    ## Product
264
265    compare_and_update_database_prod <- function(old_csv, new_csv, table_name, primary_key) {
266      # Load old and new data
267      old_data <- read.csv(old_csv)
268      new_data <- read.csv(new_csv)
269
270      # Check for differences
271      added_rows <- anti_join(new_data, old_data, by = primary_key)
272
273      if (nrow(added_rows) == 0) {
274        print("No differences found between the old and new data in Product Table. No updates ne
275        return(NULL)
276      }
```

```r
    # Function to validate product data
    validate_product_data <- function(data) {
      # Function to check if a value is decimal
      valid_decimal <- function(x) {
        !is.na(as.numeric(x))
      }

      # Function to check if a value is an integer
      valid_integer <- function(x) {
        !is.na(as.integer(x))
      }

      # Check for missing values
      na_values <- apply(is.na(data), 2, sum)

      # Ensure "product_id" values are unique
      if (length(unique(data$product_id)) != nrow(data)) {
        stop("product_id values are not unique.")
      }

      # Check length of "product_name"
      if (any(nchar(data$product_name) > 255)) {
        stop("product_name exceeds 255 characters.")
      }

      # Check if inventory and product views are integers
      if (any(!sapply(data$inventory, valid_integer)) || any(!sapply(data$product_views, valid
        stop("Inventory and product views should be integers.")
      }

      # Check if price and weight are decimal
      if (any(!sapply(data$price, valid_decimal)) || any(!sapply(data$weight, valid_decimal)))
        stop("Price and weight should be decimal values.")
      }

      # If no errors are found, return TRUE
      return(TRUE)
    }

    # Validate new product data
    if (!validate_product_data(added_rows)) {
```

```r
319        return(NULL)
320      }
321
322      # Create a database connection
323      database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
324
325      # Drop existing Product table if it exists
326      if (dbExistsTable(database, table_name)) {
327        dbExecute(database, paste0("DROP TABLE ", table_name))
328      }
329
330      # Write the new data to the database with append
331      RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
332      cat("New records added to the Product table.\n")
333
334      # Print information about the last added record
335      if (nrow(added_rows) > 0) {
336        prod_result <- RSQLite::dbGetQuery(database, "SELECT * FROM Product ORDER BY ROWID DESC 
337        print(prod_result[c("product_id", "product_name")])
338        print("Product data is valid. Data loaded into the database...")
339      }
340
341      # Close database connection
342      dbDisconnect(database)
343    }
344    # Usage example
345    compare_and_update_database_prod("Dataset/product.csv", "DatasetTest/product_data_test_no_ne
346
347    ## Discount
348
349    compare_and_update_database_discount <- function(old_csv, new_csv, table_name, primary_key) 
350      # Load old and new data
351      old_data <- read.csv(old_csv)
352      new_data <- read.csv(new_csv)
353
354      # Check for differences
355      added_rows <- anti_join(new_data, old_data, by = primary_key)
356
357      if (nrow(added_rows) == 0) {
358        print("No differences found between the old and new data in Discount Table. No updates ne
359        return(NULL)
360      }
```

```
361
362     # Validation for discount data
363     ## Function to check if date is in the desired format
364     is_datetime_format <- function(x) {
365       tryCatch({
366         as.POSIXlt(x, format = "%Y-%m-%d %H:%M:%S")
367         TRUE
368       }, error = function(e) {
369         FALSE
370       })
371     }
372
373     ## Convert discount_start_date and discount_end_date to desired format if not already in t
374     if (!all(sapply(added_rows$discount_start_date, is_datetime_format))) {
375       added_rows$discount_start_date <- as.POSIXlt(added_rows$discount_start_date, format = "%
376     }
377
378     if (!all(sapply(added_rows$discount_end_date, is_datetime_format))) {
379       added_rows$discount_end_date <- as.POSIXlt(added_rows$discount_end_date, format = "%Y-%m-
380     }
381
382     ## Check for missing values in Discount dataframe
383     na_disc <- apply(is.na(added_rows), 2, sum)
384
385     ## Validate discount_percentage, discount_start_date, and discount_end_date data types
386     valid_decimal <- function(x) {
387       !is.na(as.numeric(x))
388     }
389
390     valid_datetime <- function(x) {
391       !is.na(as.POSIXlt(x))
392     }
393
394     ## Check discount percentage range (assuming it's between 0 and 100)
395     if (any(added_rows$discount_percentage < 0 | added_rows$discount_percentage > 100) ||
396         !all(sapply(added_rows$discount_percentage, valid_decimal))) {
397       print("Invalid discount percentage.")
398     }
399
400     ## Check discount dates
401     if (any(added_rows$discount_start_date >= added_rows$discount_end_date) ||
402         !all(sapply(added_rows$discount_start_date, valid_datetime)) ||
```

```
403        !all(sapply(added_rows$discount_end_date, valid_datetime))) {
404      print("Discount start date should be before the end date.")
405    }
406
407    ## Check if discount_id is unique
408    if (any(duplicated(added_rows$discount_id))) {
409      print("Duplicate discount IDs found.")
410    }
411
412    ## Check if product_id exists in Product table
413    if (any(!added_rows$product_id %in% Product$product_id)) {
414      print("Invalid product IDs. Some product IDs do not exist in the Product table.")
415    }
416
417    ## If no errors are found, print a message indicating that the data is valid
418    if (!any(is.na(na_disc)) &&
419        all(added_rows$discount_percentage >= 0 & added_rows$discount_percentage <= 100) &&
420        all(added_rows$discount_start_date < added_rows$discount_end_date) &&
421        !any(duplicated(added_rows$discount_id)) &&
422        all(added_rows$product_id %in% Product$product_id) &&
423        all(sapply(added_rows$discount_percentage, valid_decimal)) &&
424        all(sapply(added_rows$discount_start_date, valid_datetime)) &&
425        all(sapply(added_rows$discount_end_date, valid_datetime))) {
426      print("Discount data is valid. Loaded data into the database...")
427
428      database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
429
430      # Drop existing Discount table if it exists
431      if (dbExistsTable(database, table_name)) {
432        dbExecute(database, paste0("DROP TABLE ", table_name))
433      }
434
435      RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
436      cat("New records added to the Discount table.\n")
437
438      discount_result <- RSQLite::dbGetQuery(database, paste0("SELECT * FROM ", table_name, " (
439      print(discount_result[c("discount_id", "discount_percentage")])
440      dbDisconnect(database)
441    } else {
442      print("Data is not valid. Please correct the errors.")
443    }
444  }
```

```r
445
446   # Usage example
447   compare_and_update_database_discount("Dataset/discount.csv", "DatasetTest/discount_data_no_ne
448
449   ## Order
450
451   compare_and_update_database_order <- function(old_csv, new_csv, table_name, primary_key) {
452     # Load old and new data
453     old_data <- read.csv(old_csv)
454     new_data <- read.csv(new_csv)
455
456     # Check for differences
457     added_rows <- anti_join(new_data, old_data, by = primary_key)
458
459     if (nrow(added_rows) == 0) {
460       print("No differences found between the old and new data in Order Table. No updates neede
461       return(NULL)
462     }
463
464     # Validation for order data
465     na_order <- apply(is.na(added_rows), 2, sum)
466
467     # Check quantity (assuming it should be a positive integer)
468     if (any(added_rows$quantity <= 0)) {
469       print("Invalid quantity.")
470     }
471
472     # Check customer rating (assuming it should be between 1 and 5)
473     if (any(added_rows$customer_rating < 1 | added_rows$customer_rating > 5)) {
474       print("Invalid customer rating.")
475     }
476
477     # Check if product_id exists in Product table
478     if (any(!added_rows$product_id %in% Product$product_id)) {
479       print("Invalid product IDs. Some product IDs do not exist in the Product table.")
480     }
481
482     # Check if customer_id exists in Customer table
483     if (any(!added_rows$customer_id %in% Customer$customer_id)) {
484       print("Invalid customer IDs. Some customer IDs do not exist in the Customer table.")
485     }
486
```

```r
487    # Check if shipment_id exists in Shipment table
488    if (any(!added_rows$shipment_id %in% Shipment$shipment_id)) {
489      print("Invalid shipment IDs. Some shipment IDs do not exist in the Shipment table.")
490    }
491
492    # Check uniqueness based on primary key (order_number, customer_id, product_id)
493    if (any(duplicated(added_rows[c("order_number", "customer_id", "product_id")]))) {
494      print("Duplicate records found based on order_number, customer_id, and product_id.")
495    }
496
497    # Check order date format and range
498    if (any(!is_datetime_format(Order$order_date))) {
499      # Convert order date to the desired format if not already
500      Order$order_date <- as.POSIXct(Order$order_date, format = "%Y-%m-%d %H:%M:%S", tz = "UTC"
501    }
502
503    # If any errors are found, do not proceed with writing to the database
504    if (any(is.na(na_order)) ||
505        any(added_rows$quantity <= 0) ||
506        any(added_rows$customer_rating < 1 | added_rows$customer_rating > 5) ||
507        any(!added_rows$product_id %in% Product$product_id) ||
508        any(!added_rows$customer_id %in% Customer$customer_id) ||
509        any(!added_rows$shipment_id %in% Shipment$shipment_id) ||
510        any(duplicated(added_rows[c("order_number", "customer_id", "product_id")])) &&
511        all(is_datetime_format(Order$order_date))) {
512      print("Order data is not valid. Please correct the errors.")
513      return(NULL)
514    }
515
516    # Create a database connection
517    database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
518
519    # Drop existing Order table if it exists
520    if (dbExistsTable(con, table_name)) {
521      dbExecute(con, paste0("DROP TABLE ", table_name))
522    }
523
524    # Write the new data to the database with append
525    RSQLite::dbWriteTable(con, table_name, added_rows, append = TRUE, row.names = FALSE)
526    cat("New records added to the Order table.\n")
527
528    # Print information about the last added record
```

57

```r
529    if (nrow(added_rows) > 0) {
530      order_result <- RSQLite::dbGetQuery(database, paste0("SELECT * FROM ", table_name, " ORDI
531      print(order_result[c("order_number", "product_id")])
532      print("Order data is valid. Loaded data into the database...")
533    }
534
535    # Close database connection
536    dbDisconnect(database)
537  }
538
539  # Usage example
540  compare_and_update_database_order("Dataset/order.csv", "DatasetTest/order_data_test_no_new.c
541
542
543  ## Shipment Data
544
545  compare_and_update_database_shipment <- function(old_csv, new_csv, table_name, primary_key)
546    # Load old and new data
547    old_data <- read.csv(old_csv)
548    new_data <- read.csv(new_csv)
549
550    # Check for differences
551    added_rows <- anti_join(new_data, old_data, by = primary_key)
552
553    if (nrow(added_rows) == 0) {
554      print(paste("No differences found between the old and new data in", table_name, "Table. 
555      return(NULL)
556    }
557
558    # Validation for Shipment Data
559    na_shipment <- sapply(added_rows, function(x) sum(is.na(x)))
560
561    # Ensure "shipment_id" values are unique
562    if (length(unique(added_rows$shipment_id)) != nrow(added_rows)) {
563      print("shipment_id values are not unique.")
564    }
565
566    # Validate "refund" column
567    valid_refunds <- c("Yes", "No")
568    if (!all(added_rows$refund %in% valid_refunds)) {
569      print("Invalid values in the 'refund' column.")
570    }
```

58

```
571
572     # Validate "shipment_delay_days" and "shipment_cost" columns
573     if (any(added_rows$shipment_delay_days <= 0) || any(added_rows$shipment_cost <= 0)) {
574       print("shipment_delay_days and shipment_cost should be positive numbers.")
575     }
576
577     # Ensure that "shipment_delay_days" is an integer
578     if (any(!as.integer(added_rows$shipment_delay_days) == added_rows$shipment_delay_days)) {
579       print("shipment_delay_days should be integers.")
580     }
581
582     # Ensure that all "order_number" values exist in the "Order" table
583     order_numbers <- unique(added_rows$order_number)
584     if (!all(order_numbers %in% Order$order_number)) {
585       print("Some order numbers do not exist in the 'Order' table.")
586     }
587
588     # If any errors are found, do not proceed with writing to the database
589     if (any(na_shipment != 0) ||
590         length(unique(added_rows$shipment_id)) != nrow(added_rows) ||
591         !all(added_rows$refund %in% valid_refunds) ||
592         any(added_rows$shipment_delay_days <= 0) || any(added_rows$shipment_cost <= 0) ||
593         any(!as.integer(added_rows$shipment_delay_days) == added_rows$shipment_delay_days) ||
594         !all(order_numbers %in% Order$order_number)) {
595       print(paste(table_name, "data is not valid. Please correct the errors."))
596       return(NULL)
597     }
598
599     # Create a database connection
600     database <- RSQLite::dbConnect(RSQLite::SQLite(), dbname = 'Ecommerce.db')
601
602     # Drop existing Shipment table if it exists
603     if (dbExistsTable(database, table_name)) {
604       dbExecute(database, paste0("DROP TABLE ", table_name))
605       cat(paste("Existing", table_name, "table dropped.\n"))
606     }
607
608     # Write the new data to the database with append
609     RSQLite::dbWriteTable(database, table_name, added_rows, append = TRUE, row.names = FALSE)
610     cat(paste("New records added to the", table_name, "table.\n"))
611
612     # Print information about the last added record
```

```
613    if (nrow(added_rows) > 0) {
614      shipment_result <- RSQLite::dbGetQuery(database, paste0("SELECT * FROM ", table_name, " (
615      print(shipment_result[c("shipment_id", "order_number")])
616    }
617
618    # Close database connection
619    dbDisconnect(database)
620  }
621
622  # Usage example
623  compare_and_update_database_shipment("Dataset/shipment.csv", "DatasetTest/shipment_data_no_n(
```

## Appendix 2

```
1   knitr::opts_chunk$set(python.reticulate = FALSE)
2   # install faker library
3   # pip install Faker
4
5   # csv library
6   import csv
7
8   # random library
9   import random
10
11  # datetime
12  from datetime import datetime, timedelta
13
14  # inititalise faker generator
15  from faker import Faker
16  fake = Faker()
17
18
19  # %% customer
20  country_codes = ['+44','+1', '+32', '+33', '+86'] #usa/canada belgium france china
21
22  # Function to get city information based on country code
23  def get_country_info_cust(country_code):
24      if country_code == '+44':
25          return 'United Kingdom'
26      elif country_code == '+1':
```

```
27            return 'United States'
28        elif country_code == '+32':
29            return 'Belgium'
30        elif country_code == '+33':
31            return 'France'
32        elif country_code == '+86':
33            return 'China'
34        else:
35            return 'Unknown'

37  def customer(filename, num_customers=500):
38      fieldnames = ['customer_id', 'first_name', 'last_name', 'gender', 'date_of_birth', 'emai

40      with open(filename, 'w', newline='') as csvfile:
41          csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)

43          # Write the header
44          csvwriter.writeheader()

46          # Generate and write fake data to the CSV file
47          for customer_id_num in range(1, num_customers + 1):
48              customer_id = f'c{customer_id_num:05}'  # Format customer_id as 'c' followed by !
49              gender = fake.random_element(elements=('Male', 'Female', 'Other'))
50              first_name = fake.first_name_male() if gender == 'Male' else fake.first_name_fema
51              last_name = fake.last_name()
52              email = f"{first_name.lower()}_{last_name.lower()}@gmail.com"
53              country_code = fake.random_element(elements=country_codes)
54              phone = f"({country_code}){fake.random_number(digits=10, fix_len=True)}"
55              customer_street = fake.street_address()
56              customer_country = get_country_info_cust(country_code)
57              customer_zip_code = fake.zipcode()
58              date_of_birth = fake.date_of_birth(minimum_age=35, maximum_age=60).strftime('%d/%
59              platform = fake.random_element(elements=('Facebook', 'Instagram', 'Referral', 'Ot

61              csvwriter.writerow({
62                  'customer_id': customer_id,
63                  'first_name': first_name,
64                  'last_name': last_name,
65                  'gender': gender,
66                  'date_of_birth': date_of_birth,
67                  'email': email,
68                  'phone': phone,
```

```python
69                  'customer_street': customer_street,
70                  'customer_country': customer_country,
71                  'customer_zip_code': customer_zip_code,
72                  'platform': platform
73              })
74
75  # Call the function to generate and save fake customer data to a CSV file
76  customer('customer.csv', num_customers=500)
77
78  # %% seller entity
79
80
81  country_codes = ['+44', '+1', '+32', '+33', '+86']  # United Kingdom, USA/Canada, Belgium, Fi
82
83  # Function to get city information based on country code
84  def get_country_info_seller(country_code):
85      if country_code == '+44':
86          return 'United Kingdom'
87      elif country_code == '+1':
88          return 'United States'  # You can add more cities for USA/Canada
89      elif country_code == '+32':
90          return 'Belgium'
91      elif country_code == '+33':
92          return 'France'
93      elif country_code == '+86':
94          return 'China'
95      else:
96          return 'Unknown'
97
98  # Function to generate fake data and save it to a CSV file
99  def seller(filename, num_records=500):
100     fieldnames = ['seller_id', 'company_name', 'supplier_phone', 'supplier_email', 'seller_s
101
102     with open(filename, 'w', newline='') as csvfile:
103         csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)
104
105         # Write the header
106         csvwriter.writeheader()
107
108         # Generate and write fake data to the CSV file
109         for seller_id_num in range(1, num_records + 1):
110             seller_id = f's{seller_id_num:05}'
```

```
111          company_name = fake.company().replace(',', '-High')
112          country_code = random.choice(country_codes)
113          supplier_phone = f"({country_code}){fake.random_number(digits=10, fix_len=True)}"
114          supplier_email = f"{fake.first_name().lower()}@{company_name.split()[0].lower()}
115          seller_street = fake.street_address()
116          seller_country = get_country_info_seller(country_code)
117          seller_zip_code = fake.zipcode()
118
119          csvwriter.writerow({
120              'seller_id': seller_id,
121              'company_name': company_name,
122              'supplier_phone': supplier_phone,
123              'supplier_email': supplier_email,
124              'seller_street': seller_street,
125              'seller_country': seller_country,
126              'seller_zip_code': seller_zip_code
127          })
128
129  # Call the function to generate and save fake data to a CSV file (500 records)
130  seller('seller.csv', num_records=500)
131
132
133  # %% product category
134
135  category_descriptions = {
136      "Electronics": "Explore the latest in cutting-edge technology with our electronic gadgets
137      "Home and Kitchen": "Enhance your living spaces with our stylish and functional home and
138      "Sports and Outdoors": "Gear up for outdoor adventures with our high-quality sports and
139      "Clothing and Accessories": "Stay on trend with our fashionable clothing and accessories
140      "Beauty and Personal Care": "Discover a world of beauty and personal care products to enh
141      "Health and Wellness": "Prioritize your well-being with our selection of health and well
142      "Toys and Games": "Entertain and educate with our fun and exciting toys and games for al
143      "Automotive": "Keep your vehicle running smoothly with our automotive parts and accessor
144      "Books and Literature": "Immerse yourself in captivating stories and knowledge with our
145      "Garden and Outdoor": "Create a lush and inviting outdoor space with our gardening and o
146      "Sportswear": "Elevate your active lifestyle with our stylish and high-performance sport
147      "Jewelry": "Adorn yourself with exquisite jewelry that complements your unique style.",
148      "Skincare": "Indulge in luxurious skincare products designed to nourish and rejuvenate y
149      "Health Supplements": "Boost your health and vitality with our premium selection of heal
150      "Board Games": "Gather friends and family for memorable game nights with our exciting bo
151      "Car Engine Products": "Enhance the performance and longevity of your vehicle with our h
152      "Gardening Tools": "Transform your garden into a vibrant oasis with our premium Gardening
```

```python
153  }
154
155
156  # Function to generate categories and save them to a CSV file
157  def category(filename, num_categories=17):
158      fieldnames = ['category_id','p_category_id', 'cat_name', 'cat_description']
159
160      with open(filename, 'w', newline='') as csvfile:
161          csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)
162
163          # Write the header
164          csvwriter.writeheader()
165
166          # Mapping of cat_name to p_category_id
167          p_category_mapping = {
168              "Sportswear": "pc01",
169              "Jewelry": "pc02",
170              "Skincare": "pc03",
171              "Health Supplements": "pc04",
172              "Board Games": "pc05",
173              "Car Engine Products": "pc06",
174              "Gardening Tools": "pc07"
175          }
176
177          # Generate and write category data to the CSV file
178          for category_id_num, (cat_name, cat_description) in enumerate(category_descriptions.
179              category_id = f'c{category_id_num:02}'  # Format category_id as 'c' followed by
180
181              # Get the corresponding p_category_id based on cat_name
182              p_category_id = p_category_mapping.get(cat_name, None)
183
184              # Set p_category_id to 'NULL' if not found in the mapping
185              p_category_id = 'NULL' if p_category_id is None else p_category_id
186
187              csvwriter.writerow({
188                  'category_id': category_id,
189                  'p_category_id': p_category_id,
190                  'cat_name': cat_name,
191                  'cat_description': cat_description
192              })
193
194  # Call the function to generate and save category data to a CSV file
```

```python
category('category.csv', num_categories=17)


# %% products

product_list = {
    "Electronics": [
        {"name": "SmartHome Hub", "description": "A central control hub for all your smart h
        {"name": "Wireless Bluetooth Earbuds", "description": "Enjoy wireless freedom with h
        {"name": "4K Ultra HD Smart TV", "description": "Bring the cinema experience home wi
    ],
    "Home and Kitchen": [
        {"name": "Smart Coffee Maker", "description": "Brew your favorite coffee with smart
        {"name": "Non-Stick Cookware Set", "description": "Premium cookware for your kitchen
        {"name": "Smart Mini Refrigerator", "description": "Keep your food fresh and organize
    ],
    "Sports and Outdoors": [
        {"name": "Fitness Tracker", "description": "Track your fitness activities with this
        {"name": "Camping Tent", "description": "Explore the outdoors with a durable camping
        {"name": "Waterproof Hiking Boots", "description": "Conquer any trail with our Water
    ],

    "Clothing and Accessories": [
        {"name": "Classic Leather Jacket", "description": "Make a statement with our Classic
        {"name": "Stylish Sunglasses", "description": "Shield your eyes in style with our St
        {"name": "Cozy Knit Sweater", "description": "Embrace warmth and comfort with our Co
    ],
    "Beauty and Personal Care": [
        {"name": "Luxury Concealer", "description": "Indulge in the ultimate concealer exper
        {"name": "Professional Hair Dryer", "description": "Achieve salon-quality results at
        {"name": "Relaxing Aromatherapy Candle", "description": "Unwind and de-stress with o
    ],
    "Health and Wellness": [
        {"name": "Fitness Tracker Watch", "description": "Take charge of your health journey
        {"name": "Organic Superfood Blend", "description": "Boost your nutrition with our Or
        {"name": "Meditation Pillow Set", "description": "Find tranquility and peace with ou
    ],
    "Toys and Games": [
        {"name": "Interactive Robot Toy", "description": "Spark creativity and play with our
        {"name": "Board Game Collection", "description": "Gather friends and family for game
        {"name": "Kid's Building Blocks Set", "description": "Foster imagination and creativi
    ],
```

```
237        "Automotive": [
238            {"name": "Car Care Kit", "description": "Keep your vehicle in top condition with our
239            {"name": "Portable Tire Inflator", "description": "Stay prepared on the road with our
240            {"name": "HD Dash Cam", "description": "Capture every moment on the road with our HD
241        ],
242        "Books and Literature": [
243            {"name": "Bestselling Mystery Novel", "description": "Dive into a gripping mystery wi
244            {"name": "Personal Development Guide", "description": "Embark on a journey of self-di
245            {"name": "Illustrated Children's Book", "description": "Spark imagination and joy wit
246        ],
247        "Garden and Outdoor": [
248            {"name": "Solar-Powered Garden Lights", "description": "Illuminate your garden with c
249            {"name": "Folding Outdoor Lounge Chair", "description": "Relax in style with our Fold
250            {"name": "Weather-Resistant Patio Umbrella", "description": "Enjoy outdoor living to
251        ],
252        "Sportswear": [
253            {"name": "High-Performance Running Shoes", "description": "Achieve your fitness goals
254            {"name": "Moisture-Wicking Athletic Shirt", "description": "Stay cool and dry during
255            {"name": "Compression Fit Leggings", "description": "Enhance your performance with ou
256        ],
257        "Jewelry": [
258            {"name": "Elegant Diamond Necklace", "description": "Adorn yourself with our Elegant
259            {"name": "Stylish Silver Bracelet", "description": "Complete your look with our Styli
260            {"name": "Classic Gold Hoop Earrings", "description": "Make a statement with our Clas
261        ],
262        "Skincare": [
263            {"name": "Hydrating Facial Moisturizer", "description": "Nourish and hydrate your ski
264            {"name": "Gentle Cleansing Foam", "description": "Achieve a clean and refreshed compl
265            {"name": "Anti-Aging Serum", "description": "Turn back the clock with our Anti-Aging
266        ],
267        "Health Supplements": [
268            {"name": "Multivitamin Capsules", "description": "Support your overall health with ou
269            {"name": "Omega-3 Fish Oil Softgels", "description": "Boost heart health with our Ome
270            {"name": "Immune System Booster Tablets", "description": "Enhance your immune system
271        ],
272        "Board Games": [
273            {"name": "Strategic Card Game", "description": "Engage in thrilling battles of strate
274            {"name": "Classic Chess Set", "description": "Exercise your mind with our Classic Che
275            {"name": "Family-Friendly Board Game", "description": "Create lasting memories with c
276        ],
277         "Car Engine Products": [
278            {"name": "Performance Gear Oil", "description": "Enhance your car's performance with
```

```python
            {"name": "Engine Degreaser Spray", "description": "Keep your engine clean and running
            {"name": "Heavy-Duty Transmission Fluid", "description": "Ensure optimal transmissio
        ],
        "Gardening Tools": [
            {"name": "Premium Garden Pruner", "description": "Achieve precision in your gardening
            {"name": "Durable Garden Trowel", "description": "Dig, plant, and transplant with ea
            {"name": "Extendable Telescopic Lopper", "description": "Reach high branches and tri
        ]


}
def product(filename, num_products=51):
    fieldnames = ['product_id', 'product_name', 'price', 'product_description', 'inventory',

    with open(filename, 'w', newline='') as csvfile:
        csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)

        # Write the header
        csvwriter.writeheader()

        # Generate and write product data to the CSV file
        product_id_num = 1
        product_names_used = set()

        # List to store 10 random seller_ids for duplication
        random_seller_ids = random.sample(range(1, 501), 10)

        for category, products in product_list.items():
            category_id = f'c{(product_id_num - 1) // 3 + 1:02}'  # Calculate category_id ba

            for product_data in products:
                product_name = product_data["name"]
                product_description = product_data["description"]

                # Ensure no duplicates in product names
                while product_name in product_names_used:
                    product_name = fake.word() + ' ' + fake.word()

                product_names_used.add(product_name)

                product_id = f'p{product_id_num:03}'  # Format product_id as 'p' followed by
                price = round(random.uniform(1, 150), 1)  # Use round to ensure two decimal
```

```
321                    inventory = fake.random_int(min=1, max=100)  # Random inventory between 1 and
322                    weight = round(random.uniform(1.00, 10.00), 2)  # Use random.uniform for weig
323                    seller_id = f's{int((product_id_num - 1) / 2) + 1:05}' if product_id_num <= 1
324                    product_views = fake.random_int(min=500, max=1000)  # Random product views be

326                    # Assign category_id based on the pattern
327                    csvwriter.writerow({
328                        'product_id': product_id,
329                        'product_name': product_name,
330                        'price': price,
331                        'product_description': product_description,
332                        'inventory': inventory,
333                        'weight': weight,
334                        'category_id': category_id,
335                        'seller_id': seller_id,
336                        'product_views': product_views
337                    })

339                    product_id_num += 1

341    # Call the function to generate and save product data to a CSV file
342    product('product.csv', num_products=51)


345    # %% discount
346    def discount(filename, num_discounts=500):
347        fieldnames = ['discount_id', 'discount_percentage', 'discount_start_date', 'discount_end_

349        with open(filename, 'w', newline='') as csvfile:
350            csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)

352            # Write the header
353            csvwriter.writeheader()

355            # Generate and write discount data to the CSV file
356            for discount_id_num in range(1, num_discounts + 1):
357                discount_id = f'd{discount_id_num:05}'  # Format discount_id as 'd' followed by 5
358                discount_percentage = min(round(fake.random.randint(1, 15) * 5 / 100, 2), 0.8)  #
359                start_date = fake.date_between(start_date='-1y', end_date='now')  # Start date wi
360                end_date = start_date + timedelta(days=30)  # End date is 1 month after the star

362                # Assign product_id with some random duplication for the first 50 to 100 products
```

```python
            if 50 <= discount_id_num <= 100:
                product_id = f'p{random.randint(20, 30):03}'
            else:
                product_id = f'p{random.randint(1, 30):03}'

            csvwriter.writerow({
                'discount_id': discount_id,
                'discount_percentage': discount_percentage,
                'discount_start_date': start_date.strftime('%Y-%m-%d'),
                'discount_end_date': end_date.strftime('%Y-%m-%d'),
                'product_id': product_id
            })

# Call the function to generate and save discount data to a CSV file
discount('discount.csv', num_discounts=500)

# %% order

general_reviews = [
    "Great product! Very satisfied with my purchase.",
    "Highly recommend this item. Excellent quality.",
    "Good value for money. Happy with my choice.",
    "Exactly as described. No complaints here.",
    "Impressed with the functionality. Works well.",
    "Smooth transaction. Quick delivery and good packaging.",
    "Nice product. Met my expectations.",
    "Easy to use and durable. Very pleased.",
    "Would buy again. Reliable and efficient.",
    "Overall, a positive experience with this product."
]

# Function to generate order data
def generate_order_data(filename, num_orders=500):
    fieldnames = ['order_number', 'payment_method', 'order_date', 'quantity', 'review', 'cus

    with open(filename, 'w', newline='') as csvfile:
        csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)

        # Write the header
        csvwriter.writeheader()

        # Initialize variables to keep track of the last order_number and payment_method
```

```
405            last_order_number = None
406            last_payment_method = None
407
408            # Generate unique product IDs for the first 50 orders
409            unique_product_ids = set(f'p{i:03}' for i in range(1, 52))
410
411            # Generate and write order data to the CSV file
412            for order_number_num in range(1, num_orders + 1):
413                # Generate order_number and customer_id based on the pattern for the first 100 o
414                if order_number_num <= 100:
415                    order_number = f'on{(order_number_num - 1) // 2 + 1:05}'
416                    customer_id = f'c{(order_number_num - 1) // 2 + 1:05}'
417
418                    if unique_product_ids:
419                        # Use unique product IDs for the first 50 orders
420                        product_id = unique_product_ids.pop()
421                    else:
422                        # If the set is empty, generate random product IDs
423                        product_id = f'p{random.randint(1, 51):03}'
424                else:
425                    order_number = f'on{order_number_num - 50:05}' # Start from previous order n
426                    customer_id = f'c{random.randint(1, 500):05}'
427
428                    # Generate random product IDs for orders after the first 50 orders
429                    product_id = f'p{random.randint(1, 51):03}'
430
431                quantity = fake.random_int(min=1, max=3)  # Random quantity between 1 and 3
432                review = random.choice(general_reviews) if fake.boolean(chance_of_getting_true=6
433
434                # Order_date is the same as order_number
435                order_date = fake.date_between(start_date='-1y', end_date='now').strftime('%Y-%m
436
437                # Payment method takes reference from order_number for cash transactions
438                if last_order_number != order_number:
439                    last_order_number = order_number
440                    last_payment_method = fake.random_element(elements=('Credit Card', 'PayPal',
441
442                payment_method = last_payment_method
443
444                # Generate shipment_id_num based on order_number
445                shipment_id = f'sh{order_number[2:]}'
446
```

```python
447                # Generate customer_rating
448                customer_rating = random.randint(1, 5)
449
450                csvwriter.writerow({
451                    'order_number': order_number,
452                    'payment_method': payment_method,
453                    'order_date': order_date,
454                    'quantity': quantity,
455                    'review': review,
456                    'customer_id': customer_id,
457                    'product_id': product_id,
458                    'shipment_id': shipment_id,
459                    'customer_rating': customer_rating
460                })
461
462    # Call the function to generate and save order data to a CSV file
463    generate_order_data('order.csv', num_orders=500)
464
465
466    # %% shipment
467    def shipment(filename, num_shipments=450):
468        fieldnames = ['shipment_id', 'shipment_delay_days', 'shipment_cost', 'order_number', 're
469
470        with open(filename, 'w', newline='') as csvfile:
471            csvwriter = csv.DictWriter(csvfile, fieldnames=fieldnames)
472
473            # Write the header
474            csvwriter.writeheader()
475
476            # Generate and write shipment data to the CSV file
477            for shipment_id_num in range(1, num_shipments + 1):
478                shipment_id = f'sh{shipment_id_num:05}'  # Format shipment_id as 'sh' followed by
479
480                # Reference order_number from the existing order data
481                order_number = f'on{shipment_id_num:05}'  # Assuming order_number follows the sa
482
483                shipment_delay_days = fake.random_int(min=1, max=3)
484                shipment_cost = round(random.uniform(1, 4), 1)
485
486                # 5% chance of 'Yes', 95% chance of 'No'
487                refund = 'Yes' if fake.random_int(min=1, max=100) <= 5 else 'No'
488
```

71

```
489          csvwriter.writerow({
490              'shipment_id': shipment_id,
491              'shipment_delay_days': shipment_delay_days,
492              'shipment_cost': shipment_cost,
493              'order_number': order_number,
494              'refund': refund
495          })
496
497  # Call the function to generate and save shipment data to a CSV file
498  shipment('shipment.csv', num_shipments=450)
```