

Masters Programmes: Group Assignment Cover Sheet

Student Numbers:	5518804, 2236681, 5577818, 5588817, 5523853
Module Code:	IB9HPO
Module Title:	Data Management
Submission Deadline:	20/03/2024
Date Submitted:	20/03/2024
Word Count:	2000
Number of Pages:	68 (excluding appendix)
Question Attempted:	4
Have you used Artificial Intelligence (AI) in any part of this assignment?	Yes

Academic Integrity Declaration

We're part of an academic community at Warwick. Whether studying, teaching, or researching, we're all taking part in an expert conversation which must meet standards of academic integrity. When we all meet these standards, we can take pride in our own academic achievements, as individuals and as an academic community.

Academic integrity means committing to honesty in academic work, giving credit where we've used others' ideas and being proud of our own achievements.

In submitting my work, I confirm that:

- I have read the guidance on academic integrity provided in the Student Handbook and understand the University regulations in relation to Academic Integrity. I am aware of the potential consequences of Academic Misconduct.
- I declare that this work is being submitted on behalf of my group and is all our own, , except where I have stated otherwise.
- No substantial part(s) of the work submitted here has also been submitted by me in other credit bearing assessments courses of study (other than in certain cases of a resubmission of a piece of work), and I acknowledge that if this has been done this may lead to an appropriate sanction.
- Where a generative Artificial Intelligence such as ChatGPT has been used I confirm I have abided by both the University guidance and specific requirements as set out in the Student Handbook and the Assessment brief. I have clearly acknowledged the use of any generative Artificial Intelligence in my submission, my reasoning for using it and which generative AI (or AIs) I have used. Except where indicated the work is otherwise entirely my own.
- I understand that should this piece of work raise concerns requiring investigation in relation to any of points above, it is possible that other work I have submitted for assessment will be checked, even if marks (provisional or confirmed) have been published.
- Where a proof-reader, paid or unpaid was used, I confirm that the proof-reader was made aware of and has complied with the University's proofreading policy.

Upon electronic submission of your assessment you will be required to agree to the statements above

Data Management Assignment Group 18

Table of contents

Introduction	3
Part 1 Database Design and Implementation	4
1.1 Entity-Relationship Modelling	4
1.2 SQL Database Schema Creation	9
1.2.1 Strategic Advantage of ETL over ELT	12
Part 2 Data Generation and Management	12
2.1 Synthetic Data Generation	12
2.2 Data Import and Quality Assurance	20
Reading the files	21
Creating a loop to read all files and list number of rows and columns in each file	21
Structure of data	22
Number of columns and their column names	22
Checking unique id column as the primary key in each table	22
Data Validation for Each Table	23
Referential Integrity for ensuring data integrity	42
Part 3 Data Pipeline Generation	42
3.1 GitHub Repository and Workflow Setup	42
3.2 GitHub Actions for Continuous Integration	43
Part 4 Data Analysis and Reporting	46
4.1 Advanced Data Analysis in R	46
4.1.1 Promotion Discount Trend	46
4.1.2 Promotion Count Trend	48
Promotion Analysis	49
4.1.3 Monthly Revenue Trend	49
Monthly Revenue Analysis	50

4.1.4 Monthly Best-Selling Products	51
Monthly Best-Selling Products Analysis	53
4.1.5 Monthly Shipping Efficiency	53
4.1.6 Monthly Delivery Efficiency	55
Shipping and Delivery Efficiency Analysis	58
4.2 Comprehensive Reporting with Quarto	58
4.2.1 Demographic Distribution of Customers	58
A. The Distribution of Gender across Customers	58
B. The Distribution of Age across Customers	58
C. The Distribution of Careers across Customers (Top 10)	59
D. The Distribution of Geographic Location across Customers (Top 10)	60
E. The Current Customer Referral Rate	62
Customer Analysis	63
4.2.2 Product Portfolio	63
A. The Distribution of Product Review Scores (Top 10)	63
B. The Number of Products supplied by Different Suppliers	64
C. The Product Review Scores of Different Suppliers (Best Top 5)	65
D. The Product Review Scores of Different Suppliers (Worst Top 5)	65
E. The Top 10 Best Selling Products	66
Supplier Analysis	67
Product Analysis	67
4.2.3 Sales Analysis	67
A. Order Refund Rate	67
Order Refund Analysis	68
Appendices	68

Introduction

In this project, the aim is to emphasise the ETL flow, automation and data analysis, outlining the framework for managing e-commerce data environment, covering end-to-end data management. The framework entails 4 major steps as follows.

1) Business Requirement - This involves thorough understanding of the raw data, including aspects related to business.

2) Conceptual Database Design - This includes Entity-Relationship (E-R) modelling, which refers to the creation of the entities, their attributes and the intricate relationships between them. This is the core stage for creating the Database Management System.

3) Logical Database Design - This involves translating the E-R model into the relational database, in which the entities are represented in the form of tables.

4) Physical Database Design - This step involves creation of the tables using Data Definition Language (DDL) commands to store the data within the database, including the decisions on data types of the attributes.

Part1 Database Design and Implementation

1.1 Entity-Relationship Modelling

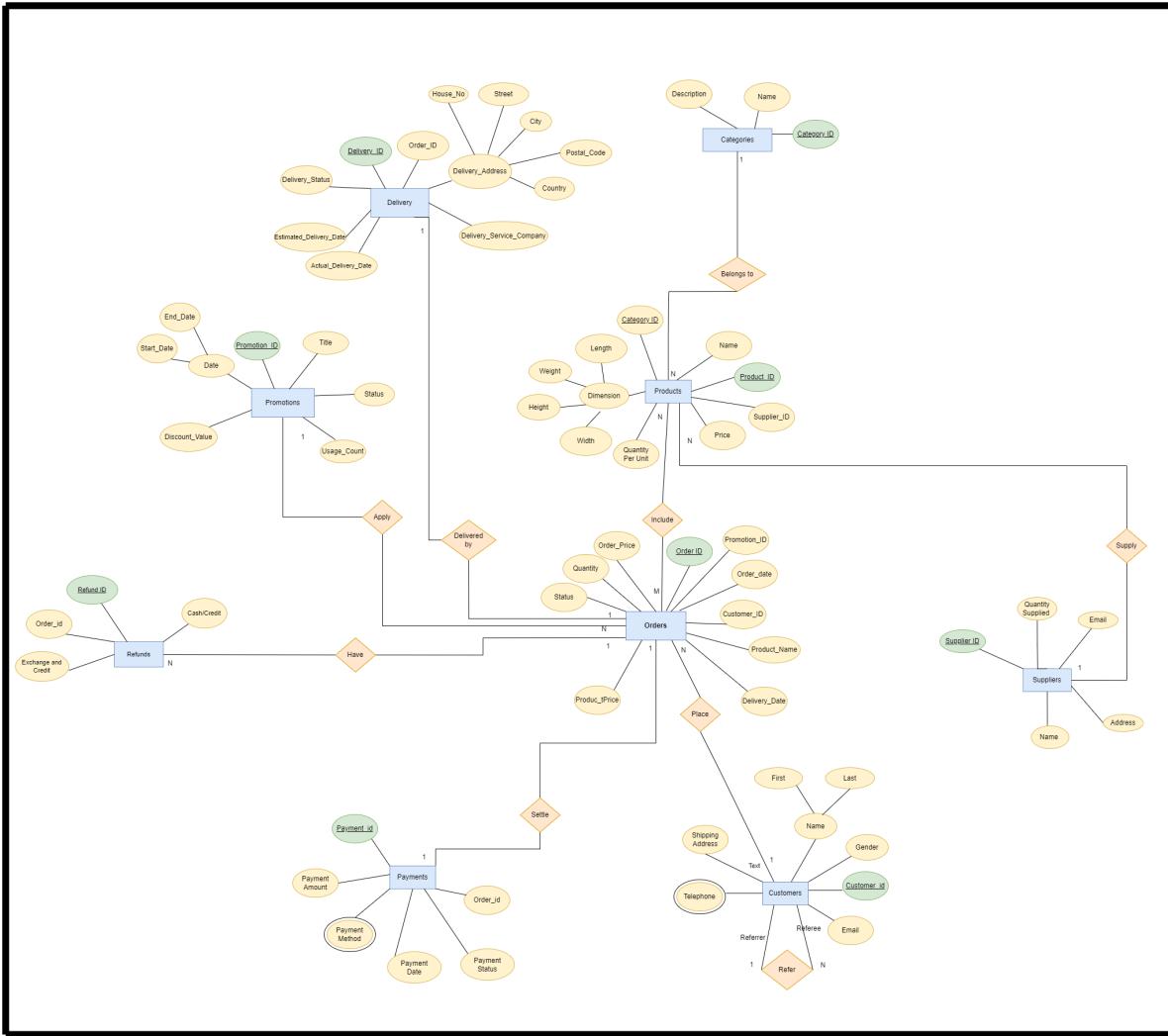


Figure 1: Initial Draft of E-R Diagram

After the initial draft of the E-R diagram, subsequent analysis led to the significant refinements aimed at enhancing the efficiency and clarity of the database structure.

1. Transformation of the “orders” entity into a relationship between “customer” and “product” is justified because of the recognition that an order represents a transaction initiated

by a customer for a specific product aligning more closely with real-world e-commerce processes. Therefore, order is an associative entity in the final E-R diagram.

2. Removal of the “payment” entity was justified by the ability to calculate the payment amounts dynamically based on the quantity and price attributes of the order relationship and product entity respectively.
 3. The “refund” entity was replaced with the refund status attribute within the order relationship. Recognising that the status of the refund is inherently tied to a specific order, thus making it as an attribute within the order relationship rather than a separate entity.
 4. The “delivery” entity was replaced with a more detailed “shipment” entity, accommodating attributes such as shipment_id, shipment_date, and delivery_date, to optimise delivery tracking.

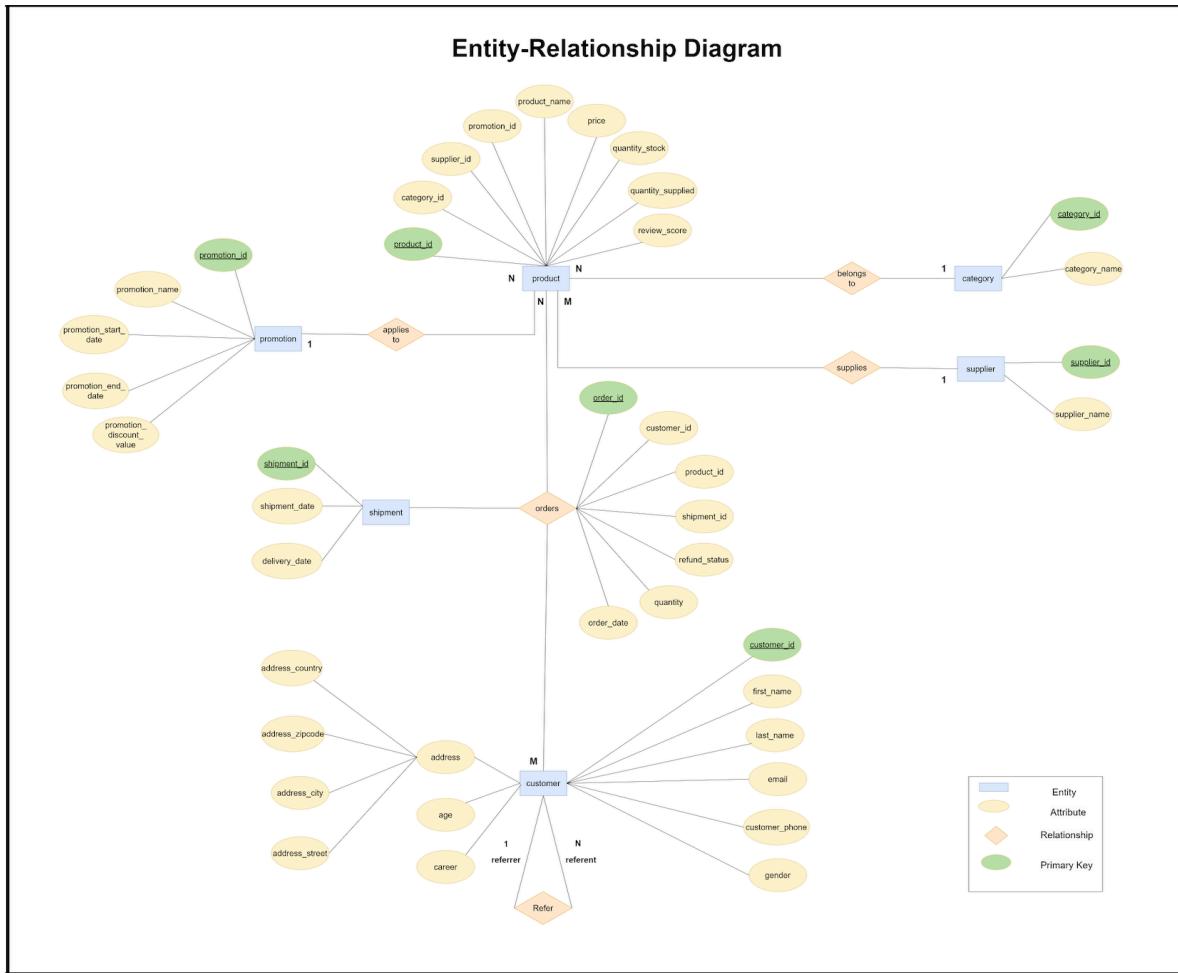


Figure 2: Final ER diagram

After making all the above significant changes the revised E-R diagram has 6 entities that are customer, category, supplier, promotion, shipment, and product. The assumptions that have been taken into consideration for the above E-R modelling are:

- Promotion is only applied to the product.
- There is only one and full final payment for one order that will be calculated using the price, quantity and promotion discount value.
- One order is being shipped and delivered at once, and one order can have many products.
- If one supplier is supplying a particular product, then that product is not going to be supplied by any other supplier.

Logical Database Schema:

The below represents the logical database schema comprising entities and their attributes derived from the E-R diagram. A single underline represents the Primary Key and a double underline represents the Foreign Key in the table.

1. customer {customer_id, first_name, last_name, email, gender, age, career, customer_phone, address_country, address_zipcode, address_city, address_street, referred_by}
2. category {category_id, category_name}
3. supplier {supplier_id, supplier_name}
4. promotion {promotion_id, promotion_name, promotion_start_date, promotion_end_date, promotion_discount_value}
5. shipment {shipment_id, shipment_date, delivery_date}
6. product {product_id, category_id, supplier_id, promotion_id, product_name, price, quantity_stock, quantity_supplied, review_score}
7. orders {order_id, customer_id, product_id, shipment_id, quantity, refund_status, order_date}

Figure 3: Logical Schema

Relationship sets:

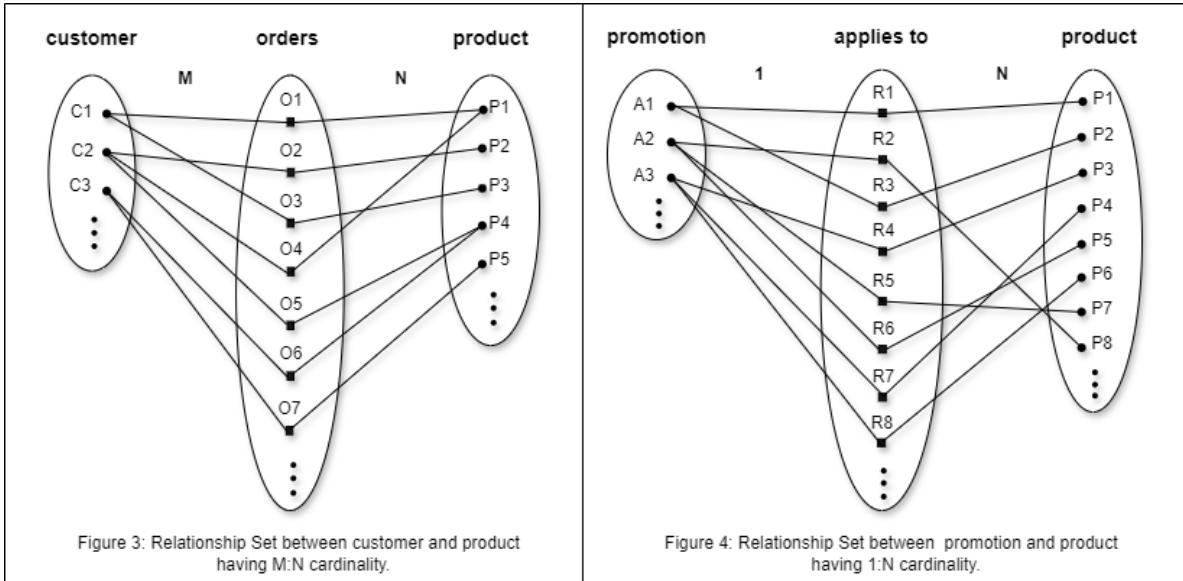


Figure 3 illustrates many-to-many (M:N) relationship between customer and product, indicating that multiple customers can order multiple products, and conversely multiple products can be ordered by various customers.

Figure 4 illustrates one-to-many (1:N) relationship between promotion and product, that a single promotion code may be applied to multiple products, whereas multiple products can be associated with a single promotion code only.

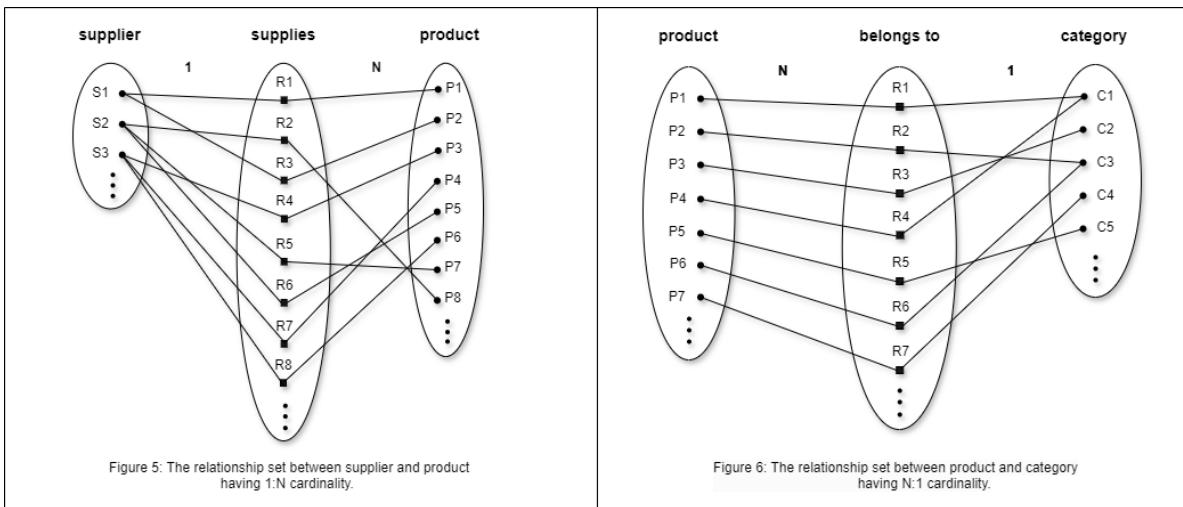


Figure 5 illustrates one-to-many (1:N) relationship between supplier and product. This indicates that a single supplier can supply multiple products, and multiple products can be supplied by a single supplier.

Figure 6 illustrates many-to-one (N:1) relationship between product and category, where multiple products can be categorised under a single category, while each category can have multiple products.

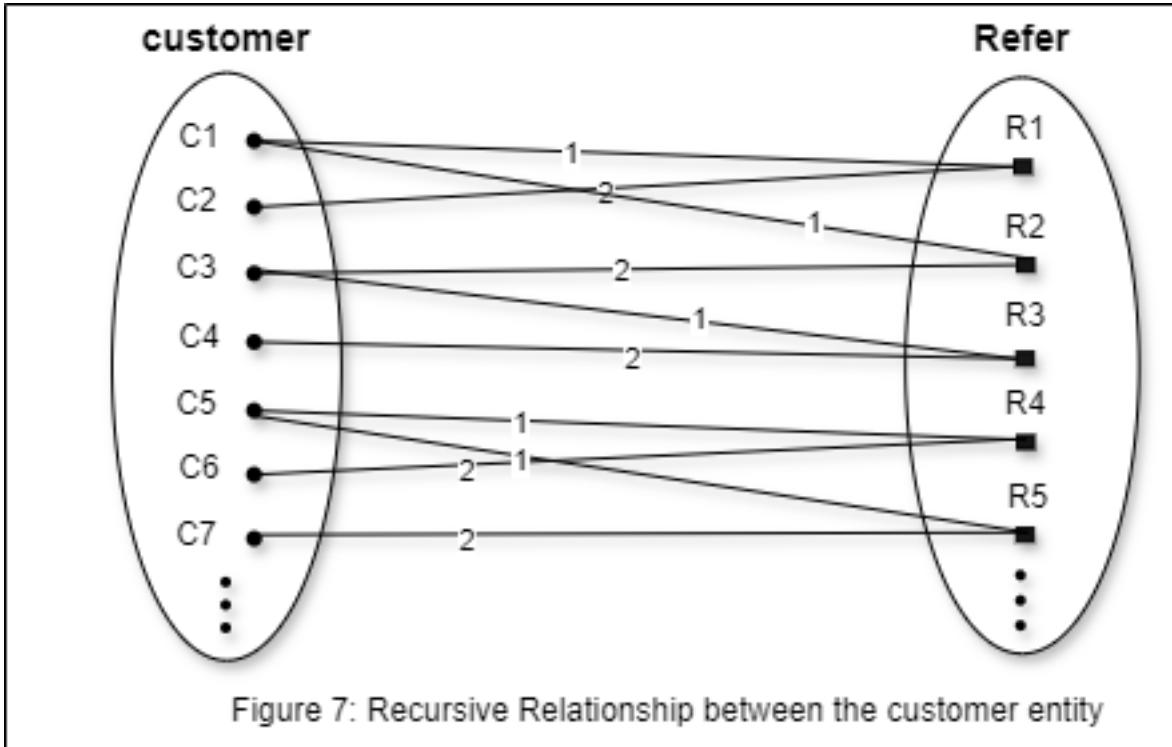


Figure 7 illustrates one-to-many (1:N) recursive relationship involving the customer entity. This indicates that while one customer can refer multiple customers, where each customer can only be referred by one customer. The number “1” above the relationship lines represents the referrer role and “2” represents the referent.

1.2 SQL Database Schema Creation

In this part, E-R diagram is translated into a functional SQL database schema.

Creating table schema for all tables

1. Creating customer table

```
1 RSQLite::dbExecute(connection,"  
2  
3 CREATE TABLE IF NOT EXISTS customer (  
4  
5     customer_id VARCHAR(10) PRIMARY KEY,  
6     first_name VARCHAR(50) NOT NULL,  
7     last_name VARCHAR(50) NOT NULL,  
8     email VARCHAR(50) NOT NULL,  
9     gender VARCHAR(20) NOT NULL,  
10    age INT NOT NULL,  
11    career VARCHAR(50) NOT NULL,  
12    customer_phone VARCHAR(20) NOT NULL,  
13    address_country VARCHAR(50) NOT NULL,  
14    address_zipcode VARCHAR(20) NOT NULL,  
15    address_city VARCHAR(20) NOT NULL,  
16    address_street VARCHAR(50) NOT NULL,  
17    referred_by VARCHAR(10) NULL  
18  
19 );  
20 ")
```

2. Creating category table

```
1 RSQLite::dbExecute(connection,"  
2  
3 CREATE TABLE IF NOT EXISTS category (  
4  
5     category_id VARCHAR(10) PRIMARY KEY,  
6     category_name VARCHAR(50) NOT NULL  
7  
8 );  
9  
10 ")
```

3. Creating supplier table

```
1 RSQLite::dbExecute(connection,"  
2 CREATE TABLE IF NOT EXISTS supplier (  
3  
4     supplier_id VARCHAR(10) PRIMARY KEY,  
5     supplier_name VARCHAR(50) NOT NULL  
6  
7 );  
8  
9 ")
```

4. Creating promotion table

```
1 RSQLite::dbExecute(connection,"  
2  
3 CREATE TABLE IF NOT EXISTS promotion (  
4  
5     promotion_id VARCHAR(10) PRIMARY KEY,  
6     promotion_name VARCHAR(20) NOT NULL,  
7     promotion_start_date DATE NOT NULL,  
8     promotion_end_date DATE NOT NULL,  
9     promotion_discount_value FLOAT NOT NULL  
10  
11 );  
12  
13 ")
```

5. Creating shipment table

```
1 RSQLite::dbExecute(connection,"  
2  
3 CREATE TABLE IF NOT EXISTS shipment (  
4  
5     shipment_id VARCHAR(10) PRIMARY KEY,  
6     shipment_date DATE NOT NULL,  
7     delivery_date DATE NOT NULL  
8 );  
9  
10 ")
```

6. Creating product table

```

1 RSQLite::dbExecute(connection,"
2
3 CREATE TABLE IF NOT EXISTS product (
4
5     product_id VARCHAR(10) PRIMARY KEY,
6     category_id VARCHAR(10) NOT NULL,
7     supplier_id VARCHAR(10) NOT NULL,
8     promotion_id VARCHAR(10) NULL,
9     product_name VARCHAR(20) NOT NULL,
10    price INT NOT NULL,
11    quantity_stock INT NOT NULL,
12    quantity_supplied INT NOT NULL,
13    review_score FLOAT NOT NULL,
14    FOREIGN KEY ('category_id') REFERENCES category('category_id'),
15    FOREIGN KEY ('supplier_id') REFERENCES supplier('supplier_id'),
16    FOREIGN KEY ('promotion_id') REFERENCES promotion('promotion_id')
17
18 );
19
20 ")

```

7. Creating orders table

```

1 RSQLite::dbExecute(connection,"
2
3 CREATE TABLE IF NOT EXISTS orders (
4
5     order_id VARCHAR(20) NOT NULL,
6     customer_id VARCHAR(10) NOT NULL,
7     product_id VARCHAR(10) NOT NULL,
8     shipment_id VARCHAR(10) NOT NULL,
9     quantity INT NOT NULL,
10    refund_status VARCHAR(20) NOT NULL,
11    order_date DATE NOT NULL,
12    FOREIGN KEY ('customer_id') REFERENCES customer('customer_id'),
13    FOREIGN KEY ('product_id') REFERENCES product('product_id'),
14    FOREIGN KEY ('shipment_id') REFERENCES shipment('shipment_id')
15    PRIMARY KEY (order_id,product_id,customer_id,shipment_id)
16
17 );
18 ")

```

After creating the schema, below lists all of the tables from the database to check what have been created.

```
1 RSQLite::dbListTables(connection)
```

1.2.1 Strategic Advantage of ETL over ELT

In the development of the database, the principle of ETL (Extract, Transform, Load) was used rather than ELT (Extract, Load, Transform) as ETL can clean and transform the data before loading it into the database tables. This indicates that the data warehouse or database has only verified data, resulting in good data quality and consistency and the transformation happens before loading. Any mistakes or inconsistencies may be discovered and handled immediately, lowering the risk of ingesting the faulty data into the database. The process employed in this project is described below -

Extract – The data was generated using both Large Language Models (LLMs) and Mockaroo, which aligns with the extraction phase, where data is collected from various sources and prepared for further processing.

Transform – The extracted data was converted into CSV files and was further performed some validation tests. This step was to ensure the data quality and consistency, which is the most crucial part of the ETL transformation. Furthermore, normalisation processes were used to standardise the data, for improving its integrity.

Load - Finally, the validated and normalised data was loaded into the tables in the database for further analysis.

Part 2 Data Generation and Management

2.1 Synthetic Data Generation

Based on the logical schema, the physical schema involved establishing the details in the database, including tables, columns, relationships, and constraints. Furthermore, the attributes of each entity ensured that the normalisation issues were resolved. Following this, the dependency between entities prioritised the order when generating the data. Thus, the data within the 1:N entity were created before the M:N entity due to the foreign key for the associative entity.

Order Entity / Associative Entity Number of Observations

1	supplier	20
2	category	20
3	promotion	30
4	customer	200
5	product	200
6	shipment	389

Based on the nature of each attribute, customer data was created through “Mockaroo” to obtain the initial data.

The screenshot shows the Mockaroo website's schema editor. At the top, there are navigation links: SCHEMAS, DATASETS, MOCK APIs, SCENARIOS, PROJECTS, and FUNCTIONS. There are also links for HELP, SIGN IN, and UPGRADE NOW. Below the header, a table lists fields with their types and configuration options:

Field Name	Type	Options
customer_id	Row Number	blank: 0 % Σ X
first_name	First Name	blank: 0 % Σ X
last_name	Last Name	blank: 0 % Σ X
email	Email Address	blank: 0 % Σ X
phone_number	Phone	format: ###-###-#### blank: 0 % Σ X
address	Street Address	blank: 0 % Σ X
city	City	blank: 0 % Σ X
state	State	restrict states... All Countries blank: 0 % Σ X
country	Country	restrict countries... blank: 0 % Σ X
zipcode	Postal Code	blank: 0 % Σ X
referred_by	Row Number	blank: 0 % Σ X

Figure 4: Data generation screenshot on mockaroo website

However, the generated data lacks authentic features, such as age, career, and address, and has no logical meaning for each customer. To make data reflect realistic patterns and distributions, further LLM was used to enhance the practical characteristics of the initial data. The following images show the ChatGPT prompts.

JA You



MOCK_DATA_modified.xlsx

Spreadsheet

Based on the attached file, can you help me improve it to be more realistic

1. Age Distribution - Mean: 35 / Standard deviation: 10

2. Career Distribution:

- 12 to 17 years old: Primarily students, which may include part-time jobs or internships.

- 18 to 22 years old: Students, interns, and young individuals who are just entering the workforce (such as entry-

level positions).

- 23 to 30 years old: A wider range of career choices, including positions that require professional degrees or

experience.

- 31 years old and above: This includes all professions, assuming that by this age, individuals have stabilized in

their career paths, possibly covering mid-level to senior positions.

1 customer phone begin with "+44" & one space & 10-digit numbers

Figure 5: Data improvement request to ChatGPT for customer table

 You

[Category Table]

I have a 3NF category table with the following fields:

*category_id: It should be a unique identifier of a product category. It should be a 10-digit length character. It is the primary key of the table. It can't be empty.

*category_name: It should be the name of the product category. It should be in character format. You should use the real and meaningful category name in an e-commerce company. It can't be empty.

Can you generate 20 records for me and provide me the result in an excel file?

 ChatGPT

I've created an Excel file with 20 records for the category table as requested. You can download the file from the link below:

[Download Category_Table.xlsx \[>\]](#)



Figure 6: Data improvement request to ChatGPT for category table

 **You**

[Supplier Table]

I have a 3NF supplier table with the following fields:

*supplier_id: It should be a unique identifier of a product supplier. It should be a 10-digit length character. It is the primary key of the table. It can't be empty.

*category_name: It should be the name of the product supplier. It should be in character format.

You should use the real and meaningful supplier name in an e-commerce company. It can't be empty.

Can you generate 20 records for me and provide me the result in an excel file?

 **ChatGPT**

I've created an Excel file with 20 records for the supplier table as requested. You can download the file from the link below:

[Download Supplier_Table.xlsx \[x-\]](#)

Figure 7: Data improvement request to ChatGPT for supplier table



You

[Promotion Table]

I have a 3NF promotion table with the following fields:

*promotion_id: It should be a unique identifier of a promotion. It should be a 10-digit length character. It is the primary key of the table. It can't be empty.

*promotion_name: It should be the name of a e-commerce promotion. It should be in character format. You should use the real and meaningful promotion name in an e-commerce environment. It can't be empty.

*promotion_start_date: It's the start date of a promotion. It should be in DD-MM-YY format. You should ensure the promotion_start_date is real and meaningful, and match with its the promotion_name. Like the Black Friday promotion should happen in November or December. It can't be empty. The time period should be from 2023 January to 2024 February.

*promotion_end_date: It's the end date of a promotion. It should be in DD-MM-YY format. You should ensure the promotion_end_date is real and meaningful, and match with its the promotion_name. Like the Black Friday promotion should end in November or December. You should also ensure each promotion_id's promotion_end_date is later than its promotion_start_date. It can't be empty. The time period should be from 2023 January to 2024 February.

*promotion_discount_value: It's the discount value of one promotion. It should be a float number in range [0,1]. The new price after promotion should be price * (1-promotion_discount_value). It can't be empty. Its value shouldn't be too high and most of them should be smaller than 0.6.

Can you generate 30 records for me and provide me the result in an excel file?



ChatGPT

I've generated 30 records for the promotion table with the updated requirements and saved them to an Excel file. You can download the file from the link below:

[Download Promotion_Table.xlsx](#)



Figure 8: Data improvement request to ChatGPT for promotion table



You

[Shipment Table]

I have a 3NF shipment table with the following fields:

*shipment_id: It should be a unique identifier of a shipment matching with an order. It should be a 10-digit length character and the first 4 digits should be "SHIP", and the other 6 digits should be numbers. It is the primary key of the table, and it can't be empty.

*shipment_date: It should be the date that an order starts to be shipped. It should be in DD-MM-YY format and the time period should be from 2023 January to 2024 March.

*delivery_date: It should be the date that an order has been received by the customer. It should be in DD-MM-YY format and the time period should be from 2023 January to 2024 March.

Can you generate 389 records for me and provide me the result in an excel file?



ChatGPT

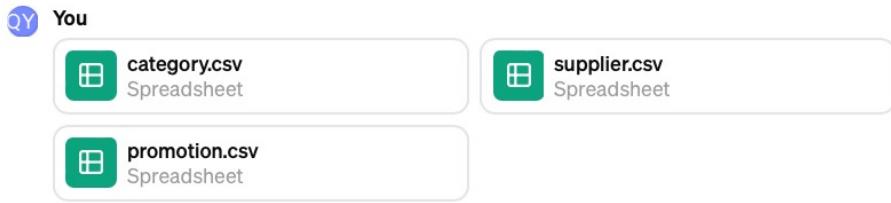
I've generated 389 records for the shipment table as requested and saved them to an Excel file.

You can download the file from the link below:

[Download Shipment_Table.xlsx](#)



Figure 9: Data improvement request to ChatGPT for shipment table



[Product Table]

I already have a 3NF category table, supplier table, and promotion table (I'll show you later). But now I have a 3NF product table with the following fields:

*product_id: It should be a unique identifier of a product, in a 10-digit length character. It is the primary key of the table. It can't be empty. When the category_id is the same, you should keep the first three digit of the product_id is the same, please use alphabets; When the supplier_id is the same, you should keep the 4-6 digits of the product_id is the same, please use numbers; For the 7-10 digits of the product_id, you can randomly give numbers to them.

*category_id: It should be the product's category_id, and you should choose values from the category table's category_id column I provided you, it should be real and meaningful, it can't be empty, all category_id in the category table should be included in this column.

*supplier_id: It should be the product's supplier's id, and you should choose values from supplier table's supplier_id I provide you, it should be real and meaningful, it can't be empty, all supplier id in supplier table should be included in this column.

*promotion_id: It should be the promotion that applied to the product, you should choose values from promotion table's promotion_id I provide you, it should be real and meaningful, it can be empty because some products may not have any promotions. All promotion_id in promotion table should be included in this column.

*product_name: It should be the product's name, you should choose names according to the category_id you choose from category table, which means this product is actually belongs to this category. It should be a real and meaningful product name that actually match with its category_name and category_id. It can't be empty and should be unique and specific. Could you please name the product_name in such format: the specific type of the product(should be a real type like Mattress, not be general like product A) - the specific product brand(should be a real type like John Lewis, not be general like brand A) - the specific product name(like silent night) , for example, [Mattress John Lewis silent night].

*price: It should be the selling price of the product, it should be real and meaningful, and can't be empty, in number format. You should generate price values according to its product_name and category_name, make it as real as possible.

*quantity_stock: It should be the current stock quantity of the product. It can't be empty and should be integer, be real and meaningful.

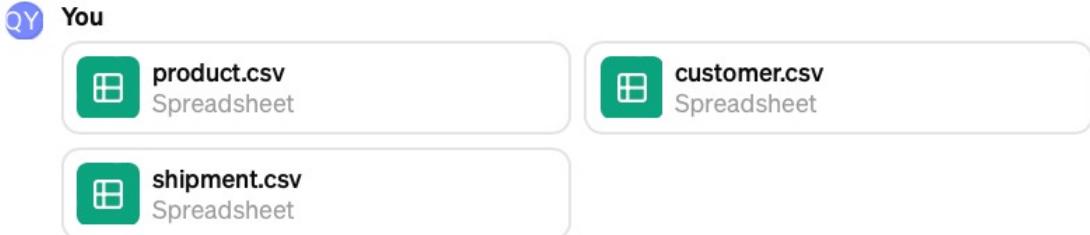
*quantity_supplied: It should be the total supply amount of the product that is supplied by its suppliers. It can't be empty and should be integer, be real and meaningful.

*average_score: it should be the product's average review score provided by its customers, it should be real and meaningful. It can't be empty, and should be integer in the range [1,5].

Can you generate 200 records for me and provide me the result in an excel file?

 ChatGPT
With the structure of the category, supplier, and promotion tables reviewed, I'll now proceed to generate 200 records for the product table. This process includes creating unique product identifiers that meet the specified requirements, assigning category IDs, supplier IDs, promotion

Figure 10: Data improvement request to ChatGPT for product table



[Order Table]

I already have a 3NF customer table, product table and shipment table(I'll show you later). But now I have a 3NF order table with the following fields:

*order_id: It should be the unique identifier, it should be in a 14 digit length character and begin with "ORDER". It's the primary key of the table and can't be empty, but order_id can be duplicated because one order may have many different products.

*customer_id: It should be the customer who places this order. You should choose values from the customer table's customer_id column I provided you, one customer_id can place many order_id, and it could not be empty, but can be duplicated because one customer can place many different orders.

*product_id: It should be the product's id that included in this order. You should choose values from the product table's product_id column I provided you. It can not be empty but it can be duplicated because one product can be included in many different orders.

*shipment_id: It should be the shipment's id of this order. You should choose values from the shipment table's shipment_id column I provided you. One order_id only match with one shipment_id. It can not be empty but can be duplicated because order_id can also be duplicated.

*quantity: It should be the quantity of the product in one order. It should be a positive integer and can't be empty. The value should be real and meaningful.

*refund_status: It should be characters and has two values "yes" "no" which means whether this order is requested to be refunded or not , the number of "No" should be a little bit more than "Yes".

*order_date: It should be the date when the order is placed. It can't be empty and should in DD-MM-YY format. You should make sure the order_date should be same or a little bit earlier than the shipment_date in shipment table.

Can you generate 1000 records for me and provide me the result in an excel file?

Figure 11: Data improvement request to ChatGPT for order table

2.2 Data Import and Quality Assurance

The ETL approach included extracting data from CSV files, transformed it through validations for data quality, referential integrity, and ensured accurate and normalised data in the database (ecomdata.db). The R scripts validate field correctness, including email formats, age ranges, and proper date formats, etc. Duplicate entries are prevented by checking existing records in

the database and within incoming files. Referential integrity for the ‘orders’, ‘product’, and ‘customer’ tables is maintained by verifying foreign keys before insertion. Additionally, hash checks before and after data loading for ‘orders’, ‘product’, and ‘shipment’ tables safeguard against data alteration.

```
Customer ID does not exist: NA
Product ID does not exist: NA
Shipment ID does not exist: NA
Referential integrity for the row is invalid.
Starting validation for checked new records.
Validation completed for new records.
Starting to insert validated data into the database. Number of records:  2
Data insertion completed successfully.
Data integrity verified: All record hashes match.
Successfully moved processed file to: data_processed/orders6_20240319205817.csv
Finished processing file: data_upload/orders6.csv
```

The output log provides the data importation and integrity verification process for orders table. It excludes the rows with referential integrity issues, confirms the successful insertion of two validated records into the database, and verifies data integrity through matching hash checks. Also, it documents the relocation of processed files to an organised directory.

Reading the files

Firstly, all the CSV files are listed before doing the validations.

```
1 all_files <- list.files("data_upload/")
2 all_files
```

Creating a loop to read all files and list number of rows and columns in each file

```
1 for (variable in all_files) {
2   filepath <- paste0("data_upload/",variable)
3   file_contents <- readr::read_csv(filepath)
4
5   number_of_rows <- nrow(file_contents)
6   number_of_columns <- ncol(file_contents)
7
8   #Printing the number of rows and columns in each file
9   print(paste0("The file: ",variable,
10             " has: ",
```

```

11     format(number_of_rows,big.mark = ","),
12     " rows and ",
13     number_of_columns," columns"))
14
15 number_of_rows <- nrow(file_contents)
16
17 print(paste0("Checking for file: ",variable))
18
19
20 #Printing True if the first column is the primary key column else printing False
21 print(paste0(" is ",nrow(unique(file_contents[,1]))==number_of_rows))
22 }
```

Structure of data

```

1 for (variable in all_files) {
2   filepath <- paste0("data_upload/",variable)
3   file_contents <- readr::read_csv(filepath)
4   str_data <-str(file_contents)
5   print(paste0(str_data,"Srtucture of the file ", variable))
6 }
```

Number of columns and their column names

```

1 for (variable in all_files) {
2   filepath <- paste0("data_upload/",variable)
3   file_contents <- readr::read_csv(filepath)
4   column_names <-colnames(file_contents)
5   print(paste0("File ", variable, " has column as ",column_names))
6 }
```

Checking unique id column as the primary key in each table

```

1 for (variable in all_files) {
2   filepath <- paste0("data_upload/", variable)
3   file_contents <- readr::read_csv(filepath)
```

```

4 primary_key <- nrow(unique(file_contents[,1])) == nrow(file_contents)
5 print(paste0("Primary key of ", variable, " is ", primary_key))
6 }

```

Data Validation for Each Table

1. Validation for customer data

```

1 fetch_existing_customer_ids <- function(connection) {
2   query <- "SELECT DISTINCT customer_id FROM customer"
3   existing_ids <- dbGetQuery(connection, query)$customer_id
4   return(existing_ids)
5
6 }
7
8 #Customer data validation and Reffered by referencial integrity
9 validate_and_prepare_customer_data <- function(data, existing_ids) {
10
11   #Validation for customer ID
12   customer_id_check <- grepl("^\w{10}$", data$customer_id)
13   data <- data[customer_id_check, ]
14
15   #Validation for email
16   email_check <- grepl("^\w+@\w+\.\w{2,3}$", data$email)
17   data<- data[email_check, ]
18
19   #Validation for gender
20   gender_check <- c("male","female","Female","Male")
21   data<- data[data$gender %in% gender_check, ]
22
23   #Validation for age
24   age_check <- 1:100
25   data <- data[data$age %in% age_check, ]
26
27   #Validation for phone number
28   phone_check <- grepl("^\+\d{10}$", data$customer_phone)
29   data <- data[phone_check, ]
30
31   #Validation for zip code
32   zipcode_check <- grepl("^\w{3} \w{3}$", data$address_zipcode)
33   data <- data[zipcode_check, ]

```

```

34
35 #Referred by check
36 unique_customer_ids <- unique(c(data$customer_id, existing_ids))
37 #Validate 'referred by' IDs
38 valid_referral_flags <- (data$referred_by == "") | data$referred_by %in% unique_customer_ids
39 data <- data[valid_referral_flags, ]
40
41 return(data)
42 }

43
44 #Fetch existing customer IDs from the database
45 existing_customer_ids <- fetch_existing_customer_ids(connection)
46
47 customer_file_paths <- list.files(path = "data_upload", pattern = "customer.*\\.csv$", full.names = TRUE)
48 #Initialising empty dataframe
49 customer_possible_data <- data.frame()
50
51 customer_primary_key <- "customer_id"
52
53 #Read each customer CSV file and check for the existence of the primary key in the database
54 for (file_path in customer_file_paths) {
55   cat("Starting processing file:", file_path, "\n")
56   #Read the current file
57   customer_data <- read.csv(file_path)
58
59   #Iterate through each row of the file
60   for (i in seq_len(nrow(customer_data))) {
61     new_record <- customer_data[i, ]
62     primary_key_value <- new_record[[customer_primary_key]]
63     conditions <- paste(customer_primary_key, "=", paste0("'", primary_key_value, "'"))
64
65     #Check if a record with the same primary key exists in the database
66     record_exists_query <- paste("SELECT COUNT(*) FROM customer WHERE", conditions)
67     record_exists_result <- dbGetQuery(connection, record_exists_query)
68     record_exists <- record_exists_result[1, 1] > 0
69
70     if(record_exists) {
71       cat("Record with primary key", primary_key_value, "already exists in the database.\n")
72     }
73     if (!record_exists) {
74       #Check if the primary key value of the new record is unique in the temporary dataframe
75       if (!primary_key_value %in% customer_possible_data[[customer_primary_key]]) {

```

```

76     customer_possible_data <- rbind(customer_possible_data, new_record)
77   }
78 }
79
80 cat("Finished processing file:", file_path, "\n")
81
82 }
83 cat("Starting validation for new records.\n")
84 customer_possible_data <- validate_and_prepare_customer_data(customer_possible_data, existing_data)
85 cat("Validation completed for new records.\n")
86 }
87
88
89 if (nrow(customer_possible_data) > 0)
90 {
91   cat("Starting to insert validated data into the database. Number of records: ", nrow(customer_possible_data))
92
93 #Ingesting prepared data to our database
94 dbWriteTable(connection, name = "customer", value = customer_possible_data, append = TRUE,
95   cat("Data insertion completed successfully.\n"))
96 } else
97 {
98   cat("No valid customer data to insert into the database.\n")
99 }
```

Recursive relationship in the customer table

In the ‘customer’ table, ‘customer_id’ serves as the primary key, that uniquely identifies each customer and ‘referred_by’ is the attribute that defines the customer who refers this customer serving as the link between different customers within the same customer table.

2. Validations for category data

```

1 validate_and_prepare_category_data <- function(data) {
2
3   # Validation for category ID
4   category_id_check <- grepl("^[A-Za-z0-9]{10}$", data$category_id)
5   data <- data[category_id_check,]
6
7   return(data)
8 }
9
```

```

10 # Fetch existing category IDs from the database
11
12 category_file_paths <- list.files(path = "data_upload", pattern = "category.*\\.csv$", full.names = TRUE)
13
14 # Define the primary key column for the category table
15 category_primary_key <- "category_id"
16
17 #Initialising empty data frame
18 category_possible_data <- data.frame()
19
20 # Read each category CSV file and check for the existence of the primary key in the database
21 for (file_path in category_file_paths) {
22
23   cat("Starting processing file:", file_path, "\n")
24
25   # Read the current file
26   category_data <- readr::read_csv(file_path)
27
28   # Iterate through each row of the file
29   for (i in seq_len(nrow(category_data))) {
30     new_record <- category_data[i, ]
31     primary_key_value <- new_record[[category_primary_key]]
32     conditions <- paste(category_primary_key, "=", paste0("'", primary_key_value, "'"))
33
34     # Check if a record with the same primary key exists in the database
35     record_exists_query <- paste("SELECT COUNT(*) FROM category WHERE", conditions)
36     record_exists_result <- dbGetQuery(connection, record_exists_query)
37     record_exists <- record_exists_result[1, 1] > 0
38
39     if(record_exists) {
40       cat("Record with primary key", primary_key_value, "already exists in the database.\n")
41     }
42     if (!record_exists) {
43       # Check if the primary key value of the new record is unique in the temporary datafram
44       if (!primary_key_value %in% category_possible_data[[category_primary_key]]) {
45         category_possible_data <- rbind(category_possible_data, new_record)
46       }
47     }
48
49   cat("Finished processing file:", file_path, "\n")
50 }

```

```

52 }
53 cat("Starting validation for new records.\n")
54 category_possible_data <- validate_and_prepare_category_data(category_possible_data)
55 cat("Validation completed for new records.\n")
56
57 if (nrow(category_possible_data) > 0) {
58   cat("Starting to insert validated data into the database. Number of records: ", nrow(categ
59
60   # Ingesting prepared data to our database
61   dbWriteTable(connection, name = "category", value = category_possible_data, append = TRUE,
62   cat("Data insertion completed successfully.\n")
63 } else
64 {
65   cat("No valid category data to insert into the database.\n")
66 }

```

3. Validations for supplier data

```

1 validate_and_prepare_supplier_data <- function(data) {
2   # Validation for supplier ID
3   supplier_id_check <- grepl("^[A-Za-z0-9]{10}$", data$supplier_id)
4   data <- data[supplier_id_check,]
5
6   return(data)
7 }
8
9 # Fetch existing supplier IDs from the database
10
11 supplier_file_paths <- list.files(path = "data_upload", pattern = "supplier.*\\.csv$", full.r
12
13 # Define the primary key column for the supplier table
14 supplier_primary_key <- "supplier_id"
15
16 #Initialising empty data frame
17 supplier_possible_data <- data.frame()
18
19 # Read each supplier CSV file and check for the existence of the primary key in the database
20 for (file_path in supplier_file_paths) {
21
22   cat("Starting processing file:", file_path, "\n")
23

```

```

24 # Read the current file
25 supplier_data <- readr::read_csv(file_path)
26
27 # Iterate through each row of the file
28 for (i in seq_len(nrow(supplier_data))) {
29   new_record <- supplier_data[i, ]
30   primary_key_value <- new_record[[supplier_primary_key]]
31   conditions <- paste(supplier_primary_key, "=", paste0("'", primary_key_value, "'"))
32
33   # Check if a record with the same primary key exists in the database
34   record_exists_query <- paste("SELECT COUNT(*) FROM supplier WHERE", conditions)
35   record_exists_result <- dbGetQuery(connection, record_exists_query)
36   record_exists <- record_exists_result[1, 1] > 0
37
38   if(record_exists) {
39     cat("Record with primary key", primary_key_value, "already exists in the database.\n")
40   }
41   if (!record_exists) {
42     # Check if the primary key value of the new record is unique in the temporary dataframe
43     if (!primary_key_value %in% supplier_possible_data[[supplier_primary_key]]) {
44       supplier_possible_data <- rbind(supplier_possible_data, new_record)
45     }
46   }
47
48   cat("Finished processing file:", file_path, "\n")
49
50 }
51
52 }
53 cat("Starting validation for new records.\n")
54 supplier_possible_data <- validate_and_prepare_supplier_data(supplier_possible_data)
55 cat("Validation completed for new records.\n")
56
57 if (nrow(supplier_possible_data) > 0) {
58   cat("Starting to insert validated data into the database. Number of records: ", nrow(supplier_
59
60   # Ingesting prepared data to our database
61   dbWriteTable(connection, name = "supplier", value = supplier_possible_data, append = TRUE,
62   cat("Data insertion completed successfully.\n")
63 } else
64 {
65   cat("No valid supplier data to insert into the database.\n")

```

66 }

4. Validations for promotion data

```
1 validate_and_prepare_promotion_data <- function(data) {  
2  
3 # Validation for promotion ID  
4 promotion_id_check <- grepl("^[A-Za-z0-9]{10}$", data$promotion_id)  
5 data <- data[promotion_id_check, ]  
6  
7 #Checking for the validation of the promotion_start_date and promotion_end_date in the promotion table.  
8 date_check <- !is.na(as.Date(data$promotion_start_date, format = "%d-%m-%Y")) &  
9 !is.na(as.Date(data$promotion_end_date, format = "%d-%m-%Y")) &  
10 as.Date(data$promotion_start_date, format = "%d-%m-%Y") < as.Date(data$promotion_end_date, format = "%d-%m-%Y")  
11  
12 #Check for the validation of the promotion_start_date and promotion_end_date in the promotion table.  
13 #promotion_start_date and promotion_end_data should be in correct form for eg 12/11/2023  
14 date_format <- "%d-%m-%Y"  
15 date_check <- !is.na(as.Date(data$promotion_start_date, format = date_format)) &  
16 !is.na(as.Date(data$promotion_end_date, format = date_format)) &  
17 as.Date(data$promotion_start_date, format = date_format) < as.Date(data$promotion_end_date, format = date_format)  
18 data <- data[date_check, ]  
19  
20  
21 #Check for the validation of the column promotion_discount_value in the promotion table.  
22 #promotion_discount_value should be <1  
23  
24 discount_value_check <- !is.na(data$promotion_discount_value) &  
25 is.numeric(data$promotion_discount_value) &  
26 data$promotion_discount_value < 1  
27 data <- data[discount_value_check, ]  
28 return(data)  
29 }  
30  
31  
32 # Fetch existing promotion IDs from the database  
33  
34 promotion_file_paths <- list.files(path = "data_upload", pattern = "promotion.*\\.csv$", full.names = TRUE)  
35  
36 # Define the primary key column for the promotion table  
37 promotion_primary_key <- "promotion_id"  
38
```

```

39 #Initialising empty data frame
40 promotion_possible_data <- data.frame()
41
42
43 # Read each promotion CSV file and check for the existence of the primary key in the database
44 for (file_path in promotion_file_paths) {
45
46 cat("Starting processing file:", file_path, "\n")
47
48 # Read the current file
49 promotion_data <- readr::read_csv(file_path)
50
51 # Iterate through each row of the file
52 for (i in seq_len(nrow(promotion_data))) {
53 new_record <- promotion_data[i, ]
54 primary_key_value <- new_record[[promotion_primary_key]]
55 conditions <- paste(promotion_primary_key, "=", paste0("'", primary_key_value, "'"))
56
57 # Check if a record with the same primary key exists in the database
58 record_exists_query <- paste("SELECT COUNT(*) FROM promotion WHERE", conditions)
59 record_exists_result <- dbGetQuery(connection, record_exists_query)
60 record_exists <- record_exists_result[1, 1] > 0
61
62 if(record_exists) {
63   cat("Record with primary key", primary_key_value, "already exists in the database.\n")
64 }
65 if (!record_exists) {
66   # Check if the primary key value of the new record is unique in the temporary dataframe
67   if (!primary_key_value %in% promotion_possible_data[[promotion_primary_key]]) {
68     promotion_possible_data <- rbind(promotion_possible_data, new_record)
69   }
70 }
71
72 cat("Finished processing file:", file_path, "\n")
73
74 }
75 }
76 cat("Starting validation for new records.\n")
77 promotion_possible_data <- validate_and_prepare_promotion_data(promotion_possible_data)
78 cat("Validation completed for new records.\n")
79
80

```

```

81 if (nrow(promotion_possible_data) > 0)
82 {
83   cat("Starting to insert validated data into the database. Number of records: ", nrow(promotio
84 # Digesting prepared data to our database
85 dbWriteTable(connection, name = "promotion", value = promotion_possible_data, append = TRUE,
86 cat("Data insertion completed successfully.\n")
87 } else
88 {
89   cat("No valid promotion data to insert into the database.\n")
90 }

```

5. Validations for shipment data

```

1 validate_and_prepare_shipment_data <- function(data) {
2   # # Validation for shipment ID
3   shipment_id_check <- grepl("^SHIP[0-9]{6}$", data$shipment_id)
4   data <- data[shipment_id_check,]
5
6   # Convert dates from character to Date object
7   data$shipment_date <- as.Date(data$shipment_date, format = "%d-%m-%Y")
8   data$delivery_date <- as.Date(data$delivery_date, format = "%d-%m-%Y")
9
10  # Validation for shipment_date and delivery_date format
11  date_format_check <- !is.na(data$shipment_date) & !is.na(data$delivery_date)
12
13  # Keep only rows with valid date formats
14  data <- data[date_format_check,]
15
16  # Validation for logical order of shipment and delivery dates
17  logical_date_order_check <- data$shipment_date <= data$delivery_date
18
19  # Keep only rows with logical date order
20  data <- data[logical_date_order_check,]
21
22  return(data)
23}
24
25 # Fetch existing shipment IDs from the database
26 existing_shipment_ids <- dbGetQuery(connection, "SELECT shipment_id FROM shipment")$shipme
27 # List all files that have shipment in the title
28 shipment_file_paths <- list.files(path = "data_upload", pattern = "shipment.*\\.csv$", ful
29 # Destination directory to move after processing the file

```

```

30 destination_dir <- "data_processed"
31 # Ensure the destination directory exists before processing files
32 if (!dir.exists(destination_dir)) {
33   dir.create(destination_dir)
34 }
35
36 # Check if there are new files for shipment table
37 if (length(shipment_file_paths) > 0) {
38   #Initialising empty data frame
39   shipment_possible_data <- data.frame(shipment_id = character(), stringsAsFactors = FALSE)
40   # Read each shipment CSV file and check for the existence of the primary key in the data
41   for (file_path in shipment_file_paths) {
42     cat("Starting processing file:", file_path, "\n")
43
44     # Read the current file
45     shipment_data <- readr::read_csv(file_path, show_col_types = FALSE)
46
47     # Filter out records with existing shipment_id in the database
48     unique_shipment_data <- shipment_data[!shipment_data$shipment_id %in% existing_shipment_ids]
49
50     # Combine unique records into the possible data frame
51     shipment_possible_data <- rbind(shipment_possible_data, unique_shipment_data)
52
53     # Construct the destination file path
54     file_name <- basename(file_path)
55     dest_file_path <- file.path(destination_dir, file_name)
56
57     # Construct the destination file path with a timestamp to ensure uniqueness
58     file_name <- basename(file_path)
59     timestamp <- format(Sys.time(), "%Y%m%d%H%M%S") # Get current timestamp
60     file_ext <- tools::file_ext(file_name) # Extract file extension
61     base_name <- sub(paste0("\\.", file_ext, "$"), "", file_name) # Remove extension from file name
62     new_file_name <- paste0(base_name, "_", timestamp, ".", file_ext) # Create new filename
63     dest_file_path <- file.path(destination_dir, new_file_name) # Construct new destination file path
64
65     # Move the processed file to the destination directory
66     if (file.rename(file_path, dest_file_path)) {
67       cat("Successfully moved processed file to:", dest_file_path, "\n")
68     } else {
69       cat("Failed to move file:", file_path, "\n")
70     }
71     cat("Finished processing file:", file_path, "\n")

```

```

72 }
73
74 # Ensure shipment_possible_data contains only unique shipment_ids, no repetition of new
75 shipment_possible_data <- shipment_possible_data %>% distinct(shipment_id, .keep_all = TRUE)
76
77
78 if (nrow(shipment_possible_data) > 0) {
79   cat("Starting validation for new records.\n")
80   shipment_possible_data <- validate_and_prepare_shipment_data(shipment_possible_data)
81   cat("Validation completed for new records.\n")
82 }
83
84 # Order data to ensure consistent hashing
85 shipment_possible_data <- shipment_possible_data[order(shipment_possible_data$shipment_id)]
86
87 if (!is.null(dim(shipment_possible_data)) && dim(shipment_possible_data)[1] > 0) {
88   # Generate pre-load hashes
89   pre_load_hashes <- sapply(1:nrow(shipment_possible_data), function(i) {
90     record <- as.character(unlist(shipment_possible_data[i, ]))
91     digest(paste(record, collapse = "|"), algo = "md5")
92   })
93   cat("Starting to insert validated data into the database. Number of records: ", nrow(shipment_possible_data))
94   tryCatch({
95     dbWriteTable(connection, name = "shipment", value = shipment_possible_data, append = TRUE)
96     cat("Data insertion completed successfully.\n")
97   }, error = function(e) {
98     cat("Error inserting data into the database: ", e$message, "\n")
99   })
100  # Fetch the loaded data back for post-load hash comparison
101  loaded_shipment_ids <- sprintf("%s", shipment_possible_data$shipment_id)
102  query <- sprintf("SELECT * FROM shipment WHERE shipment_id IN (%s) ORDER BY shipment_id", loaded_shipment_ids)
103
104  retrieved_shipments <- dbGetQuery(connection, query)
105
106  post_load_hashes <- sapply(1:nrow(retrieved_shipments), function(i) {
107    record <- as.character(unlist(retrieved_shipments[i, ]))
108    digest(paste(record, collapse = "|"), algo = "md5")
109  })
110
111  if (nrow(retrieved_shipments) > 0) {
112
113    # Compare hashes

```

```

114     identical_hashes <- all(pre_load_hashes == post_load_hashes)
115     if (identical_hashes) {
116       cat("Data integrity verified: All record hashes match.\n")
117     } else {
118       cat("Data integrity check failed: Record hashes do not match.\n")
119     }
120   }
121 } else {
122   cat("No valid shipment data to insert into the database.\n")
123 }
124 } else {
125   cat("No new files to process.\n")
126 }
```

6. Validations for product data

```

1 # function to validate and check referential integrity
2 validate_and_prepare_product_data <- function(data, connection) {
3
4   # Validation for product ID
5   product_id_check <- grepl("^[A-Za-z0-9]{10}$", data$product_id)
6   data <- data[product_id_check, ]
7   # Performing validation for review score
8   data <- data[data$review_score >= 1 & data$review_score <= 5, ]
9
10  # Fetch existing IDs from reference tables
11  valid_category_ids <- dbGetQuery(connection, "SELECT category_id FROM category")$category_id
12  valid_supplier_ids <- dbGetQuery(connection, "SELECT supplier_id FROM supplier")$supplier_id
13  valid_promotion_ids <- c(dbGetQuery(connection, "SELECT promotion_id FROM promotion")$promotion_id)
14
15  # Referential integrity checks
16  data <- data[data$category_id %in% valid_category_ids, ]
17  data <- data[data$supplier_id %in% valid_supplier_ids, ]
18  data <- data[is.na(data$promotion_id) | data$promotion_id %in% valid_promotion_ids, ]
19
20  # Convert quantity_stock and quantity_supplied to numeric to ensure
21  # non-numeric values are handled. This introduces NA for any non-numeric values.
22  data$quantity_stock <- as.numeric(data$quantity_stock)
23  data$quantity_supplied <- as.numeric(data$quantity_supplied)
24
25  # Now proceed with the validations for non-negative, non-NA, and integer for quantity_stock
26  data <- data[!is.na(data$quantity_stock) & data$quantity_stock >= 0 & (data$quantity_stock %in% valid_quantity_stock)]
```

```

27 data <- data[!is.na(data$quantity_supplied) & data$quantity_supplied >= 0 & (data$quantity_...
28
29 # Validation for positive price values
30 data <- data[data$price > 0 & !is.na(data$price), ]
31
32 # Validation to ensure product_name is not empty
33 data <- data[data$product_name != "" & !is.na(data$product_name), ]
34
35 return(data)
36 }
37
38
39 # Fetch existing product IDs from the database into a vector
40 existing_product_ids <- tryCatch({
41   dbGetQuery(connection, "SELECT product_id FROM product")$product_id
42 }, error = function(e) {
43   cat("Error fetching existing product IDs: ", e$message, "\n")
44   NULL
45 })
46
47 # List files that have 'product' in title
48 product_file_paths <- list.files(path = "data_upload", pattern = "product.*\\.csv$", full.names = ...
49 # Initialising temporary empty data frame for future loading
50 product_possible_data <- data.frame()
51
52 # Read each product CSV file
53 for (file_path in product_file_paths) {
54   cat("Starting processing file:", file_path, "\n")
55   # Read the current file
56   product_data <- tryCatch({
57     readr::read_csv(file_path)
58   }, error = function(e) {
59     cat("Error reading file", file_path, ":", e$message, "\n")
60     next # Skip to the next iteration of the loop
61   })
62
63 # Iterate through each row of the file and check for uniqueness of primary key in database
64 for (i in seq_len(nrow(product_data))) {
65   new_record <- product_data[i, ]
66   primary_key_value <- new_record[["product_id"]]
67
68   if (!primary_key_value %in% existing_product_ids) {

```

```

69     if (!primary_key_value %in% product_possible_data[["product_id"]]) {
70         product_possible_data <- rbind(product_possible_data, new_record)
71     }
72 } else {
73     cat("Record with primary key", primary_key_value, "already exists in the database.\n")
74 }
75 }
76 cat("Finished processing file:", file_path, "\n")
77 }

78 # Call for validation and integrity check function
79 cat("Starting validation for new records.\n")
80 product_possible_data <- validate_and_prepare_product_data(product_possible_data, connection)
81 cat("Validation completed for new records.\n")

82 # Implementing data integrity check for validated and ready to loading data for each row
83 if ("product_id" %in% names(product_possible_data) && nrow(product_possible_data) > 0) {
84     product_possible_data <- product_possible_data[order(product_possible_data$product_id), ]
85 } else {
86     cat("product_id column not found or product_possible_data is empty.\n")
87 }
88 pre_load_hashes <- sapply(1:nrow(product_possible_data), function(i) {
89     record <- as.character(unlist(product_possible_data[i, ]))
90     digest(paste(record, collapse = "|"), algo = "md5")
91 })
92

93 # Ingesting into database command
94 if (nrow(product_possible_data) > 0) {
95     cat("Starting to insert validated data into the database. Number of records: ", nrow(product_
96
97     # Ingesting prepared data to our database
98     tryCatch({
99         dbWriteTable(connection, name = "product", value = product_possible_data, append = TRUE,
100         cat("Data insertion completed successfully.\n")
101     }, error = function(e) {
102         cat("Error inserting data into the database: ", e$message, "\n")
103
104         # Additional error handling logic here
105     })
106
107 }
108
109 # Fetch the loaded data back for post-load hash comparison
110 loaded_product_ids <- sprintf("'%s'", product_possible_data$product_id)

```

```

111 query <- sprintf("SELECT * FROM product WHERE product_id IN (%s) ORDER BY product_id", paste0(retrieved_products))
112 retrieved_products <- dbGetQuery(connection, query)
113
114 # Creating hashes for retrieved data from db
115 post_load_hashes <- sapply(1:nrow(retrieved_products), function(i) {
116   record <- as.character(unlist(retrieved_products[i, ]))
117   digest(paste(record, collapse = "|"), algo = "md5")
118 })
119
120 # Compare hashes for pre-loaded and retrieved-loaded data
121 identical_hashes_product <- all(pre_load_hashes == post_load_hashes)
122 if (identical_hashes_product) {
123   cat("Data integrity verified: All record hashes match.\n")
124 } else {
125   cat("Data integrity check failed: Record hashes do not match.\n")
126 }
127 } else {
128   cat("No valid product data to insert into the database.\n")
129 }
```

7. Validation for orders data

```

1 #Validation for orders data
2
3 # Function to check if given foreign IDs of orders data exist in their respective tables
4 check_id_exists <- function(connection, table_name, column_name, ids) {
5   query_template <- "SELECT DISTINCT %s FROM %s WHERE %s IN (%s)"
6   query <- sprintf(query_template, column_name, table_name, column_name, paste0("'", ids,
7   existing_ids <- dbGetQuery(connection, query)[[1]]
8   all(ids %in% existing_ids)
9 }
10
11 # Function to validate that foreign keys of 'orders' are existing in respective tables
12 validate_referential_integrity <- function(new_record) {
13   customer_exists <- check_id_exists(connection, "customer", "customer_id", new_record$customer_id)
14   product_exists <- check_id_exists(connection, "product", "product_id", new_record$product_id)
15   shipment_exists <- check_id_exists(connection, "shipment", "shipment_id", new_record$shipment_id)
16
17   if(product_exists & customer_exists & shipment_exists) {
18     return(TRUE)
19   } else {
20     if(!customer_exists) cat("Customer ID does not exist:", new_record$customer_id, "\n")
```

```

21     if(!product_exists) cat("Product ID does not exist:", new_record$product_id, "\n")
22     if(!shipment_exists) cat("Shipment ID does not exist:", new_record$shipment_id, "\n")
23     return(FALSE)
24   }
25 }

26

27
28 # Referential integrity of orders table
29 # Function to check if given foreign IDs of orders data exist in their respective tables
30 check_id_exists <- function(connection, table_name, column_name, ids) {
31   query_template <- "SELECT DISTINCT %s FROM %s WHERE %s IN (%s)"
32   query <- sprintf(query_template, column_name, table_name, column_name, paste0("'", ids,
33   existing_ids <- dbGetQuery(connection, query)[[1]]
34   all(ids %in% existing_ids)
35 }

36

37 # Function to validate that foreign keys of 'orders' are existing in respective tables
38 validate_referential_integrity <- function(new_record) {
39   customer_exists <- check_id_exists(connection, "customer", "customer_id", new_record$cu
40   product_exists <- check_id_exists(connection, "product", "product_id", new_record$produ
41   shipment_exists <- check_id_exists(connection, "shipment", "shipment_id", new_record$sh
42
43   if(product_exists & customer_exists & shipment_exists) {
44     return(TRUE)
45   } else {
46     if(!customer_exists) cat("Customer ID does not exist:", new_record$customer_id, "\n")
47     if(!product_exists) cat("Product ID does not exist:", new_record$product_id, "\n")
48     if(!shipment_exists) cat("Shipment ID does not exist:", new_record$shipment_id, "\n")
49     return(FALSE)
50   }
51 }

52

53

54
55 # Validation function of orders table
56
57 validate_and_prepare_orders_data <- function(data){
58   # Checking format of order id
59   order_id_check <- grepl("^ORDER[0-9]{9}$", data$order_id)
60   data <- data[order_id_check, ]
61
62   customer_id_check <- grepl("^CUST[0-9]{6}$", data$customer_id)

```

```

63   data <- data[customer_id_check, ]
64
65   product_id_check <- grepl("^[A-Za-z0-9]{10}$", data$product_id)
66   data <- data[product_id_check, ]
67
68   shipment_id_check <- grepl("^SHIP[0-9]{6}$", data$shipment_id)
69   data <- data[shipment_id_check, ]
70
71   # Convert order_date strings to Date objects
72   converted_dates <- as.Date(data$order_date, format = "%d-%m-%Y")
73
74   # Efficient validation for order_date format and range
75   valid_dates <- !is.na(converted_dates) & converted_dates >= as.Date("01-01-2023")
76
77   # Filter data based on the valid_dates check
78   data <- data[valid_dates, ]
79
80   # Validation for order_date format
81   order_date_format_check <- !is.na(as.Date(data$order_date, format = "%d-%m-%Y"))
82   data <- data[order_date_format_check, ]
83
84
85   return(data)
86 }
87
88 # List all files that have 'orders' in name
89 orders_file_paths <- list.files(path = "data_upload", pattern = "orders.*\\.csv$", full.names = TRUE)
90 # Create a temporary df for data to possibly upload
91 orders_possible_data <- data.frame()
92 # Set a data_processed directory for transferring the file
93 destination_dir <- "data_processed"
94
95 # Ensure the destination directory exists
96 if (!dir.exists(destination_dir)) {
97   dir.create(destination_dir)
98 }
99
100 # Check if there are new files for the orders table
101 if (length(orders_file_paths) > 0) {
102   # Read each orders CSV file and check for the existence of the composite primary key in
103   for (file_path in orders_file_paths) {
104     orders_data <- readr::read_csv(file_path, show_col_types = FALSE)

```

```

105 cat("Starting processing file:", file_path, "\n")
106 # Iterate through each row of the file
107 for (i in seq_len(nrow(orders_data))) {
108   new_record <- orders_data[i, ]
109
110   # Check for the Referential integrity
111   if(validate_referential_integrity(new_record)) {
112     # Construct the condition to check the composite primary key (order_id, product_id)
113     conditions <- sprintf("order_id = '%s' AND product_id = '%s' AND customer_id = '%s'",
114                           new_record$order_id, new_record$product_id, new_record$customer_id)
115
116     # Check if a record with the same composite primary key exists in the database
117     record_exists_query <- paste("SELECT COUNT(*) FROM orders WHERE", conditions)
118     record_exists_result <- dbGetQuery(connection, record_exists_query)
119     record_exists <- record_exists_result[1, 1] > 0
120
121
122   if(!record_exists) {
123     # Construct a unique identifier for the composite primary key for potential new records
124     composite_key <- paste(new_record$order_id, new_record$product_id, new_record$customer_id)
125
126     # Function to check if the composite primary key is unique in the temporary data frame
127     existing_keys <- sapply(1:nrow(orders_possible_data), function(i) {
128       paste(orders_possible_data[i, "order_id"], orders_possible_data[i, "product_id"],
129             orders_possible_data[i, "shipment_id"], sep = "-"))
130     })
131
132     if (!composite_key %in% existing_keys) {
133       # If new potential row has unique primary key then move it to temporary df
134       orders_possible_data <- rbind(orders_possible_data, new_record)
135     } else {
136       cat("Record with composite primary key already exists in temporary data.\n")
137     }
138   } else {
139     cat("Record with composite primary key already exists in the database.\n")
140   }
141 } else {
142   cat("Referential integrity for the row is invalid.\n")
143 }
144
145 cat("Starting validation for checked new records.\n")
146 orders_possible_data <- validate_and_prepare_orders_data(orders_possible_data)

```

```

147 cat("Validation completed for new records.\n")
148
149 # Sorting the data to maintain ordering
150 if (nrow(orders_possible_data) > 0) {
151   orders_possible_data <- orders_possible_data[order(orders_possible_data$order_id, or
152                                                 orders_possible_data$shipment_id)]
153
154 # Hashing the ready to upload data for data integrity check
155 pre_load_hashes <- sapply(1:nrow(orders_possible_data), function(i) {
156   record <- as.character(unlist(orders_possible_data[i, ]))
157   digest(paste(record, collapse = "|"), algo = "md5")
158 })
159 }
160
161 # Orders ingestion
162 if (nrow(orders_possible_data) > 0) {
163   cat("Starting to insert validated data into the database. Number of records: ", nro
164   # Attempt to ingest prepared data to our database
165   tryCatch({
166     dbWriteTable(connection, name = "orders", value = orders_possible_data, append = T
167     cat("Data insertion completed successfully.\n")
168   }, error = function(e) {
169     cat("Error inserting data into the database: ", e$message, "\n")
170   })
171   # Fetch the loaded data back for post-load hash comparison
172   loaded_order_ids <- sprintf("%s", orders_possible_data$order_id)
173   query <- sprintf("SELECT * FROM orders WHERE order_id IN (%s) ORDER BY order_id, pr
174   retrieved_orders <- dbGetQuery(connection, query)
175   # Creating hashes for retrieved data from db
176   post_load_hashes <- sapply(1:nrow(retrieved_orders), function(i) {
177     record <- as.character(unlist(retrieved_orders[i, ]))
178     digest(paste(record, collapse = "|"), algo = "md5")
179   })
180
181   # Compare hashes for pre-loaded and retrieved-loaded data
182   identical_hashes <- all(pre_load_hashes == post_load_hashes)
183   if (identical_hashes) {
184     cat("Data integrity verified: All record hashes match.\n")
185   } else {
186     cat("Data integrity check failed: Record hashes do not match.\n")
187   }
188 } else {

```

```

189     cat("No valid orders data to insert into the database.\n")
190 }
191 # Move the processed file to 'data_processed' directory with a timestamp
192 file_name <- basename(file_path)
193 timestamp <- format(Sys.time(), "%Y%m%d%H%M%S")
194 file_ext <- tools::file_ext(file_name)
195 base_name <- sub(paste0("\\.", file_ext, "$"), "", file_name)
196 new_file_name <- paste0(base_name, "_", timestamp, ".", file_ext)
197 dest_file_path <- file.path(destination_dir, new_file_name)

198 if (file.rename(file_path, dest_file_path)) {
199   cat("Successfully moved processed file to:", dest_file_path, "\n")
200 } else {
201   cat("Failed to move file:", file_path, "\n")
202 }
203 cat("Finished processing file:", file_path, "\n")
204 }
205 } else {
206   # If there are no new files, skip processing and notify the user
207   cat("No new files to process.\n") }
```

Referential Integrity for ensuring data integrity

Foreign key is the crucial component of a database that enforces referential integrity, ensuring that a value appearing in one entity should also exist in another entity as a foreign key. Basically, it ensures that a value referenced in one table exists in another table, maintaining the integrity and consistency of the data.

In the ‘category’ table, ‘category_id’ serves as the primary key, that uniquely identifies each category of the product. On the other hand, in products table ‘category_id’ serves as the foreign key. It means if one of the category ids is removed from the category table then it should be removed from the product table as well.

Part 3 Data Pipeline Generation

3.1 GitHub Repository and Workflow Setup

In this part, a GitHub repository named “DM_Group_18” was created and connected it to the Posit cloud. It also acted as the central hub of the project which helped the team members to collaborate effectively. In the repository the CSV files were uploaded in the folder named

“data_upload”. Also, the R scripts for database schema creation, validation and analysis could be found in the folder named “R”.

The URL of the repository - https://github.com/AkarshaShrivastava19/DM_group_18

The screenshot shows a GitHub repository page for 'DM_group_18'. At the top, there's a navigation bar with links for Code, Issues, Pull requests, Actions, Projects, Security, Insights, and Settings. Below the navigation bar, the repository name 'DM_group_18' is displayed, along with a 'Public' badge. A search bar and a 'Type' input field are also present. The main content area shows a list of commits in the 'main' branch. The commits are as follows:

Author	Commit Message	Time Ago
AkarshaShrivastava19	Update ecomdata.db	5b09019 · 6 minutes ago
.github/workflows	Update etl.yaml	14 hours ago
R	Update the Analysis R Script	11 minutes ago
data_upload	Update Shipment Dataset	44 minutes ago
rdadata	Committing the database	4 days ago
.gitignore	Initial commit	3 weeks ago
README.md	Initial commit	3 weeks ago
Report.qmd	update the analysis file and the report file	2 days ago

On the right side of the repository page, there are sections for 'About' (with a note: 'No description, website, or topics provided.'), 'Readme', 'Activity' (0 stars, 1 watching, 0 forks), 'Releases' (No releases published, Create a new release), and 'Packages'.

Figure 12: Created GitHub Repository

3.2 GitHub Actions for Continuous Integration

Here, GitHub actions were implemented to automate different stages of data pipeline such as database updates, data validation and data analysis. The workflow was configured to trigger each time new changes were pushed to the main branch from Posit cloud ensuring that the automated tasks were executed in response to the relevant changes.

It runs multiple jobs seamlessly on the latest Ubuntu Environment as shown below.

The screenshot shows a build summary for a workflow named 'etl.yaml #67'. The build was triggered 'yesterday' and succeeded in 4m 49s. The build steps are listed on the right, each with a status icon and duration:

- > Set up job: 10s
- > Checkout code: 0s
- > Setup R environment: 40s
- > Cache R packages: 2s
- > Install packages: 3m 41s
- > Execute R script for Initial table creation and validation: 3s
- > Execute R script for automation analysis: 5s
- > Add files: 0s
- > Check if ecomdata.db has changed: 0s
- < Commit and push if ecomdata.db has changed: 0s
- < Push changes: 0s
- > Post Cache R packages: 0s
- > Post Checkout code: 0s
- > Complete job: 0s

So, after the detection of any changes the workflow activates in response to any “pull” and “push” request. We configured to automatically update the database, run data validations and analysis.

By implementing this, we significantly reduced manual intervention required for this process thus increasing the efficiency and rapidly detecting the issues.

The screenshot shows a summary of a build triggered via push 14 hours ago. The build was successful with a total duration of 4m 56s. The build step 'build' completed successfully in 4m 49s.

Triggered via push 14 hours ago	Status	Total duration	Artifacts
Ilyas1210K pushed -> b2218c9 main	Success	4m 56s	-

etl.yaml
on: push

build 4m 49s

Figure 13: Successfully Build Workflow

The yaml code is added below:

```

name: ETL Workflow Group 18

on: # schedule: # - cron: '0 /3 * *' # Run every 3 hours push: branches: [ main ]

jobs: build: runs-on: ubuntu-latest
steps:
  - name: Checkout code
    uses: actions/checkout@v2

  - name: Setup R environment
    uses: r-lib/actions/setup-r@v2
    with:
      r-version: '4.2.0'

  - name: Cache R packages
    uses: actions/cache@v2
    with:
      path: ${{ env.R_LIBS_USER }}
      key: ${{ runner.os }}-r-${{ hashFiles('**/lockfile') }}
      restore-keys: ${{ runner.os }}-r-

  - name: Install packages
    if: steps.cache.outputs.cache-hit != 'true'
    run: |
      Rscript -e
      'install.packages(c("ggplot2","dplyr","readr","RSQLite","DBI","lubridate","tidyverse",
      "knitr", "digest"))'

  - name: Execute R script for Initial table creation and validation
    run: |
      Rscript R/data_workflow.R

  - name: Execute R script for automation analysis
    run: |
      Rscript R/Analysis.R

  - name: Add files
    run: |
      git config --global user.email "akarsha.shrivastava@warwick.ac.uk"
      git config --global user.name "AkarshaShrivastava19"

  - name: Check if ecomdata.db has changed
    id: check-db
    run: |
      git diff --exit-code --name-only | grep 'ecomdata.db' && echo "::set-output
      name=changed::true" || echo "::set-output name=changed::false"

  - name: Commit and push if ecomdata.db has changed
    if: steps.check-db.outputs.changed == 'true'
    run: |
      git add ecomdata.db
      git commit -m "Update ecomdata.db"
      git push

  - name: Push changes
    if: steps.check-db.outputs.changed == 'true'
    uses: ad-m/github-push-action@v0.6.0
    with:
      github_token: ${{ secrets.GITHUB_TOKEN }}
      branch: main

```

Part 4 Data Analysis and Reporting

The following data analysis was performed on the generated e-commerce data.

4.1 Advanced Data Analysis in R

```
1 # Retrieve data from the database
2 customer <- dbGetQuery(connection, "SELECT * FROM customer")
3 product <- dbGetQuery(connection, "SELECT * FROM product")
4 supplier <- dbGetQuery(connection, "SELECT * FROM supplier")
5 category <- dbGetQuery(connection, "SELECT * FROM category")
6 shipment <- dbGetQuery(connection, "SELECT * FROM shipment")
7 promotion <- dbGetQuery(connection, "SELECT * FROM promotion")
8 orders <- dbGetQuery(connection, "SELECT * FROM orders")
```

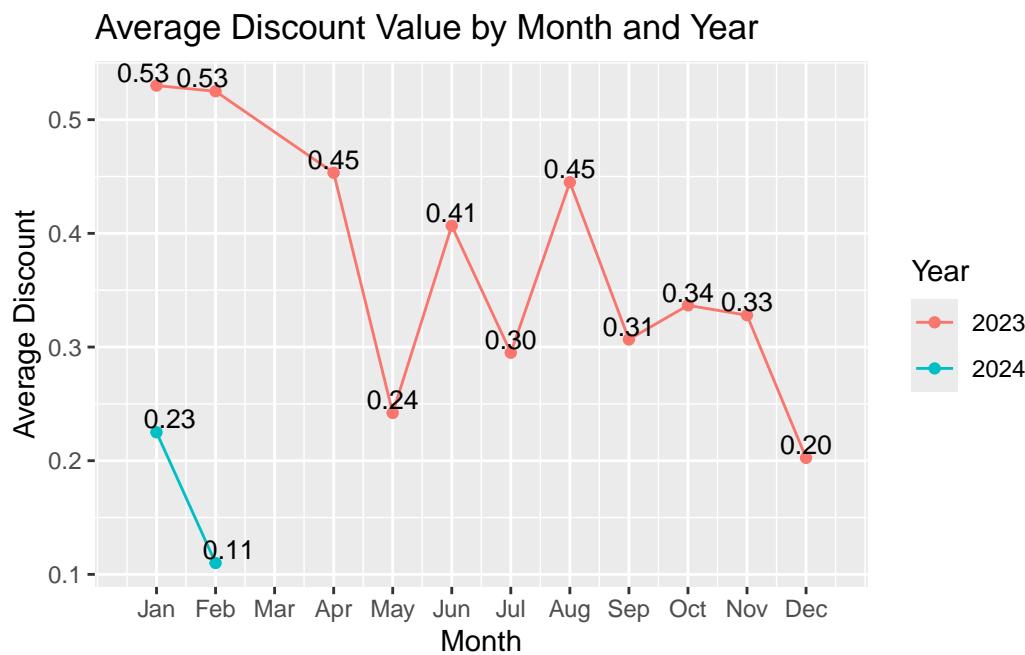
4.1.1 Promotion Discount Trend

```
1 # Convert Date Format
2 promotion <- promotion %>%
3   mutate(promotion_start_date = dmy(promotion_start_date),
4         promotion_end_date = dmy(promotion_end_date))
5
6 # Generate records for each month that each promotion spans
7 data_expanded <- promotion %>%
8   rowwise() %>%
9   mutate(months = list(seq(from = promotion_start_date,
10                      to = promotion_end_date,
11                      by = "month")))) %>%
12   unnest(months) %>%
13   mutate(year = year(months), month = month(months)) %>%
14   group_by(promotion_id, year, month) %>%
15   summarise(promotion_discount_value = mean(promotion_discount_value), .groups = 'drop')
16
17 # Calculate the average discount value for each month
18 average_discounts <- data_expanded %>%
19   group_by(year, month) %>%
20   summarise(average_discount = mean(promotion_discount_value))
21
22 # Specify the dimensions of the plot
```

```

23 width <- 12
24 height <- 8
25
26 # Visualise the average discount value for different months and years
27 g_promotionvalue <- ggplot(average_discounts, aes(x = month, y = average_discount, group = year))
28   geom_line() +
29   geom_point() +
30   scale_x_continuous(breaks = 1:12, labels = month.abb) +
31   geom_text(aes(label = sprintf("%.2f", average_discount)), position = position_dodge(width = 0.9))
32   labs(title = "Average Discount Value by Month and Year",
33       x = "Month",
34       y = "Average Discount",
35       color = "Year")
36 print(g_promotionvalue)

```



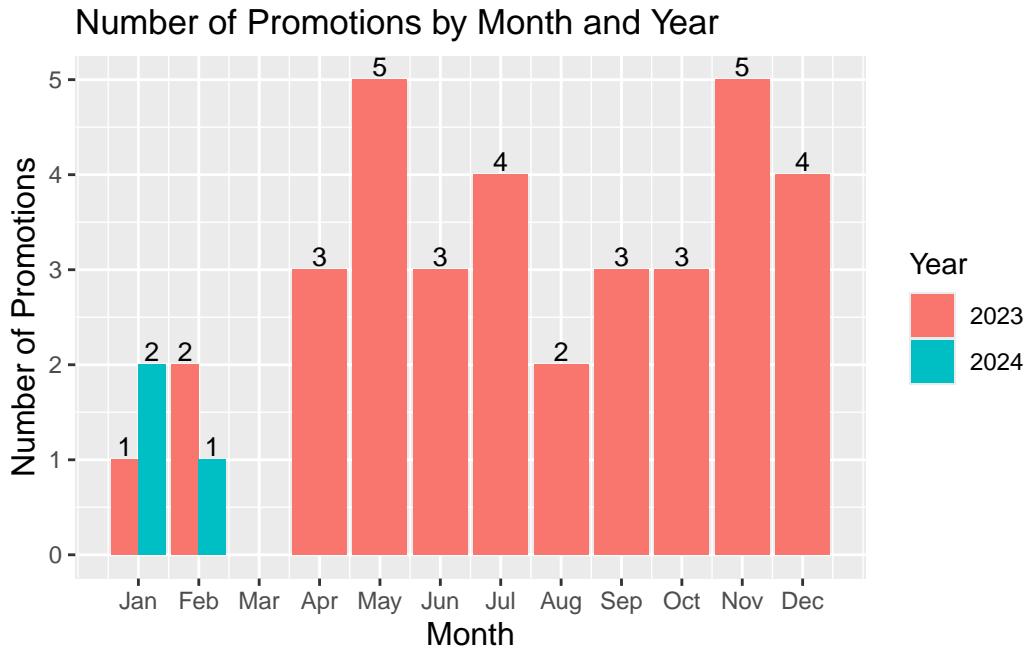
```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("promotion_discount_trend_", format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
3
4 # Save the plot with the dynamic filename
5 ggsave(filename, plot = g_promotionvalue, width = width, height = height)

```

4.1.2 Promotion Count Trend

```
1 # Calculate the number of times a promotion appears in each month
2 promotion_counts <- data_expanded %>%
3   count(year, month)
4
5 # Visualise the number of times promotions appear in different years and months
6 g_promotioncount <- ggplot(promotion_counts, aes(x = month, y = n, fill = as.factor(year))) +
7   geom_bar(stat = "identity", position = "dodge") +
8   scale_x_continuous(breaks = 1:12, labels = month.abb) +
9   geom_text(aes(label = n), position = position_dodge(width = 0.9), vjust = -0.2, size = 3.5)
10  theme(axis.title = element_text(size = 12)) +
11  labs(title = "Number of Promotions by Month and Year",
12    x = "Month",
13    y = "Number of Promotions",
14    fill = "Year")
15 print(g_promotioncount)
```



```
1 # Dynamically generate filename with current date and time
2 filename <- paste0("promotion_number_trend_", format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
3
4 # Save the plot with the dynamic filename
5 ggsave(filename, plot = g_promotioncount, width = width, height = height)
```

Promotion Analysis

The above graphs illustrate the highest discounts that were offered during January and February, i.e. 53% off on the product price.

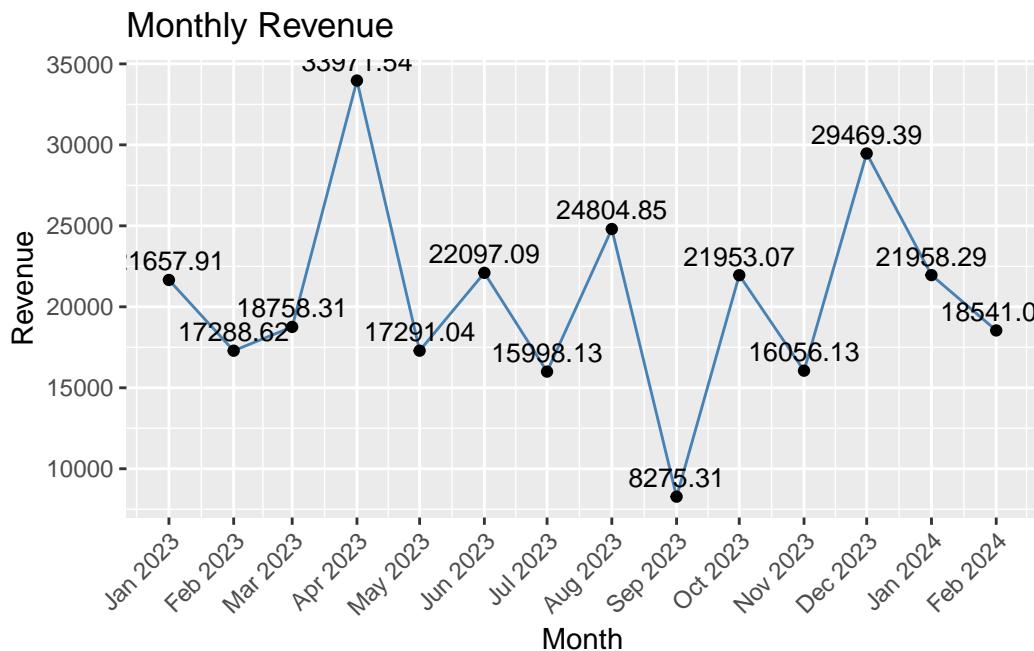
4.1.3 Monthly Revenue Trend

```
1 # Preprocessed date formats
2 orders <- orders %>% mutate(order_date = dmy(order_date))
3 promotion <- promotion %>%
4   mutate(start_date = promotion_start_date,
5         end_date = promotion_end_date,
6         promotion_discount_value = if_else(is.na(promotion_discount_value), 0, promotion_di
7
8 # Merge orders with products for pricing information
9 order_products <- orders %>%
10   left_join(product, by = "product_id")
11
12 # Make sure there are no missing prices or quantities
13 order_products <- order_products %>%
14   mutate(price = if_else(is.na(price), 0, price),
15         quantity = if_else(is.na(quantity), 0, quantity))
16
17 # Combine orders, products and promotions to take into account discounts during promotions
18 order_products_promotions <- order_products %>%
19   left_join(promotion, by = "promotion_id") %>%
20   mutate(is_promotion = if_else(order_date >= start_date & order_date <= end_date, TRUE, FALSE),
21         revenue = price * quantity * if_else(is_promotion, 1 - promotion_discount_value, 1))
22
23 # Remove any missing income values generated in the calculation
24 order_products_promotions <- order_products_promotions %>%
25   filter(!is.na(revenue))
26
27 # Calculation of gross monthly income
28 monthly_revenue <- order_products_promotions %>%
29   mutate(month = floor_date(order_date, "month")) %>%
30   group_by(month) %>%
31   summarize(total_revenue = sum(revenue, na.rm = TRUE))
32
33 # Visualisation of monthly income
34 (g_monthlyrevenue <- ggplot(monthly_revenue, aes(x = month, y = total_revenue)) +
```

```

35 geom_line(color = "steelblue") +
36 geom_point() +
37 scale_x_date(date_labels = "%b %Y", date_breaks = "1 month") +
38 geom_text(aes(label = sprintf("%.2f", total_revenue)), position = position_dodge(width = 0
39 theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
40 labs(title = "Monthly Revenue", x = "Month", y = "Revenue"))

```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("monthly_revenue_",
3                     format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
4
5 # Save the plot with the dynamic filename
6 ggsave(filename, plot = g_monthlyrevenue, width = width, height = height)

```

Monthly Revenue Analysis

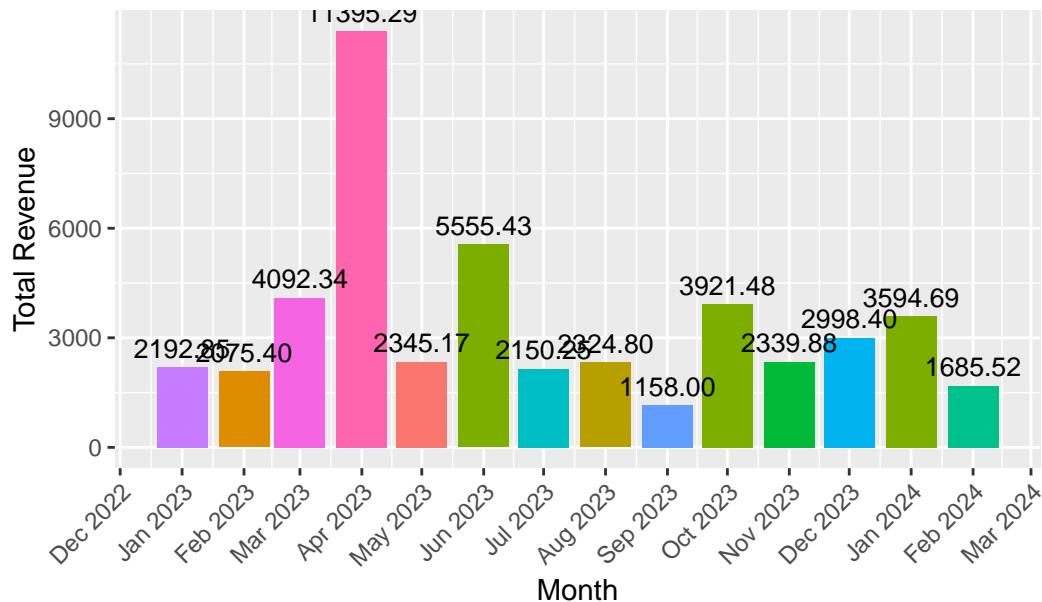
It is observed that the month of April in 2023 records the highest revenue, credited to the promotion event with a high discount rate, followed by revenue in December 2023. Although no attractive discount offers were applied in the month of December 2023, still an outstanding revenue was observed because of the festive season in the United Kingdom.

4.1.4 Monthly Best-Selling Products

Graph for Monthly Best-Selling Products

```
1 # Calculate total monthly revenue per product
2 monthly_product_revenue <- order_products_promotions %>%
3   mutate(month = floor_date(order_date, "month")) %>%
4   group_by(month, product_id) %>%
5   summarize(total_revenue_product = sum(revenue, na.rm = TRUE))
6
7 # Select top earning products per month
8 best_selling_products_each_month <- monthly_product_revenue %>%
9   group_by(month) %>%
10  slice_max(total_revenue_product, n = 1) %>%
11  ungroup() %>%
12  select(month, product_id, total_revenue_product)
13
14 best_selling_products_each_month <- merge(best_selling_products_each_month, product[, c("prod
15
16 # Visualise the top earning products and their revenues per month
17 g_bestseller_product_monthly <- ggplot(best_selling_products_each_month, aes(x = month, y =
18   geom_col(show.legend = FALSE) +
19   geom_text(aes(label = sprintf("%.2f", total_revenue_product)), position = position_dodge(w
20   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
21   scale_x_date(date_labels = "%b %Y", date_breaks = "1 month") +
22   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
23   labs(title = "Best-Selling Products by Month", x = "Month", y = "Total Revenue")
24
25 print(g_bestseller_product_monthly)
```

Best-Selling Products by Month



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("monthly_bestseller_product_",
3                     format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
4
5 # Save the plot with the dynamic filename
6 ggsave(filename, plot = g_bestseller_product_monthly, width = width, height = height)

```

Table for Monthly Best-Selling Products

```

1 best_selling_products_each_month$month <- as.Date(best_selling_products_each_month$month, "%"
2
3 best_selling_products_each_month$YearMonth <- format(best_selling_products_each_month$month,
4
5 best_selling_products_each_month <- best_selling_products_each_month %>%
6   arrange(YearMonth)
7
8 table_to_display <- best_selling_products_each_month %>%
9   select(YearMonth, product_name, total_revenue_product) %>%
10  rename('Total Revenue' = total_revenue_product)
11
12 # Display the table with kable
13 kable(table_to_display, caption = "Monthly Best-Selling Products", col.names = c("Time", "Pr

```

Table 1: Monthly Best-Selling Products

Time	Product Name	Total Revenue
2023-01	Wall Paint Behr Premium Plus Ultra	2192.85
2023-02	Garden Tool Corona Bypass Pruner	2075.40
2023-03	Wall Paint Kilz Complete Coat	4092.34
2023-04	Wall Paint Olympic One	11395.29
2023-05	Facial Cleanser Cetaphil Gentle Skin Cleanser	2345.17
2023-06	Laptop Acer Swift 5	5555.43
2023-07	Pasta Barilla Linguine	2150.25
2023-08	Guitar Epiphone SG Standard	2324.80
2023-09	Stroller Baby Jogger City Mini GT2	1158.00
2023-10	Laptop Acer Swift 5	3921.48
2023-11	Office Chair Boss Office Products Executive	2339.88
2023-12	Refrigerator Samsung RS27T5561SR	2998.40
2024-01	Laptop Acer Swift 5	3594.69
2024-02	Office Chair Serta Big and Tall	1685.52

Monthly Best-Selling Products Analysis

It is observed that Wall Paint Olympic One has become the top best-selling product in 2023 April which creates an extremely high revenue.

4.1.5 Monthly Shipping Efficiency

```

1 ## Monthly Shipping Efficiency
2 # Convert dates to Date objects
3 shipment$shipment_date <- as.Date(shipment$shipment_date, origin = "1970-01-01")
4
5 shipment_unique <- shipment %>% distinct(shipment_id, .keep_all = TRUE)
6
7 # Merge orders and shipment data on order_id
8 combined_data <- merge(orders, shipment, by = "shipment_id")
9
10 # Calculate shipping duration in days
11 combined_data$shipping_duration <- as.numeric(difftime(combined_data$shipment_date, combined_
12
13 # Calculate monthly statistics
14 monthly_stats <- combined_data %>%
15   mutate(month = floor_date(order_date, "month")) %>%

```

```

16 group_by(month) %>%
17 summarise(Average_Shipping_Duration = round(mean(shipping_duration),2),
18           Min_Shipping_Duration = min(shipping_duration),
19           Max_Shipping_Duration = max(shipping_duration))
20
21 table_to_display <- monthly_stats %>%
22   mutate(Time = format(month, "%Y-%m")) %>%
23   select(Time, Average_Shipping_Duration, Min_Shipping_Duration, Max_Shipping_Duration)
24
25 # Display the table with kable
26 kable(table_to_display, caption = "Monthly Shipping Duration", col.names = c("Time", "Averag

```

Table 2: Monthly Shipping Duration

Time	Average Duration	Min Duration	Max Duration
2023-01	0.56	0	1
2023-02	0.48	0	1
2023-03	0.48	0	1
2023-04	0.77	0	1
2023-05	0.39	0	1
2023-06	0.48	0	1
2023-07	0.30	0	1
2023-08	0.39	0	1
2023-09	0.64	0	1
2023-10	0.44	0	1
2023-11	0.43	0	1
2023-12	0.62	0	1
2024-01	0.59	0	1
2024-02	0.51	0	1

```

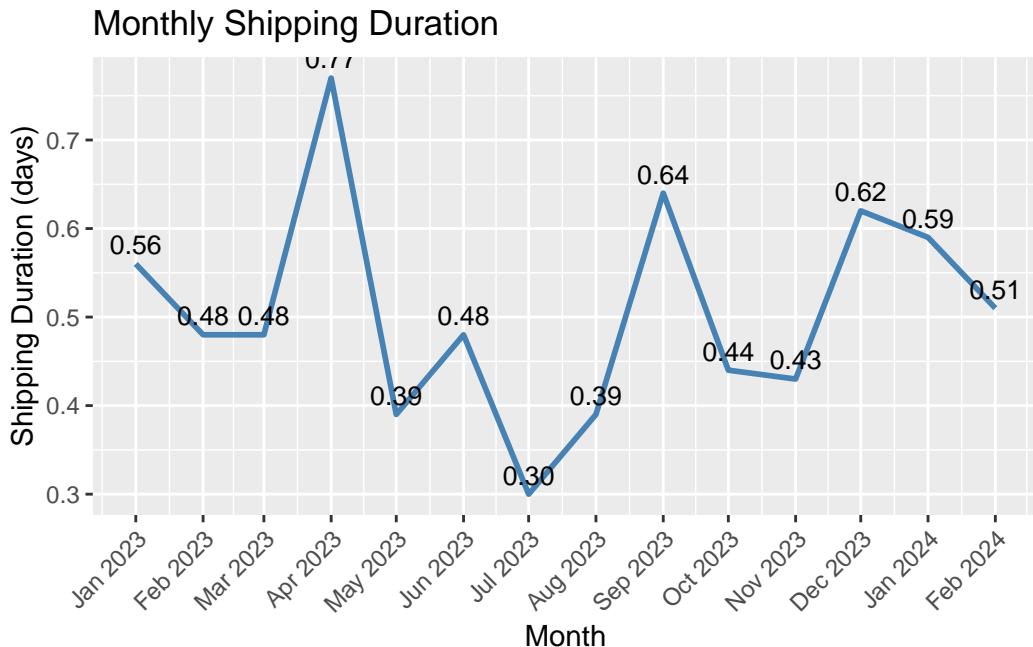
1 # Visualize the statistics
2 g_shipping_efficiency <- ggplot(monthly_stats, aes(x = month)) +
3   geom_line(aes(y = Average_Shipping_Duration), color = "steelblue", size = 1) +
4   geom_text(aes(y = Average_Shipping_Duration,
5                 label = sprintf("%.2f", Average_Shipping_Duration)),
6             position = position_dodge(width = 0.9), vjust = -0.5, size = 3.5) +
7   scale_x_date(date_labels = "%b %Y",
8               date_breaks = "1 month",
9               limits = c(min(monthly_stats$month), max(monthly_stats$month))) +
10  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +

```

```

11   labs(title = "Monthly Shipping Duration",
12       x = "Month", y = "Shipping Duration (days)")
13
14 print(g_shipping_efficiency)

```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("monthly_shipping_efficiency_",
3                     format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
4
5 # Save the plot with the dynamic filename
6 ggsave(filename, plot = g_shipping_efficiency, width = width, height = height)

```

4.1.6 Monthly Delivery Efficiency

```

1 ## Monthly Delivery Efficiency
2 # Convert dates to Date objects
3 shipment$delivery_date <- as.Date(shipment$delivery_date, origin = "1970-01-01")
4
5 shipment_unique <- shipment %>% distinct(shipment_id, .keep_all = TRUE)
6
7 # Merge orders and shipment data on order_id

```

```

8 combined_data <- merge(orders, shipment, by = "shipment_id")
9
10 # Calculate delivery duration in days
11 combined_data$delivery_duration <- as.numeric(difftime(combined_data$delivery_date, combined_
12
13 # Calculate monthly statistics
14 monthly_stats <- combined_data %>%
15   mutate(month = floor_date(order_date, "month")) %>%
16   group_by(month) %>%
17   summarise(Average_Delivery_Duration = round(mean(delivery_duration), 2),
18             Min_Delivery_Duration = min(delivery_duration),
19             Max_Delivery_Duration = max(delivery_duration))
20
21 table_to_display <- monthly_stats %>%
22   mutate(Time = format(month, "%Y-%m")) %>%
23   select(Time, Average_Delivery_Duration, Min_Delivery_Duration, Max_Delivery_Duration)
24
25 # Calculate the overall average delivery duration
26 overall_avg_delivery <- mean(combined_data$delivery_duration, na.rm = TRUE)
27
28 # Display the table with kable
29 kable(table_to_display, caption = "Monthly Delivery Duration", col.names = c("Time", "Average_
```

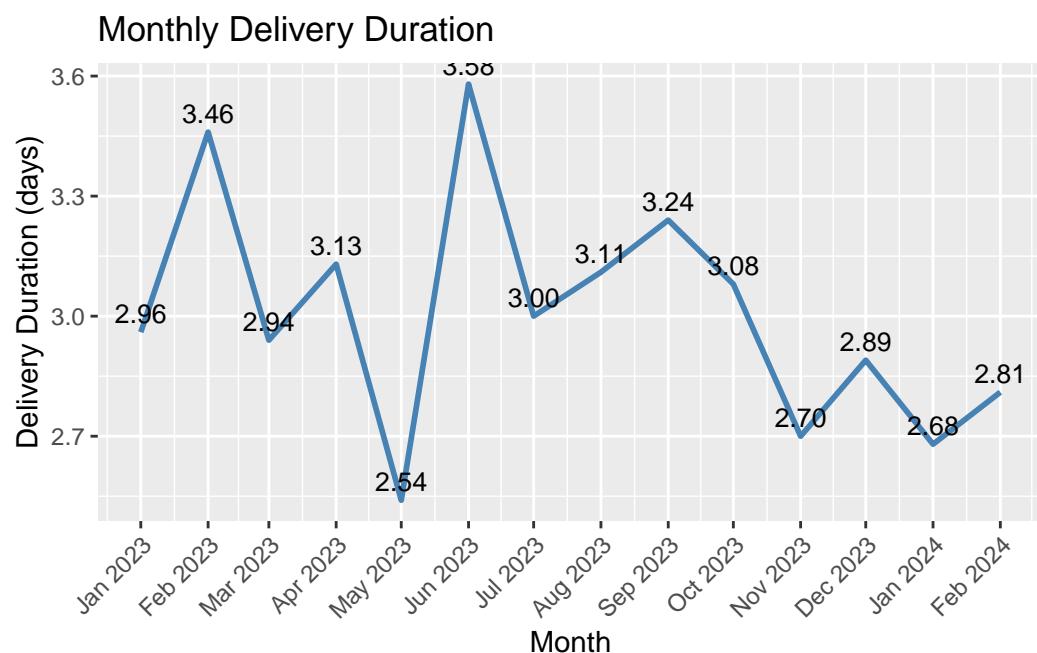
Table 3: Monthly Delivery Duration

Time	Average Duration	Min Duration	Max Duration
2023-01	2.96	1	5
2023-02	3.46	1	5
2023-03	2.94	1	5
2023-04	3.13	1	5
2023-05	2.54	1	5
2023-06	3.58	1	5
2023-07	3.00	1	5
2023-08	3.11	1	5
2023-09	3.24	1	5
2023-10	3.08	1	5
2023-11	2.70	1	5
2023-12	2.89	1	5
2024-01	2.68	1	5
2024-02	2.81	1	5

```

1 # Visualize the statistics
2 g_delivery_efficiency <- ggplot(monthly_stats, aes(x = month)) +
3   geom_line(aes(y = Average_Delivery_Duration), color = "steelblue", size = 1) +
4   geom_text(aes(y = Average_Delivery_Duration,
5                 label = sprintf("%.2f", Average_Delivery_Duration)),
6             position = position_dodge(width = 0.9), vjust = -0.5, size = 3.5) +
7   scale_x_date(date_labels = "%b %Y",
8               date_breaks = "1 month",
9               limits = c(min(monthly_stats$month), max(monthly_stats$month))) +
10  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
11  labs(title = "Monthly Delivery Duration",
12        x = "Month", y = "Delivery Duration (days)")
13 print(g_delivery_efficiency)

```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("monthly_delivery_efficiency_",
3                     format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
4
5 # Save the plot with the dynamic filename
6 ggsave(filename, plot = g_delivery_efficiency, width = width, height = height)

```

Shipping and Delivery Efficiency Analysis

It is observed that the average duration ranges between 0.30 and 0.77 days for shipment efficiency and ranges between 2.54 and 3.58 days for delivery efficiency, indicating the efficient capability of processing parcels for the customers.

4.2 Comprehensive Reporting with Quarto

4.2.1 Demographic Distribution of Customers

A. The Distribution of Gender across Customers

```
1 SELECT
2     gender,
3     COUNT(*) AS GenderCount,
4     CONCAT(ROUND((COUNT(*) * 100.0) / (SELECT COUNT(*) FROM customer), 2), '%') AS Percentage
5 FROM
6     customer
7 GROUP BY
8     gender
9 ORDER BY
10    Percentage DESC;
```

Table 4: 2 records

gender	GenderCount	Percentage
Male	110	55.28%
Female	89	44.72%

B. The Distribution of Age across Customers

```
1 SELECT
2     AgeGroup,
3     COUNT(*) AS Count,
4     CONCAT(ROUND((COUNT(*) * 100.0) / (SELECT COUNT(*) FROM customer), 2), '%') AS Percentage
5 FROM (
6     SELECT
7         customer_id,
8         CASE
9             WHEN age >= 0 AND age < 18 THEN '0-18'
```

```

10      WHEN age >= 18 AND age < 30 THEN '19-30'
11      WHEN age >= 30 AND age < 40 THEN '31-40'
12      WHEN age >= 40 AND age < 50 THEN '41-50'
13      WHEN age >= 50 AND age < 60 THEN '51-60'
14      WHEN age >= 60 AND age < 70 THEN '61-70'
15      WHEN age >= 70 THEN '71+'
16      ELSE 'Unknown'
17  END AS AgeGroup
18  FROM customer
19 ) AS AgeCategories
20 GROUP BY AgeGroup
21 ORDER BY Count DESC;

```

Table 5: 6 records

AgeGroup	Count	Percentage
31-40	68	34.17%
41-50	58	29.15%
19-30	54	27.14%
51-60	16	8.04%
0-18	2	1.01%
61-70	1	0.5%

C. The Distribution of Careers across Customers (Top 10)

```

1 SELECT
2     career,
3     COUNT(*) AS CareerCount,
4     CONCAT(ROUND((COUNT(*) * 100.0) / (SELECT COUNT(*) FROM customer), 2), '%') AS Percentage
5 FROM
6     customer
7 GROUP BY
8     career
9 ORDER BY
10    CareerCount DESC
11 LIMIT 10;

```

Table 6: Displaying records 1 - 10

career	CareerCount	Percentage
Senior Nurse	17	8.54%
Accountant	17	8.54%
Sales Director	16	8.04%
Senior Data Analyst	15	7.54%
Product Manager	13	6.53%
Graphic Designer	12	6.03%
Senior Teacher	11	5.53%
Software Engineer	10	5.03%
HR Manager	10	5.03%
Sales Manager	9	4.52%

D. The Distribution of Geographic Location across Customers (Top 10)

```

1  SELECT
2      address_city,
3      COUNT(*) AS CityCount,
4      CONCAT(ROUND((COUNT(*) * 100.0) / (SELECT COUNT(*) FROM customer), 2), '%') AS Percentage
5  FROM
6      customer
7  GROUP BY
8      address_city
9  ORDER BY
10     CityCount DESC
11    LIMIT 10;

```

Table 7: Displaying records 1 - 10

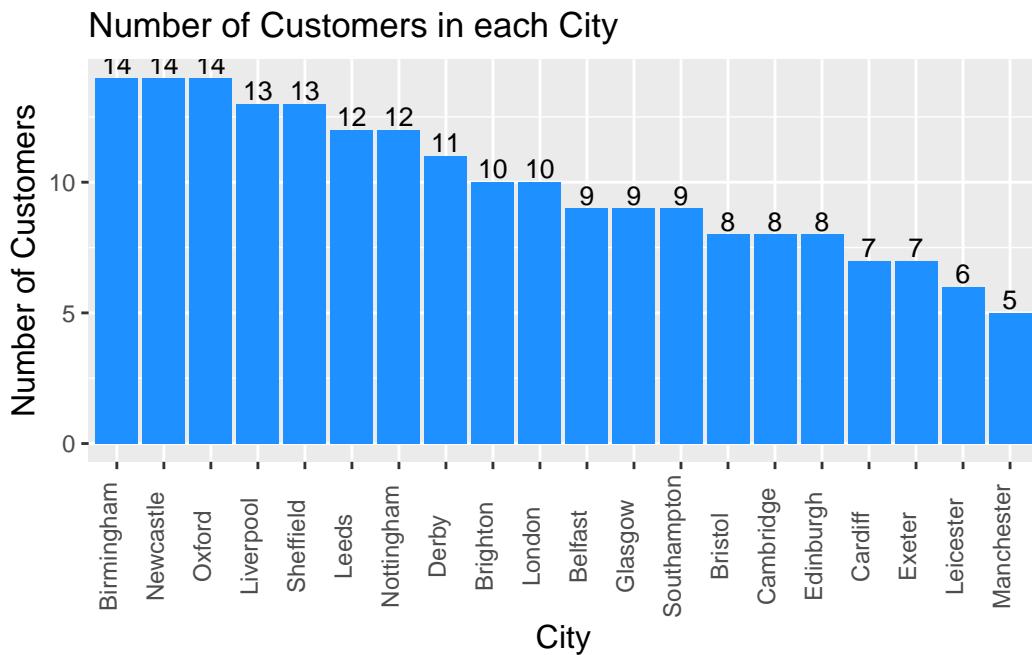
address_city	CityCount	Percentage
Oxford	14	7.04%
Newcastle	14	7.04%
Birmingham	14	7.04%
Sheffield	13	6.53%
Liverpool	13	6.53%
Nottingham	12	6.03%
Leeds	12	6.03%
Derby	11	5.53%
London	10	5.03%
Brighton	10	5.03%

address_city	CityCount	Percentage
--------------	-----------	------------

```

1 # Group by city and count the number of customer in each city
2 customer_city_count <- customer %>%
3   group_by(address_city) %>%
4   summarise(number_of_customers = n()) %>%
5   arrange(desc(number_of_customers))
6
7 # Specify the dimensions of the plot
8 width <- 12
9 height <- 8
10
11 # Use ggplot to create a bar chart showing the number of customers in each city
12 g_customer <- ggplot(customer_city_count,
13   aes(x = reorder(address_city, -number_of_customers),
14     y = number_of_customers)) +
15   geom_col(fill = "dodgerblue") +
16   geom_text(aes(label = number_of_customers),
17     position = position_dodge(width = 0.9), vjust = -0.2,
18     size = 3.5) +
19   theme(axis.text.x = element_text(angle = 90, hjust = 0.5, vjust = 0),
20     axis.title = element_text(size = 12)) +
21   labs(title = "Number of Customers in each City",
22     x = "City",
23     y = "Number of Customers")
24 print(g_customer)

```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("geographical distribution of customers_",
3                     format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
4
5 # Save the plot with the dynamic filename
6 ggsave(filename, plot = g_customer, width = width, height = height)

```

E. The Current Customer Referral Rate

```

1 SELECT
2     COUNT(CASE WHEN referred_by != '' AND referred_by IS NOT NULL THEN 1 END) AS Customer_with_Referral,
3     COUNT(*) AS total_customer,
4     CONCAT(ROUND((COUNT(CASE WHEN referred_by != '' AND referred_by IS NOT NULL THEN 1 END) /
5 FROM
6     customer;

```

Table 8: 1 records

Customer_with_Referral	total_customer	referral_rate
74	199	37.19%

Customer Analysis

The gender of the customers is evenly distributed, with 55% males and 45% females. Their age is mainly located in the senior group, who are in the age group of 31-40 years and 41-50 years, accounting for 34% and 29%, of males and females respectively.

The career distribution states that customers are employed in diverse industries and have different job positions.

The customers currently live in big cities around the United Kingdom, represents people living in big cities who shop online more frequently.

It is observed that 37% of customers are referred by existing customers, showing a moderate satisfaction level from our customers.

4.2.2 Product Portfolio

A. The Distribution of Product Review Scores (Top 10)

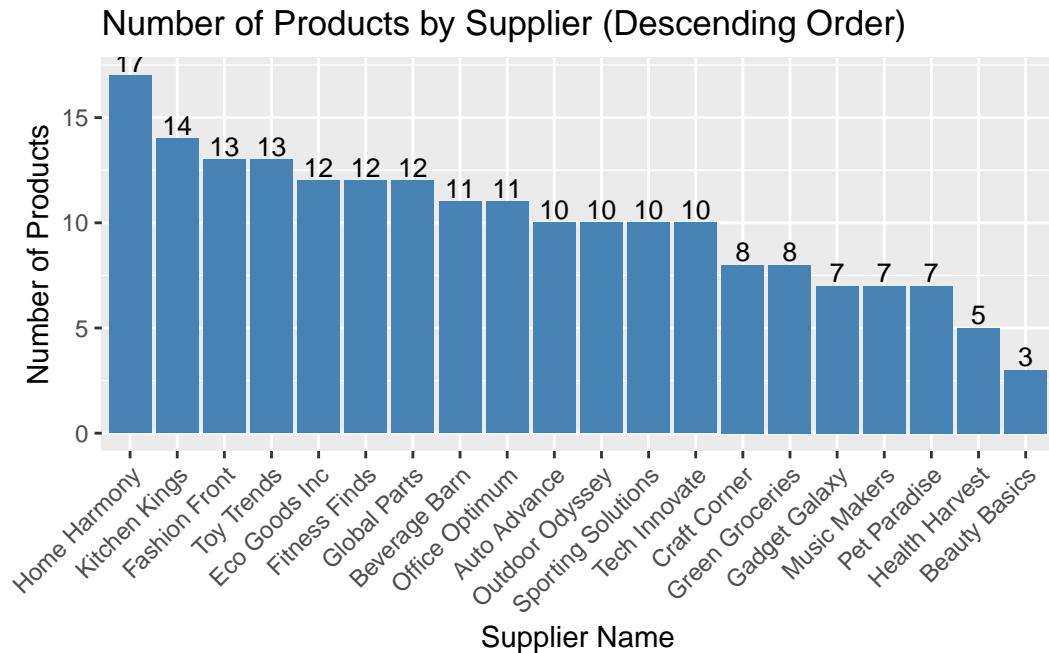
```
1 SELECT
2     product_name, review_score
3 FROM
4     product
5 ORDER BY
6     review_score DESC
7 LIMIT 10;
```

Table 9: Displaying records 1 - 10

product_name	review_score
Office Chair Herman Miller Aeron	4.98
Craft Kit Faber-Castell Young Artist Essentials Set	4.96
Bed Frame Wayfair Zinus Upholstered Platform Bed	4.94
Smartphone Samsung Galaxy S21 Ultra	4.90
Book Penguin Classics Pride and Prejudice	4.89
Laptop HP Spectre x360	4.89
Lipstick NARS Orgasm	4.87
Drill Milwaukee Brushless Cordless Drill	4.84
Lipstick Revlon Super Lustrous Lipstick	4.82
Sofa IKEA Kivik	4.82

B. The Number of Products supplied by Different Suppliers

```
1 # Perform an inner join to combine 'product' with 'supplier' on 'supplier_id'
2 joint_supplier_product <- inner_join(product, supplier, by = "supplier_id")
3
4 # Group by supplier_name and count the number of products for each supplier
5 product_count_by_supplier <- joint_supplier_product %>%
6   group_by(supplier_name) %>%
7   summarise(number_of_products = n())
8
9 # Specify the dimensions of the plot
10 width <- 12
11 height <- 8
12
13 # Use ggplot to create a bar chart showing the number of products for each supplier
14 g_supplier <- ggplot(product_count_by_supplier, aes(x = reorder(supplier_name, -number_of_products),
15   geom_col(fill = "steelblue") +
16   geom_text(aes(label = number_of_products), position = position_dodge(width = 0.9), vjust =
17   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
18   labs(title = "Number of Products by Supplier (Descending Order)",
19       x = "Supplier Name",
20       y = "Number of Products")
21 print(g_supplier)
```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("the product number supplied by supplier_", format(Sys.time(), "%Y%m%d_%H"))
3
4 # Save the plot with the dynamic filename
5 ggsave(filename, plot = g_supplier, width = width, height = height)

```

C. The Product Review Scores of Different Suppliers (Best Top 5)

```

1 SELECT s.supplier_name, ROUND(AVG(p.review_score), 2) AS average_review_score
2 FROM product p
3 JOIN supplier s ON p.supplier_id = s.supplier_id
4 GROUP BY s.supplier_name
5 ORDER BY average_review_score DESC
6 LIMIT 5;

```

Table 10: 5 records

supplier_name	average_review_score
Global Parts	3.90
Pet Paradise	3.81
Office Optimum	3.76
Fashion Front	3.73
Music Makers	3.62

D. The Product Review Scores of Different Suppliers (Worst Top 5)

```

1 SELECT s.supplier_name, ROUND(AVG(p.review_score), 2) AS average_review_score
2 FROM product p
3 JOIN supplier s ON p.supplier_id = s.supplier_id
4 GROUP BY s.supplier_name
5 ORDER BY average_review_score ASC
6 LIMIT 5;

```

Table 11: 5 records

supplier_name	average_review_score
Fitness Finds	2.96
Auto Advance	2.99

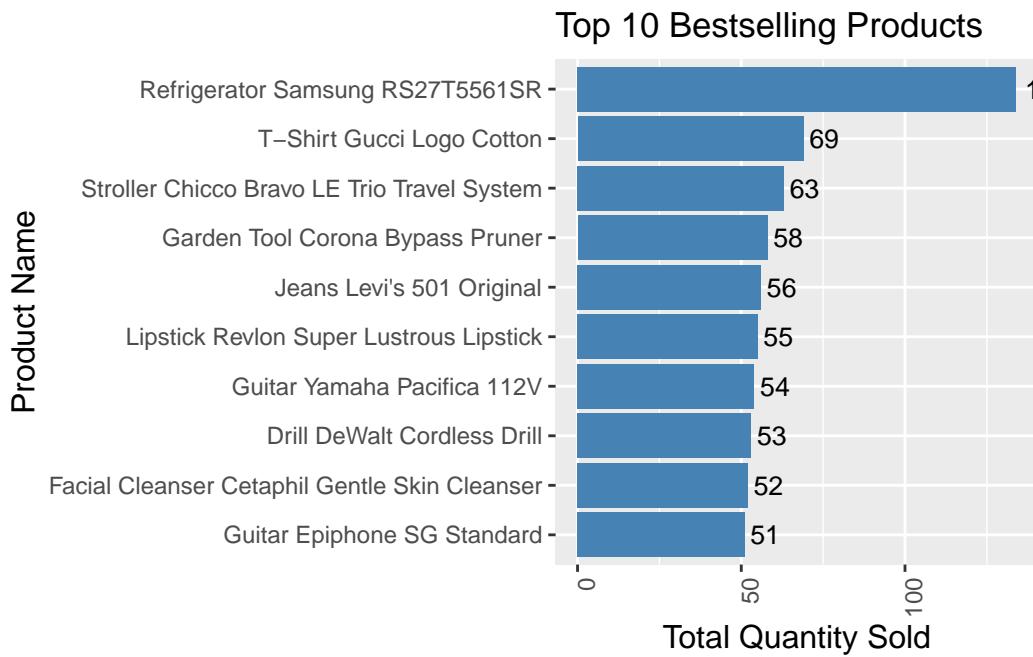
supplier_name	average_review_score
Craft Corner	3.02
Kitchen Kings	3.04
Health Harvest	3.05

E. The Top 10 Best Selling Products

```

1 # Perform an inner join to combine 'orders' with 'product' on 'product_id'
2 joint_order_product <- inner_join(orders, product, by = "product_id")
3
4 # Calculate the total quantity sold for each product
5 product_sales_volume <- joint_order_product %>%
6   group_by(product_name) %>%
7   summarise(total_quantity_sold = sum(quantity)) # Assuming 'quantity' exists in your orders
8 # Processing the text of product name
9 product_sales_volume$product_name <- iconv(product_sales_volume$product_name, "UTF-8", "ASCII")
10
11 # Choose only the top 10 products based on total quantity sold
12 top_product_sales_volume <- product_sales_volume %>%
13   arrange(desc(total_quantity_sold)) %>%
14   slice_head(n = 10)
15
16 # Specify the dimensions of the plot
17 width <- 12
18 height <- 8
19
20 # Use ggplot to create a bar chart showing the total quantity sold for each product
21 g_topproduct <- ggplot(top_product_sales_volume, aes(x = reorder(product_name, total_quantity_sold),
22   geom_col(fill = "steelblue") +
23   geom_text(aes(label = total_quantity_sold), position = position_dodge(width = 0.9), hjust = 0.5) +
24   coord_flip() +
25   theme(axis.text.x = element_text(angle = 90, hjust = 1),
26         axis.title = element_text(size = 12)) +
27   labs(title = "Top 10 Bestselling Products",
28        x = "Product Name",
29        y = "Total Quantity Sold")
30 print(g_topproduct)

```



```

1 # Dynamically generate filename with current date and time
2 filename <- paste0("top10_products_by_quantity_", format(Sys.time(), "%Y%m%d_%H%M%S"), ".png")
3
4 # Save the plot with the dynamic filename
5 ggsave(filename, plot = g_topproduct, width = width, height = height)

```

Supplier Analysis

It shows that the products are supplied by various suppliers, ranging from household items, fashion stores, and toy wholesalers to fitness equipment suppliers. Breaking down into the suppliers' satisfaction, the top five have more than 3.5 out of 5 scores, and the worst five have approximately 3 scores.

Product Analysis

The reviewed scores shows that the top 10 products are mainly from the product category electronic devices and household items such as refrigerators, T-shirt,etc.

4.2.3 Sales Analysis

A. Order Refund Rate

```

1 SELECT
2     COUNT(CASE WHEN refund_status = 'yes' THEN 1 END) AS refund_orders,
3     COUNT(*) AS total_orders,
4     CONCAT(ROUND((COUNT(CASE WHEN refund_status = 'yes' THEN 1 END) * 1.0 / COUNT(*)) * 100,
5 FROM (
6     SELECT DISTINCT order_id, refund_status
7     FROM orders
8 ) AS unique_orders;

```

Table 12: 1 records

refund_orders	total_orders	refund_rate_percentage
151	387	39.02%

Order Refund Analysis

The order refund rate is 39%, higher than the IMRG's 2020 ([benchmarking data](#)); the average e-commerce business saw a 15% return rate. It is a signal for customer reduction since when customers see high refund rates they may be less likely to purchase the products.

Appendices

We asked AI how we can use cat() function in our code. We used cat() in R script to check how the validation and insertion process is going.

 You

do you think that I can implement in my code somewhere the same cat() function to track the progress of validation of new customers data?



 ChatGPT

Absolutely, using `cat()` function to log progress and key information during the validation and data processing steps is an excellent way to track progress and debug your script. It can provide you with immediate feedback in the R console about what your script is doing, which can be particularly helpful for long-running scripts or when processing large datasets.

Here are some places you might consider adding `cat()` statements in your script to log progress and important milestones:

1. Start and End of Processing Each File:

```
r                                     Copy code

for (file_path in customer_file_paths) {
  cat("Starting processing file:", file_path, "\n")
  # Read and process the file...

  cat("Finished processing file:", file_path, "\n")
}
```

2. Before and After Validation:

Inside the loop where you're iterating through each row of the customer data:

We asked AI how to make our validation code for tables more efficient, before it was checking separately which was inefficient. After that prompt, we updated our validation parts and created functions for each of table validations.

 ChatGPT

To make the date validation more efficient, particularly when working with large datasets, you can minimize the number of operations and avoid extracting individual date components unless necessary. One approach is to directly filter out dates that do not meet the criteria using vectorized operations. This reduces the computational overhead by leveraging R's internal optimizations for handling dates and logical vectors.

Here's a revised version of the `validate_and_prepare_orders_data` function with a more streamlined date validation process:

```
r
validate_and_prepare_orders_data <- function(data){
  # Checking format of order id
  order_id_check <- grep("ORDER[0-9]{9}", data$order_id)
  data <- data[order_id_check,]

  customer_id_check <- grep("CUST[0-9]{6}", data$customer_id)
  data <- data[customer_id_check,]

  product_id_check <- grep("[A-Za-z0-9]{10}", data$product_id)
  data <- data[product_id_check,]

  shipment_id_check <- grep("SHIP[0-9]{6}", data$shipment_id)
  data <- data[shipment_id_check,]

  # Convert order_date strings to Date objects
  converted_dates <- as.Date(data$order_date, format = "%d/%m/%Y")

  # Efficient validation for order_date format and range
  valid_dates <- !is.na(converted_dates) & converted_dates >= as.Date("2023-01-01")
}
```

We asked AI to improve my code for orders table validation function, so I can add one more validation check that days are in range of 1-31, months in 1-12 and years starts from 2023 (established of the e-commerce shop). I used the provided code in orders table validation part/

ChatGPT 4 ~

 You

Can you make more code more efficient by adding to the date validation part that days can only

be in 1-31 range, months 1-12 and year is starting from 2023. validate_and_prepare_orders_data <-

function(data){

Checking format of order id

order_id_check <- grep("ORDER[0-9]{9}", data\$order_id)

data <- data[order_id_check,]

customer_id_check <- grep("CUST[0-9]{6}", data\$customer_id)

data <- data[customer_id_check,]

product_id_check <- grep("[A-Za-z0-9]{10}", data\$product_id)

data <- data[product_id_check,]

shipment_id_check <- grep("SHIP[0-9]{6}", data\$shipment_id)

data <- data[shipment_id_check,]

Validation for order_date format

order_date_format_check <- !is.na(as.Date(data\$order_date, format = "%d/%m/%Y"))

data <- data[order_date_format_check]

return(data)

}

This is the reply from AI, we provided him my code for products table from R script and asked him how he can increase efficiency of the code. He suggested to use batch checking method to improve it. Based on that we constructed more efficient code for creating a batch of existing 'product_id' in database and then compare new upcoming primary keys with this vector.

Optimizing Record Existence Checks

Your loop through each product record to check its existence in the database before insertion is effective but could be optimized, especially for large datasets:

1. **Batch Checking:** Instead of individual queries for each `product_id`, consider fetching existing `product_id`'s once into a vector and then checking existence in R. This reduces database queries and can significantly speed up the process.
2. **Temporary Dataframe Uniqueness:** The check for uniqueness in `product_possible_data` is done in an efficient manner. However, for very large dataframes, consider alternatives like `dplyr::distinct()` or maintaining a set of seen keys to avoid the performance cost of subsetting on each iteration.

We provided AI my chunk for products table and asked him to help me integrate referential integrity. We used and adjusted provided code in Product table chunk, the code verifies that in new potential data foreign keys are valid and are presented in associated tables.



ChatGPT

To implement referential integrity checks in your script for the `product` table, you'll need to verify that each `product` record's `category_id`, `supplier_id`, and `promotion_id` exist in their respective reference tables (`category`, `supplier`, `promotion`) before attempting to insert the record into the `product` table. This process ensures that all foreign keys in the `product` table point to existing records, maintaining the integrity of your database relationships.

Here's a way to integrate referential integrity checks into your `validate_and_prepare_product_data` function:

Enhanced Validation Function with Referential Integrity Checks

```
r Copy code

validate_and_prepare_product_data <- function(data, connection) {
  # Validation for product ID
  product_id_check <- grepl("^[A-Za-z0-9]{10}$", data$product_id)
  data <- data[product_id_check, ]

  # Validation for review_score range
  data <- data[data$review_score >= 1 & data$review_score <= 5, ]

  # Fetch existing IDs from reference tables
  valid_category_ids <- dbGetQuery(connection, "SELECT category_id FROM category")$ca
  valid_supplier_ids <- dbGetQuery(connection, "SELECT supplier_id FROM supplier")$su
  valid_promotion_ids <- c(dbGetQuery(connection, "SELECT promotion_id FROM promotion"))

  # Referential integrity checks
  data <- data[data$category_id %in% valid_category_ids, ]
  data <- data[data$supplier_id %in% valid_supplier_ids, ]
  data <- data[is.na(data$promotion_id) | data$promotion_id %in% valid_promotion_ids,
             ]

  return(data)
```

We asked to improve validation function of product table to incorporate more validation constraints on quantity stock, quantity supplied and price. The code was used in validation function of product table.

KE

You

Help me to produce a validation code for 'quantity_stock', "quantity_supplied" that they are non-negative and integer number (fractional numbers are not allowed). also, 'price' should be only positive values. And 'product_name' cannot be empty

< 2 / 2 >



ChatGPT

To extend the `validate_and_prepare_product_data` function with additional validations for `quantity_stock`, `quantity_supplied`, `price`, and ensuring `product_name` is not empty, you can incorporate these checks using logical indexing in R. Here's how you might update your function to include these validations:

Updated Validation Function

```
r Copy code

validate_and_prepare_product_data <- function(data, connection) {
  # Initial validation for product ID and review_score range
  product_id_check <- grepl("^[A-Za-z0-9]{10}$", data$product_id)
  data <- data[product_id_check, ]
  data <- data[data$review_score >= 1 & data$review_score <= 5, ]

  # Fetch existing IDs from reference tables for referential integrity checks
  valid_category_ids <- dbGetQuery(connection, "SELECT category_id FROM category")$category_id
  valid_supplier_ids <- dbGetQuery(connection, "SELECT supplier_id FROM supplier")$supplier_id
  valid_promotion_ids <- c(dbGetQuery(connection, "SELECT promotion_id FROM promotion"))

  # Referential integrity checks
  data <- data[data$category_id %in% valid_category_ids, ]
  data <- data[data$supplier_id %in% valid_supplier_ids, ]
  data <- data[is.na(data$promotion_id) | data$promotion_id %in% valid_promotion_ids,
             ]
  ↓
  # Validation for non-negative and integer quantity_stock and quantity_supplied
```

We asked AI about TryCatch function and how it can be used in our code. We used it R scripts to handle errors and warnings.

KE You

do I need to incorporate tryCatch in my code above?



ChatGPT

Incorporating `tryCatch` into your script can significantly enhance its robustness, especially when dealing with external resources like databases or files where unexpected issues might arise.

`tryCatch` in R is a mechanism for handling errors and warnings gracefully, allowing your script to continue running or to perform specific actions when encountering issues, rather than halting execution unexpectedly.

Here are a few scenarios in your script where `tryCatch` could be particularly beneficial:

1. Database Queries

When fetching existing `product_id`s or performing any database operation, there's always a chance of encountering database connection issues, query syntax errors, or timeouts. Wrapping these operations in `tryCatch` blocks can help handle such errors more gracefully.

```
r Copy code

existing_product_ids <- tryCatch({
  dbGetQuery(connection, "SELECT product_id FROM product")$product_id
}, error = function(e) {
  cat("Error fetching existing product IDs: ", e$message, "\n")
  NULL # Return NULL or an appropriate value to indicate failure
})
```

2. Reading Files

Reading files with `readr::read_csv` could fail due to issues like file access permissions, incorrect file paths, or malformed CSV content. Using `fatch` here allows you to log the error and skip the problematic file without stopping the script.

We asked Chatgpt to use data integrity checks using hashing. Originally we wanted to verify the data integrity before and after loading for the loading data frame. However, due to data type changes and other possible issues they were not validated even if data frames' values were similar.



ChatGPT

To incorporate data integrity checks using hashing before loading and after loading the validated data into the database, you can use the `digest` package in R. Here's how you can modify your script to include these steps:

Step 1: Install and Load the `digest` Package

If you haven't already, install the `digest` package and then load it in your script. This package provides functions for creating hash digests of arbitrary R objects.

```
r
Copy code

if (!requireNamespace("digest", quietly = TRUE)) {
  install.packages("digest")
}
library(digest)
```

Step 2: Compute Hash Before Validation and Loading

After the validation function but before the data loading into the database, compute a hash of `product_possible_data` to capture its state post-validation but pre-loading:

```
r
Copy code

# Ensure product_possible_data is sorted by a consistent key to ensure identical order
product_possible_data <- product_possible_data[order(product_possible_data$product_id)]

# Compute hash of product_possible_data after validation but before loading
pre_load_hash <- digest(product_possible_data, algo = "md5")
cat("Pre-load hash of product_possible_data: ", pre_load_hash, "\n")
```

This prompt helped us to construct the code that is used to hash each row of validated data before loading and after loading. It is done to ensure that loading data and loaded data are the same and was not corrupted during the loading process. The similar codes was used in shipment, orders and product tables data integrity checks.

 You

it is not working anyway. Where else in my code I can use hashing to ensure data integrity?

 ChatGPT

If directly comparing pre-load and post-load data using hashes isn't yielding the expected results due to subtle discrepancies that are hard to control or identify (like minor floating-point differences, data type conversions, or class attribute issues), you might consider alternative approaches to ensure data integrity. Here are some strategies focusing on hashing that might work better in your context:

Hash Individual Records Before Loading and After Retrieval

Instead of hashing the entire dataset at once, consider hashing individual records (rows) before loading them into the database and then hashing the retrieved records again to compare. This approach can help pinpoint exactly where discrepancies occur.

1. **Hash Each Record Before Loading:** Iterate over each row in `product_possible_data`, convert it to a character string or a consistent representation, hash it, and store the hash in a new vector or column.
2. **Load Data into the Database:** Proceed with loading your data into the database as you normally would.
3. **Retrieve and Hash Each Record After Loading:** After loading, retrieve the data (ensuring the same order or using identifiers to match records), hash each retrieved record in the same manner as before, and compare individual hashes.

Implementing Individual Record Hashing

```
r                                     Copy code

# Example for hashing individual records before loading
pre_load_hashes <- sapply(1:nrow(product_possible_data), function(i) {
  record <- as.character(unlist(product_possible_data[i, ]))
  digest(paste(record, collapse = "|"), algo = "md5")
```

We prompted Chatgpt to improve code for orders table using batching processing. We provided our initial code and then based on AI suggestions and iterative prompts process, we improved code to be more efficient and used it in orders table chunks.

 You

it seems that it is sorted by order id but the values with the same order then have different orders in possible and retrieved data.

 ChatGPT

If the sorted `orders_possible_data` and the retrieved data from the database do not maintain the same order for rows with the same `order_id` due to other varying fields, you need a more stable sorting mechanism. One way to address this issue is to sort the data by multiple columns that collectively ensure a unique ordering for every row. This approach would keep the order of rows consistent between the pre-load and post-load datasets.

Given the structure of your `orders` data, it's likely that sorting by `order_id` alone isn't sufficient because each `order_id` could be associated with multiple product, customer, and shipment IDs. Sorting by all these fields should ensure that the order of the data remains consistent between the pre-load phase and the post-load phase.

Here's how you can adjust the sorting mechanism to sort by multiple columns, ensuring consistent order in both `orders_possible_data` and the retrieved data:

```
# Adjust the sorting to include all key columns that might affect the order
orders_possible_data <- orders_possible_data[order(orders_possible_data$order_id, ord
pre_load_hashes <- sapply(1:nrow(orders_possible_data), function(i) {
  record <- as.character(unlist(orders_possible_data[i, ]))
  digest(paste(record, collapse = "|"), algo = "md5")
})

if (nrow(orders_possible_data) > 0) {
  cat("Starting to insert validated data into the database. Number of records: ", nro
```

We prompted AI to help us implement the code for data integrity check for orders table code. After iterative process of prompting and adjusting the code, we implemented data integrity check. So, the code now creates a hash number for validated rows that are ready to be included to data base. After that it retrieves the same loaded data and creates a hash for them. After that it is being compared.



ChatGPT

Your current script processes each order record individually, which is straightforward but might not be the most efficient way, especially if there are many records to process. Batch processing can significantly improve performance by reducing the number of database queries and transactions. Here's how you can modify your code to use batch processing:

Step 1: Collect All IDs First and Perform Bulk Checks

Instead of checking the existence of each ID in the database one by one, collect all unique IDs first and then check their existence in bulk. This can drastically reduce the number of database queries.

```
r Copy code

# Example for batch checking product IDs
unique_product_ids <- unique(orders_data$product_id)
existing_product_ids <- check_id_exists(connection, "product", "product_id", unique_p

# Similarly for customer_ids and shipment_ids
```

Step 2: Validate Referential Integrity in Batches

After getting the existing IDs in bulk, you can check each record against these lists to validate referential integrity without querying the database each time.

```
r Copy code

# Example modification of validate_referential_integrity function to work with pre-f
validate_referential_integrity_batch <- function(new_record, existing_product_ids, ex
product_exists <- new_record$product_id %in% existing_product_ids
customer_exists <- new_record$customer_id %in% existing_customer_ids
shipment_exists <- new_record$shipment_id %in% existing_shipment_ids
```

In this prompt we asked AI to help improve our shipment table code. We wanted to improve checking for existing primary keys in table, so it will query the database once and store it as a vector. It helped to create more efficient code that queries the database once instead of iterative one by one checking.

You

Improve my code for shipment table. So instead of iterative process of querying the existing database to check for uniqueness of primary id. So, fetch existing in current database primary keys to increase performance.

```
validate_and_prepare_shipment_data <- function(data) {  
  # Validation for shipment ID  
  shipment_id_check <- grepl("SHIP[0-9]{6}", data$shipment_id)  
  data <- data[shipment_id_check,]  
  
  # Validation for shipment_date and delivery_date format  
  date_format_check <- !is.na(as.Date(data$shipment_date, format = "%d/%m/%Y")) &  
    !is.na(as.Date(data$delivery_date, format = "%d/%m/%Y"))  
  
  data <- data[date_format_check,]  
  
  #Validation for loogical order of shipment and delivery dates  
  logical_date_order_check <- as.Date(data$shipment_date) < as.Date(data$delivery_date)  
  data <- data[logical_date_order_check,]  
  return(data)  
}  
  
# Fetch existing shipment IDs from the database  
  
shipment_file_paths <- list.files(path = "data_upload", pattern = "shipment.*\\.csv$", full.names =  
  TRUE)  
  
# Define the primary key column for the shipment table  
shipment_primary_key <- "shipment_id"  
  
#Initialising empty data frame  
shipment_possible_data <- data.frame()  
  
# Read each shipment CSV file and chec, ↓ . the existence of the primary key in the database  
before appending
```



ChatGPT

Improving your code for handling the `shipment` table by fetching existing primary keys from the database in one go, rather than checking for uniqueness iteratively for each row, is a great idea for enhancing performance. This approach reduces the number of database queries, thereby minimizing I/O overhead and speeding up the process. Here's how you can modify your code to incorporate this improvement:

1. **Fetch Existing Primary Keys:** Before processing the CSV files, query the database to fetch all existing `shipment_id`'s. Store these IDs in a vector or similar data structure for quick lookup.
2. **Check for Uniqueness Against the Fetched List:** As you process each new record, check if its `shipment_id` is in the list of fetched IDs to determine its uniqueness, instead of querying the database for each record.

Here's the revised code incorporating these improvements:

```
r
Copy code

library(readr)
library(RSQLite)
library(dplyr)

validate_and_prepare_shipment_data <- function(data) {
  # Validation for shipment ID
  shipment_id_check <- grepl("SHIP[0-9]{6}$", data$shipment_id)
  data <- data[shipment_id_check, ]

  # Validation for shipment_date and delivery_date format
  date_format_check <- !is.na(as.Date(data$shipment_date, format = "%d/%m/%Y")) &
    !is.na(as.Date(data$delivery_date, format = "%d/%m/%Y"))
  data <- data[date_format_check, ]
```

```

data <- data[logical_date_order_check, ]

return(data)
}

# Fetch existing shipment IDs from the database
existing_shipment_ids <- dbGetQuery(connection, "SELECT shipment_id FROM shipment")$shipment_id

shipment_file_paths <- list.files(path = "data_upload", pattern = "shipment.*\\.csv$")

# Define the primary key column for the shipment table
shipment_primary_key <- "shipment_id"

# Initialising empty data frame
shipment_possible_data <- data.frame()

for (file_path in shipment_file_paths) {
  cat("Starting processing file:", file_path, "\n")

  # Read the current file
  shipment_data <- readr::read_csv(file_path)

  # Filter out records with existing shipment_id in the database
  unique_shipment_data <- shipment_data[!shipment_data[[shipment_primary_key]] %in% existing_shipment_ids]

  # Combine unique records into the possible data frame
  if(nrow(unique_shipment_data) > 0){
    shipment_possible_data <- rbind(shipment_possible_data, unique_shipment_data)
  }

  cat("Finished processing file:", file_path, "\n")
}

```

We also asked AI to improve the code to ensure data integrity before and after loading the validated data to the database. The AI did not provided the usable code, so it took many iterations of prompting and also applying the existing programming language knowledge.

KE

You

Incorporate in my shipment table code the data integrity check before and after loading the table as I did for product table. The example code is provided: Implementing data integrity check for validated and ready to loading data for each row

```
product_possible_data <- product_possible_data[order(product_possible_data$product_id), ]  
pre_load_hashes <- sapply(1:nrow(product_possible_data), function(i) {  
  record <- as.character(unlist(product_possible_data[i, ]))  
  digest(paste(record, collapse = "|"), algo = "md5")  
})  
  
# Ingesting into database command  
if (nrow(product_possible_data) > 0) {  
  cat("Starting to insert validated data into the database. Number of records: ",  
  nrow(product_possible_data), "\n")  
  
  # Ingesting prepared data to our database  
  tryCatch({  
    dbWriteTable(connection, name = "product", value = product_possible_data, append = TRUE,  
    row.names = FALSE)  
    cat("Data insertion completed successfully.\n")  
  }, error = function(e) {  
    cat("Error inserting data into the database: ", e$message, "\n")  
    # Additional error handling logic here  
  })  
  
  # Fetch the loaded data back for post-load hash comparison  
  loaded_product_ids <- sprintf("%s", product_possible_data$product_id)  
  query <- sprintf("SELECT * FROM product WHERE product_id IN (%s) ORDER BY product_id",  
  paste(loaded_product_ids, collapse = ", "))  
  retrieved_products <- dbGetQuery(connection, query)  
  # Creating hashes for retrieved data from db  
  post_load_hashes <- sapply(1:nrow(retrieved_products), function(i) {  
    record <- as.character(unlist(retrieved_products[i, ]))  
    digest(paste(record, collapse = "|"), algo = "md5")  
  })
```

```

pre_load_hashes <- sapply(1:nrow(shipment_possible_data), function(i) {
  record <- as.character(unlist(shipment_possible_data[i, ]))
  digest(paste(record, collapse = "|"), algo = "md5")
})

if (nrow(shipment_possible_data) > 0) {
  cat("Starting to insert validated data into the database. Number of records: ", nro
}

# Attempt to insert data into the database
tryCatch({
  dbWriteTable(connection, name = "shipment", value = shipment_possible_data, appen
  cat("Data insertion completed successfully.\n")
}, error = function(e) {
  cat("Error inserting data into the database: ", e$message, "\n")
})

# Fetch the loaded data back for post-load hash comparison
loaded_shipment_ids <- sprintf("%s", shipment_possible_data$shipment_id)
query <- sprintf("SELECT * FROM shipment WHERE shipment_id IN (%s) ORDER BY shipmen
retrieved_shipments <- dbGetQuery(connection, query)

# Generate post-load hashes
post_load_hashes <- sapply(1:nrow(retrieved_shipments), function(i) {
  record <- as.character(unlist(retrieved_shipments[i, ]))
  digest(paste(record, collapse = "|"), algo = "md5")
})

# Compare hashes
identical_hashes <- all(pre_load_hashes == post_load_hashes)
if (identical_hashes) {
  cat("Data integrity verified: All record hashes match.\n")
} else {
  cat("Data integrity check failed. Record hashes do not match.\n")
}

```



ChatGPT

To incorporate a data integrity check for the `shipment` table similar to the approach you've taken with the `product` table, we'll follow the same steps: generating pre-load hashes for each row of validated data, inserting the data into the database, retrieving the inserted data, and finally comparing post-load hashes to the pre-load hashes to verify data integrity. Here's how you can adapt your shipment table code to include these steps:

1. **Generate Pre-Load Hashes** for the validated `shipment_possible_data`.
2. **Insert the Data** into the `shipment` table.
3. **Retrieve the Inserted Data** based on the primary keys (`shipment_id`).
4. **Generate Post-Load Hashes** for the retrieved data.
5. **Compare Pre-Load and Post-Load Hashes** to verify data integrity.

Here's the adapted code:

```
r

library(digest)
library(RSQLite)
library(readr)
library(dplyr)

# Assuming connection and other initial setup is done as before

# Validate and prepare data (keep your validation logic here)

# Fetch existing shipment IDs from the database
existing_shipment_ids <- dbGetQuery(connection, "SELECT shipment_id FROM shipment")$s

# List all files, read, filter, and combine unique records into shipment_possible_data
# ... Your existing code for processing CSV files ...

# Ensure shipment_possible_data contains only unique shipment_ids, no repetition
```

We prompted AI to improve the code cause we noticed that our code was not checking for duplicates of rows inside of the new upcoming files. AI guided us to construct a code that also validates for uniqueness inside of new potential files. The similar code was used in all other tables validation part.

KE You

okay, but does the code checks if there are multiple new data but they have the same primary key?

ChatGPT

Your current implementation checks for the uniqueness of `shipment_id` (the primary key) against existing IDs in the database but does not account for the possibility of encountering multiple new records with the same `shipment_id` within the batch of data being processed. This scenario could lead to potential duplicates within `shipment_possible_data` before it's inserted into the database, especially if these duplicates span across different CSV files or are present within a single file.

To ensure each `shipment_id` in `shipment_possible_data` is unique before attempting to insert the data into the database, you should add an additional step to filter out duplicates from `shipment_possible_data` after it's been populated from the CSV files but before the validation step. A straightforward way to do this is to use the `distinct()` function from the `dplyr` package, which allows you to retain only unique rows based on specified columns (in this case, `shipment_id`).

Here's how you can modify your code to include this step:

```
r
Copy code

library(dplyr)

# After combining all unique records into shipment_possible_data
shipment_possible_data <- rbind(shipment_possible_data, unique_shipment_data)

# Ensure shipment_possible_data contains only unique shipment_ids
# before proceeding with validation and insertion.
shipment_possible_data <- shipment_possible_data %>% distinct(shipment_id, .keep_all)

# Now proceed with validation and ↓ther processing
shipment_possible_data <- validate_and_prepare_shipment_data(shipment_possible_data)
```

We also decided to implement a feature that is going to move the source files after the data is processed. We asked AI if it is possible to do and also he provided a code that could be used to implement this adjustment.

KE

You

I have one more question to you. So, right now all new data files (csv) are being uploaded to 'data_upload' folder. And after loading these files to database, they are still saved in this folder. Is it possible to make a code that after ingestion of data, it transfers the used csv file to another folder?



ChatGPT

Yes, it's definitely possible and a good practice to move processed files to another directory after their data has been successfully ingested into the database. This helps in keeping your working directory organized and avoids re-processing the same files inadvertently. In R, you can use the `'file.rename()'` function to move files from one directory to another. Here's how you can integrate this into your existing code:

1. **Define Your Destination Directory:** Before the loop that processes each CSV file, define the path to the directory where you want to move the processed files. If this directory does not exist, you should create it.
2. **Move Processed Files:** After a file has been successfully processed and its data has been ingested into the database, use `'file.rename()'` to move it to the defined directory.

Here's an adapted version of your script with these additions:

```
r

# Assuming existing setup and libraries are already loaded

# Define the source and destination directories
source_dir <- "data_upload"
destination_dir <- "data_processed"

# Create the destination directory if it does not exist
if (!dir.exists(destination_dir)) {  
    ↓
```

Copy code



ChatGPT

If you want every processed file to be uniquely named by appending a timestamp and then moved to the 'data_processed' folder, you can slightly adjust your script to always generate a new filename with a timestamp, regardless of whether the file already exists in the destination directory. This approach ensures that each file is uniquely identified and preserved without overwriting any existing files.

Here's how you can modify your code to achieve this:

```
r Copy code

for (file_path in shipment_file_paths) {
  cat("Starting processing file:", file_path, "\n")

  # Read the current file
  shipment_data <- readr::read_csv(file_path, show_col_types = FALSE)

  # [Your processing logic here]

  # Construct the destination file path with a timestamp to ensure uniqueness
  file_name <- basename(file_path)
  timestamp <- format(Sys.time(), "%Y%m%d%H%M%S")  # Get current timestamp
  file_ext <- tools::file_ext(file_name)  # Extract file extension
  base_name <- sub(paste0("\\.", file_ext, "$"), "", file_name)  # Remove extension
  new_file_name <- paste0(base_name, "_", timestamp, ".", file_ext)  # Create new file name
  dest_file_path <- file.path(destination_dir, new_file_name)  # Construct new destination file path

  # Move the processed file to the destination directory
  if (file.rename(file_path, dest_file_path)) {
    cat("Successfully moved processed file to:", dest_file_path, "\n")
  } else {
    cat("Failed to move file:", file_path, "\n")
  }
  cat("Finished processing file: ", file_path, "\n")
```

KE

You

I want that every file after processing is added with time stamp and moved to 'data processed' folder dest_file_path <- file.path(destination_dir, file_name)

```
# Check if the processed file already exists in the destination directory
if (file.exists(dest_file_path)) {
  # Generate a new file name to avoid overwriting
  # For simplicity, append a timestamp or a unique identifier
  timestamp <- format(Sys.time(), "%Y%m%d%H%M%S")
  file_ext <- tools::file_ext(file_name)
  base_name <- sub(paste0("\\.", file_ext, "$"), "", file_name)
  new_file_name <- paste0(base_name, "_", timestamp, ".", file_ext)
  dest_file_path <- file.path(destination_dir, new_file_name)
}

# Move the processed file to the destination directory
if (file.rename(file_path, dest_file_path)) {
  cat("Successfully moved processed file to:", dest_file_path, "\n")
} else {
  cat("Failed to move file:", file_path, "\n")
}
cat("Finished processing file:", file_path, "\n")
```



ChatGPT

If you want every processed file to be uniquely named by appending a timestamp and then moved to the 'data_processed' folder, you can slightly adjust your script to always generate a new filename with a timestamp, regardless of whether the file already exists in the destination directory. This approach ensures that each file is uniquely identified and preserved without overwriting any existing files.

Here's how you can modify your code to achieve this:





ChatGPT

If you want every processed file to be uniquely named by appending a timestamp and then moved to the 'data_processed' folder, you can slightly adjust your script to always generate a new filename with a timestamp, regardless of whether the file already exists in the destination directory. This approach ensures that each file is uniquely identified and preserved without overwriting any existing files.

Here's how you can modify your code to achieve this:

```
r
Copy code

for (file_path in shipment_file_paths) {
  cat("Starting processing file:", file_path, "\n")

  # Read the current file
  shipment_data <- readr::read_csv(file_path, show_col_types = FALSE)

  # [Your processing logic here]

  # Construct the destination file path with a timestamp to ensure uniqueness
  file_name <- basename(file_path)
  timestamp <- format(Sys.time(), "%Y%m%d%H%M%S")  # Get current timestamp
  file_ext <- tools::file_ext(file_name)  # Extract file extension
  base_name <- sub(paste0("\\.", file_ext, "$"), "", file_name)  # Remove extension
  new_file_name <- paste0(base_name, "_", timestamp, ".", file_ext)  # Create new file name
  dest_file_path <- file.path(destination_dir, new_file_name)  # Construct new destination file path

  # Move the processed file to the destination directory
  if (file.rename(file_path, dest_file_path)) {
    cat("Successfully moved processed file to:", dest_file_path, "\n")
  } else {
    cat("Failed to move file:", file_path, "\n")
  }
  cat("Finished processing file: ", file_path, "\n")
```

AI also helped us to solve the issue which we encountered when our workflow failed to run. The problem was in applying mathematical function on non-numeric values. Also, another problem was in product table that if there was no data existing in temporary data frame, it gave the error when tried to order the frame by primary key. Based on the suggestions provided by AI, changes were implemented to restore the run process.

KE

You

look at this error and wait for my code: I want you after that to solve this issue Record with primary key SOH1230566 already exists in the database.
Finished processing file: data_upload/product.csv
Starting validation for new records.
Error in floor(data\$quantity_stock):
non-numeric argument to mathematical function
Calls: validate_and_prepare_product_data -> [-> [.data.frame
Execution halted
Error: Process completed with exit code 1.



ChatGPT

The error message you've shared indicates two primary issues in your code or process:

1. **Duplicate Primary Key Error:** "Record with primary key SOH1230566 already exists in the database." This suggests that your process attempted to insert a record into the database with a primary key ('SOH1230566') that already exists in the target table. Primary keys must be unique in a database table, so you cannot have two records with the same primary key value.
2. **Non-numeric Argument Error:** "non-numeric argument to mathematical function." This error occurs in the `floor(data\$quantity_stock)` function call, indicating that `data\$quantity_stock` contains non-numeric values, whereas `floor()` requires numeric input to round down to the nearest whole number.

To address these issues, consider the following solutions:

For the Duplicate Primary Key Error:

- **Check for Existing Records:** Before inserting data into the database, check if a record with the same primary key already exists. If it does, decide whether to skip the insertion, update the existing record, or handle it differently based on your requirements.
- **Ensure Unique Primary Keys:** If you're generating primary keys, ensure the generation logic

Revised Numeric Validation Check

Before performing the ``floor()`` operation, convert ``quantity_stock`` and ``quantity_supplied`` to numeric values, ensuring any non-numeric entries are handled gracefully. This can be done using ``as.numeric()``, which will convert non-numeric values to `NA`. After this conversion, you can proceed with your existing checks, but it's also wise to add a step to handle or filter out `NA` values resulting from non-numeric to numeric conversion.

```
# Before the validation checks for quantity_stock and quantity_supplied
# Convert to numeric, handling non-numeric values gracefully
data$quantity_stock <- as.numeric(data$quantity_stock)
data$quantity_supplied <- as.numeric(data$quantity_supplied)

# Any conversion issues will result in NA values. You might want to handle these NAs
# For example, you can decide to filter out these NAs or replace them with a default
# Here we're just ensuring the remaining validation checks only apply to non-NA, numeric values

data <- data[!is.na(data$quantity_stock) & data$quantity_stock >= 0 & (data$quantity_stock <= 100)]
data <- data[!is.na(data$quantity_supplied) & data$quantity_supplied >= 0 & (data$quantity_supplied <= 100)]
```

Handling Conversion Errors

When converting to numeric, ``as.numeric()`` will turn any non-numeric strings to `NA`. This means if ``quantity_stock`` or ``quantity_supplied`` contains such values, they will become `NA` after conversion. You should decide how to handle these cases—whether to filter them out or replace them with a default value (e.g., 0) depends on your specific requirements and data integrity policies.

KE

You

```
> # Implementing data integrity check for validated and ready to loading data for each row  
> product_possible_data <- product_possible_data[order(product_possible_data$product_id), ]  
Error in order(product_possible_data$product_id):  
  argument 1 is not a vector
```



ChatGPT

The error message you're encountering, `argument 1 is not a vector`, suggests that the attempt to order your `product_possible_data` dataframe by `product_id` is failing because `order(product_possible_data\$product_id)` is not receiving a vector as expected. This issue can arise if `product_possible_data` is not properly initialized as a dataframe or if `product_id` does not exist as a column in your dataframe at the point where the ordering is attempted.

There are a couple of ways this issue might be resolved, depending on the state of `product_possible_data` at the time the error occurs:

Ensure `product_possible_data` is Properly Initialized as a Dataframe

Make sure that `product_possible_data` is initialized as a dataframe with the expected columns before you attempt to add rows to it or perform operations like ordering. For example, if `product_possible_data` is expected to have columns matching those in your CSV files, you could initialize it like this:

```
r  
Copy code  
  
# Assuming 'product_data_columns' contains the names of the columns expected in 'prod  
product_possible_data <- data.frame(matrix(ncol = length(product_data_columns), nrow  
colnames(product_possible_data) <- product_data_columns
```

This ensures `product_possible_data` is a frame with the correct structure.

Check for the Existence of `product_id` Before Ordering

Before attempting to order `product_possible_data` by `product_id`, ensure that `product_id` exists as a column. This can be particularly important if `product_possible_data` could be empty or if its structure might vary:

```
r

# Implementing data integrity check for validated and ready to loading data for each
if ("product_id" %in% names(product_possible_data) && nrow(product_possible_data) > 0) {
  product_possible_data <- product_possible_data[order(product_possible_data$product_id)]
} else {
  cat("product_id column not found or product_possible_data is empty.\n")
}
```

This approach checks that `product_id` is a column within `product_possible_data` and that there are rows to order. If these conditions are not met, it prints a message indicating the issue rather than attempting the operation that would lead to an error.

Summary

- **Initialization:** Ensure `product_possible_data` is correctly initialized as a dataframe with the expected structure from the start.
- **Conditional Ordering:** Perform ordering operations conditionally, based on the presence of the column you're ordering by and the existence of rows in your dataframe.

These steps should help you avoid the error and ensure that your data manipulation operations proceed smoothly.