# An Application of the DAMA-DMBOK Data Governance Framework

## Team members:

1. Fatima Aghapourasl
2. Kareem Albeetar
3. Chun Lam Cheung
4. Penelope De Freitas
5. Eltigani Hamadelniel Elfatih Hamadelniel
6. Yi Liu
7. Jeffrey Ching Lok Ng
8. Ivy Tsai

# Table of Contents

# 1.0 Introduction

For this project, we were tasked with applying the DAMA-DMBOK Data Governance framework to the ecommerce industry (Future Shoppers Superstore). In particular, we chose three (3) core Data Governance (DG) sub-domains: Data Quality, Data Modeling and Design, and Data Integration and Interoperability. In addition, we created sample policies and procedures for each of the aforementioned sub-domains.

Another aspect of this project involved conducting analyses and producing plans, and models for managing the superstore database. Furthermore, two (2) reports (Operational and Executive) were produced to assist Regional Managers and the COO of the Future Shoppers Superstore in decision-making.

The following sections are a compilation of the various pieces that were produced by our team in an effort to apply Data Governance using the Future Shoppers Superstore scenario.

## 2.0 The Application of Data Governance (DG)

In this section, we identify the three (3) core Data Governance (DG) sub-domains that our group chose to work on for the project. We also provide two (2) examples of principles, policies, and procedures for each of the identified sub-domains.

## 2.1 Core DG sub-domains

The following core DG subdomains were identified by our team:

    2.1.1 **Data Quality**

        The data quality sub-domain is relevant to this project since it affects the accuracy of business analytics. Reliable data is crucial for understanding trends in sales, regional, and product performance, which inform business decisions.

    2.1.2 **Data Modeling and Design -**

        This subdomain is relevant to this project since it is key to defining a common vocabulary around data serving as the common source of truth for communication within the business.

    2.1.3 **Data Integration and Interoperability**

        This subdomain was chosen because in modern information systems and data-driven organizations, it is critical to ensure that the data can be combined, shared within various sources, and used across different systems.

## 2.2 Core DG sub-domains' Principles

In this section we identify two sample principles for each of the following sub-domains: Data Quality, Data Modeling and Design, and Data Integration and Interoperability.

    2.2.1 Data Quality
- a) **Objective Measurement and Transparency:** share relevant data with stakeholders about data quality to promote transparency (e.g., units, calculations).
- b) **Error Prevention:** strive to prevent data errors and conditions that diminish data usability

    2.2.2 Data Modeling and Design
- a) **Formalization:** strive for the creation of data entities that provide a complete scope and definition of integral processes within the business.
- b) **Data Reusability And Flexibility:** strive for the creation of an adaptable data model that can accommodate changes and updates over time.

    2.2.3 Data Integration and Interoperability
- a) **Balance Local Data Needs with Enterprise Data Needs**: Data should possess the capability to be shared and manipulated to generate diverse reports that meet the future decision-making requirements of the enterprise.
- b) **Ensure Business Accountability:** Ensuring business units take ownership of data quality, compliance, and governance as per our operational and executive reports.

## 2.3 Core DG sub-domains' Policies

Now we will expand the aforementioned principles by presenting two associated sample policies:

    2.3.1 Data Quality
- a) **Objective Measurement and Transparency**
  1. Document and preserve the history of data use and manipulation during data quality management.
  2. Being transparent about the standards, requirements, and specifications for data quality management.

**b) Error Prevention**
1. Implement guidelines for data entry to prevent errors from occurring.
2. Define procedures for data validation and validation checks to identify and address issues before they affect data quality.

2.3.2 Data Modeling and Design
**a) Formalization**
1. All entities and attributes created must document explicit knowledge about business systems and processes.
2. All data types used for specific entities must be standardized across data models in different business units.

**b) Data Reusability and Flexibility**
1. All entities in the data model must be flexible to accommodate circumstantial business processes.
2. All data model components must be designed to be reusable for different business units and processes.

2.3.3 Data Integration and Interoperability
**a) Balance Local Data Needs with Enterprise Data Needs**
1. Data must be shared between departments to ensure it can be integrated and used effectively across the enterprise.
2. Extract data and generate reports for individuals overseeing enterprise revenues, making future enterprise-related decisions, or as requested.

**b) Ensure Business Accountability:**
1. Uphold accuracy and integrity in data collection, management, and reporting across all departments. Correct discrepancies or errors promptly.
2. Regular monitoring of profit margins to identify areas for improvement and ensure financial accountability is crucial. This includes a review of profit and loss figures as presented in the dataset, aiding in a clearer understanding of the financial health and operational efficiency across different regions and segments.

# 2.4 Core DG sub-domains' Procedures

In this section, we further extend the aforementioned principles and policies, by providing sample procedures:

2.4.1 Data Quality
**a) Objective Measurement and Transparency**
1. Document and preserve the history of data use and manipulation during data quality management.
   a. All decisions that were made pertaining to the manipulation of data (e.g., cleaning and creation) should be reported and provided with an explanation
   b. Data that has been cleaned should be recorded as metadata in order to keep a good data lineage. We may have been able to improve this by keeping a metadata table to record who, when, and what was changed.

2. Being transparent about the standards, requirements, and specifications for data quality management.
   a. All calculated fields should be noted in all reports as determined by the data governance committee (e.g., monthly and quarterly).
   b. Verify that calculation methods and units align with industry standards and guidelines to prevent any confusion.

**b) Error Prevention**

1. Guidelines for data entry must be implemented to prevent errors from occurring.
   a. Clear and detailed data entry standards that cover data formats, naming conventions, units of measurement, and any specific guidelines relevant to the subject are developed and acknowledged.
   b. A review process is performed in which a second person in the team checks the entered data for accuracy and adherence to data entry standards before it's finalized.

2. Procedures for data validation and validation checks must be defined to identify and address issues before they affect data quality.
   a. We developed validation guidelines before processing, such as format validation, range checks, and data type verification to guarantee the accuracy and consistency of the data.
   b. We used consistent formulas for all the calculations, then combined that with manual checks to fix any errors, while maintaining logs to keep a complete data lineage.

2.4.2 Data Modeling and Design
   **a) Formalization**
   1. All entities and attributes created must document explicit knowledge about business systems and processes.
      a. All data entities must have clear attributes that clearly communicate their corresponding scope and definition. We created new attributes within the order table to indicate whether or not an order has been returned.
      b. All relationships must clearly communicate all constraints between physical entities. We created a weak entity "Product_Order" to clearly illustrate the relationship between the entities "Product" and "Order" which it defines the products purchased within each order and their corresponding characteristics.

   2. All data types used for specific entities must be standardized across data models in different business units.
      a. All attributes used across different entities must comply with their corresponding standardized data type. We set the data type of the "Order ID" attribute as CHAR(50) across all entities where "Order ID" is a foreign key.
      b. All primary keys must use the data type that offers the best performance. We set the data type of integer for all primary keys for all entities except for those primary keys that show specific patterns that relate to other attributes.

   **b) Data Reusability and Flexibility**
   1. All entities in the data model must be flexible to accommodate circumstantial business processes.
      c. We created a "Return" attribute in the "Order" entity to accommodate any orders that are returned by a customer.
      d. We created a metadata table to provide more context about all orders that have been indicated as a return order in the "Order" entity. Thus providing more flexibility to our data model without hurting its referential integrity.

   2. All data model components must be designed to be reusable for different business units and processes.
      a. We implement the use of a reference product table that can reused for multiple orders. The reference product table features information about the characteristics of a product which will be reused each time a customer orders that product.

        b. We apply data normalization to our data to eliminate redundancy and maintain referential integrity. We eliminated the presence of many-to-many relationships thus removing any partial redundancy.

2.4.3 Data Integration and Interoperability
  a) **Balance Local Data Needs with Enterprise Data Needs**
    1. Data must be shared between departments to ensure it can be integrated and used effectively across the enterprise.
      a. Create data backups before sharing to account for the possibility of unauthorized modifications to the data.
      b. Save the requested data as an Excel file and share it with other departments while ensuring that unauthorized individuals cannot gain access to the database system.

    2. Extract data and generate reports for individuals overseeing enterprise revenues, making future enterprise-related decisions or as requested.
      a. Generate monthly operational reports in the form of Excel files for the regional sales managers to overview their monthly sales performances and make adjustments accordingly.
      b. Generate quarterly executive reports in the form of Excel files for the COO to observe and keep track of overall performances between various quarters and targets in order to make decisions and set goals for the enterprise.

  b) **Ensure Business Accountability:**
    1. Uphold accuracy and integrity in data collection, management, and reporting across all departments. Correct discrepancies or errors promptly.
      a. Implement data validation and verification during data integration processes, such as ensuring the 'Quantity' column of the original dataset cannot be less than one.
      b. Record all changes to the data schemas, integration workflows, etc., and inform all relevant stakeholders for version controls.

    2. Provide training and guidelines.
      a. Create guidelines outlining how to interpret and present data from our dataset accurately, including calculating metrics like Order Return Rate, Average Transaction Value, and Gross Profits. Providing clear guidelines and training will help minimize errors and ensure that data from different regions and segments can be easily consolidated and compared.
      b. Create a centralized digital repository for easy access to guidelines and other reference materials.
        Establish a review process for draft reports ensuring adherence to guidelines and accuracy in data representation.
        adhering to assumptions, aiming to minimize errors, and ensuring consistency in data handling and reporting across different regions and segments.

# 3.0 Final Reports

## 3.1 Operational Report

**Monthly Regional Sales Team Operations Report 1.0**
*04-01-2019 to 04-30-2019*
*Page 1 of 3*

| Region | Segment | Category | (#) Orders | Order Return Rate (%) | (#) Discounted Products | (#) Quantity | Sales (USD $) | Profit (USD $) | Average Transaction Value (USD $) |
|---|---|---|---|---|---|---|---|---|---|
| Central | Consumer | Furniture | 6 | 0.00% | 11 | 14 | $1,453.98 | -$168.14 | $242.33 |
| | | Office Supplies | 22 | 0.00% | 53 | 84 | $2,579.72 | -$1,544.26 | $117.26 |
| | | Technology | 9 | 0.00% | 28 | 31 | $2,841.97 | $763.14 | $315.77 |
| | Consumer Total | | 37 | 0.00% | 92 | 129 | $6,875.67 | -$949.25 | $185.83 |
| | Corporate | Furniture | 5 | 20.00% | 1 | 18 | $717.25 | $304.87 | $143.45 |
| | | Office Supplies | 8 | 12.50% | 15 | 18 | $132.16 | -$0.76 | $16.52 |
| | | Technology | 5 | 20.00% | 8 | 21 | $881.63 | $175.30 | $176.33 |
| | Corporate Total | | 18 | 16.67% | 24 | 57 | $1,731.04 | $479.42 | $96.17 |
| | Home Office | Furniture | 1 | 0.00% | 5 | 5 | $213.12 | -$15.22 | $213.12 |
| | | Office Supplies | 6 | 33.33% | 8 | 18 | $4,205.81 | $1,996.34 | $700.97 |
| | Home Office Total | | 7 | 28.57% | 13 | 23 | $4,418.93 | $1,981.12 | $631.28 |
| **Central Total** | | | **62** | **8.06%** | **129** | **209** | **$13,025.64** | **$1,511.29** | **$210.09** |

**Monthly Regional Sales Team Operations Report 1.0**
*04-01-2019 to 04-30-2019*
*Page 2 of 3*

| Region | Segment | Category | (#) Orders | (#) Discounted Products | (#) Quantity | Order Return Rate (%) | Sales (USD $) | Profit (USD $) | Average Transaction Value (USD $) |
|---|---|---|---|---|---|---|---|---|---|
| East | Consumer | Furniture | 1 | 0 | 3 | 0.00% | $1,267.53 | $316.88 | $1,267.53 |
| | | Office Supplies | 11 | 19 | 48 | 0.00% | $1,191.19 | $421.15 | $108.29 |
| | | Technology | 5 | 5 | 18 | 0.00% | $2,609.62 | $475.13 | $521.92 |
| | Consumer Total | | 17 | 24 | 69 | 0.00% | $5,068.34 | $1,213.16 | $298.14 |
| | Corporate | Furniture | 1 | 2 | 2 | 0.00% | $127.76 | $2.84 | $127.76 |
| | | Office Supplies | 6 | 23 | 23 | 0.00% | $133.00 | -$38.32 | $22.17 |
| | | Technology | 1 | 3 | 3 | 0.00% | $118.78 | -$27.72 | $118.78 |
| | Corporate Total | | 8 | 28 | 28 | 0.00% | $379.54 | -$63.19 | $47.44 |
| | Home Office | Furniture | 4 | 10 | 12 | 0.00% | $1,015.57 | -$254.73 | $253.89 |
| | | Office Supplies | 9 | 27 | 30 | 0.00% | $820.63 | $79.48 | $91.18 |
| | | Technology | 2 | 7 | 7 | 0.00% | $56.84 | -$13.01 | $28.42 |
| | Home Office Total | | 15 | 44 | 49 | 0.00% | $1,893.04 | -$188.27 | $126.20 |
| **East Total** | | | **40** | **96** | **146** | **0.00%** | **$7,340.93** | **$961.71** | **$183.52** |
| South | Consumer | Furniture | 9 | 16 | 37 | 22.22% | $3,364.34 | $167.41 | $373.82 |
| | | Office Supplies | 6 | 6 | 10 | 16.67% | $189.01 | $15.42 | $31.50 |
| | | Technology | 3 | 0 | 13 | 33.33% | $1,146.09 | $321.23 | $382.03 |
| | Consumer Total | | 18 | 22 | 60 | 22.22% | $4,699.44 | $504.06 | $261.08 |
| | Corporate | Office Supplies | 1 | 1 | 1 | 0.00% | $157.79 | -$115.72 | $157.79 |

**Link to the spreadsheet with the [formatted Operational Report](#).**

## Assumptions:

- Monetary values are represented in US dollars.
- The customer's 'Region' was associated with the 'Region' of the product
- All rows with duplicate order IDs were dropped from the database.
- Negative profits were retained in the database and were interpreted as a loss.
- Negative profits were due to discounts and products being sold at a price lower than their original cost or big discounts due to loyalty programs.
- Products that have a negative profit on discount were either purchased through store loyalty points or during a clearance sale to clear out inventory.
- Discounts were expressed as the percentage of the sales that were reduced (expressed as a ratio)
- The 'Order ID' matching between 'Orders' and 'Returns' sheets determines which specific products were returned. Consequently, when a product is returned, we can assume that all the products from that order were returned and are available for resale.
- Each region aims to hit its business goal of a 2% increase in gross profit per year.
- Numbers formatted in the colour "Red" correspond to poor performance (e.g. negative profits or region did not meet the target KPI)

# 3.1.1 Insights Garnered

- All regions earned a profit for April 2019.
- Out of all of the regions, the "Central" one gained the most profit for April 2019, whereas the "South" region gained the lowest profit.
- For the "Central" region, the 'Home Office' segment gained the highest profit; for the "East" region, the Consumer segment gained the highest profit; for the "West" region, the Home Office segment gained the highest profit; for the "South" region, the Consumer segment gained the highest profit.
- The worst performing segments for the "Central", "East", "West" and "South" regions were Consumer, Home Office, Corporate, and Home Office respectively.
- The "Central" region had the most effective promotions with the highest number of discounted products and the highest revenue among all regions
- The "South" region had the highest order return rate, this high return rate could be an indication of why the region produced the least profits.
- Customers shopping in the "South" region have the most spending as indicated by the region's significant average transaction value.

# 3.2 Executive Report

**Quarterly Executive Report 1.0**

*Gross Profits for Q1 and Q2 (2019) vs Q2 of 2018 and 2019*

*Page 1 of 1*

| Region | 2019 Gross Profits USD ($) | | | KPI: Q2 (%) To Target | Q2 Gross Profits USD ($) | | KPI: 2019 (%) To Target |
|---|---|---|---|---|---|---|---|
| | Q1 | Q2 | Q2 VS Q1 | | 2018 | 2018 vs 2019 (%) | |
| West | $3,666.56 | $4,001.37 | 109.13% | 82.55% | $4,752.04 | 84.20% | 34.90% |
| Central | -$84.17 | $3,270.44 | 3885.60% | 151.07% | $2,122.45 | 154.09% | 163.94% |
| East | -$156.05 | $3,588.37 | 2299.46% | 650.86% | $540.51 | 663.88% | 21.76% |
| South | $5,782.80 | $1,330.74 | 23.01% | 34.32% | $3,801.12 | 35.01% | 63.59% |
| *Grand Total* | *$9,209.14* | *$12,190.92* | *132.38%* | *106.56%* | *$11,216.13* | *109.00%* | *71.05%* |

**Link to the spreadsheet with the [formatted Executive Report.](#)**

## Assumptions

- Monetary values are represented in US dollars.
- The customer's 'Region' was associated with the 'Region' of the product
- All rows with duplicate order IDs were dropped from the database.
- Negative profits were retained in the database and were interpreted as a loss.
- Negative profits were due to discounts and products being sold at a price lower than their original cost or big discounts due to loyalty programs.
- Products that have a negative profit on discount were either purchased through store loyalty points or during a clearance sale to clear out inventory.
- Discounts were expressed as the percentage of the sales that were reduced (expressed as a ratio)
- The 'Order ID' matching between 'Orders' and 'Returns' sheets determines which specific products were returned. Consequently, when a product is returned, we can assume that all the products from that order were returned and are available for resale.
- Each region aims to hit its business goal of a 2% increase in gross profit per year.
- Numbers formatted in the colour "Red" correspond to poor performance (e.g. negative profits or region did not meet the target KPI)

## 3.2.1 Insights Garnered

- Moving from Q1 to Q2 of 2019, it was observed that all regions except the "South" experienced an increase in profits
- The "Central" and "East" regions surpassed the 2% target (KPI) for Q2 of 2019, while the "West" and "South" regions failed to reach the 2% target (KPI) for Q2 2019
- Comparing Q2 gross profits for 2018 and 2019, the "Central" and 'South" regions surpassed the 2% target (KPI)

# 4.0 Annex

## 4.1 Lab Exercise 1 - submission A

### 4.1.1 Introduction

**Future Shoppers Superstore** is an e-commerce brand that offers a range of furniture, office supplies, and technology. They market products to three customer segments: Consumers, Home Offices, and Corporate companies.

In this exercise, the Data Analytics department (Group 6) was tasked with examining and analyzing orders, returns, and people data for the Future Shoppers Superstore. In addition, we were expected to create report templates for communicating useful insights regarding regional sales-related information to the Regional Sales Managers and Chief Operating Officer.

The following report highlights our findings from analyzing the superstore's data, the context and assumptions of our reporting, and the recommended template structures.

### 4.1.2 An Analysis of the Sample Superstore Dataset

In this section, we report on the observations that were made from the exploration and analysis of the Superstore dataset, based on completeness, inconsistencies, redundancies, and duplicates.

- **Completeness**:
    1. 9994 records exist in the 'Orders' datastore
    2. Out of the 9994 'Orders' data, 5009 maps to unique Order IDs
    3. All 5009 orders were placed between 2018-01-03 and 2021-12-30.
    4. A total of 11 missing values were found and they all pertained to the 'Postal Code' column. The missing postal code rows all had Burlington for 'City' and Vermont for 'State'. However, the same city and state combination may have more than one postal code.
    5. A 'Customer ID' maps to one 'Customer Name' and vice versa. (one-to-one relationship).
    6. Each sub-category belongs to only one category (many-to-one relationship)
    7. Each state belongs to only one region (many-to-one relationship).
    8. Some cities were found under different states, and a state has multiple cities (many-to-many relationship).
    9. There are 4 ship modes: 'Same Day', 'First Class', 'Second Class', and 'Standard Class'.
    10. There are 3 customer segments: 'Consumer', 'Corporate' and 'Home Office'.
    11. There are 531 cities and 49 states in total.
    12. There are 4 regions where each has a regional manager: 'West', 'East', 'Central', and 'South'.
    13. There are 3 product categories: 'Office Supplies', 'Furniture' and 'Technology', and a total of 17 product sub-categories.
    14. In the 'Returns' sheet, there are only two columns, 'Returned' and 'Order ID.' However, it's important to note that the 'Order ID' in both 'Orders' and 'Returns' sheets had the same occurrences. This indicates that entire packages with the same 'Order ID' were indeed returned. Nevertheless, it is recommended to consider adding more features to the 'Returns' sheet for comprehensive record-keeping purposes.

- **Inconsistencies**:
    1. There were multiple product names sharing the same product ID, however, each name should only be associated with one ID. (one-to-one relationship)

2. The date differences between 'Ship Date' and 'Order Date' vary upon 'Ship Mode' where some may have occurred due to early and late delivery. For instance, the date differences for ship mode 'First Class' ranges from 1-4 days, but only one record had a 4-day difference. Although this also occurred in 'Second Class' and 'Standard Class' only 3 other records had this issue.
3. Negative values were found in column 'Profit'. With further investigation, we noticed that 1022 orders had negative profits, but only 54 of those orders got returned. Furthermore, the total number of returned orders was 296. Our team assumed that some possible explanations for having more negative profit orders than returned orders may be loyalty points being used, extra costs for missing packages or damaged products during delivery, and last but not least clearance sales for excessive inventories.

- **Redundancies**:
    1. Since column 'Country/Region' has the same value as 'United States' for all rows, we can consider dropping this column during data cleaning to reduce the size of the dataset.
    2. In the 'Returns' sheet column 'Returned' was all 'Yes' which was meaningless and redundant. This column should be removed.
    3. The 'Orders' sheet data should be broken down into separate tables. Such as product, customer, category, order, etc. Each table will only store features that are related to that particular table. For example, the 'Products' table will contain columns 'Product ID' and 'Product 'Name'. By doing so we would reduce redundancies and follow data normalization rules for better data management.

- **Duplicates**:
    1. When setting 'Row ID' as the index in the 'Orders' sheet, values for row IDs 3406 and 3407 were exactly the same. It was most likely that the same Order ID containing the same product was generated twice accidentally; therefore we should consider dropping the row with ID 3407 for data cleaning.
    2. We identified multiple duplicate rows in the 'Returns' sheet. Upon cross-referencing the 'Order ID' column values with the 'Orders' sheet, we found that the number of duplicate rows in the 'Returns' sheet matches the number in the 'Orders' sheet where customers returned the entire package. This suggests that customers return the entire package in all instances of a return.

## 4.1.3 Operational and Executive Report Planning

## <u>Target Audience</u>

The target audience and intended use of the Operational and Executive reports are as follows:

**Operational Report**
- Target Audience: Regional Sales Managers of West, East, South and Central regions
- Intended use:
    - To assess the regional branches' monthly performance on Office Supplies, Furniture, and Technology product categories.
    - To develop customer segmentation strategies.
    - To perform monthly analyses of store performance over all of the regions
    - To see how discounts offered impact product profitability.
    - To make a decision on what products should go on sale in the upcoming month.
    - To examine the impact of return rate on profitability.

**Executive Report**
- Target Audience: COO (Chief Operating Officer) of the Future Shoppers Superstore
- Intended use:
    - To observe Quarter-on-Quarter performance across multiple years
    - To observe if the business is on track to achieve its quarterly goal.
    - To make decisions regarding the improvement of company profits

## Context and Additional Assumptions

We created Operational and Executive reports for the Regional Managers and COO of **Future Shoppers Superstore** who are preparing for the new quarter's sales (Q2) of 2020.

To support the process, we examined the profits and the product trends for April 2019 as well as comparisons for Q1 and Q2 of 2018 and 2019. These reports will aid the administrative officers in making decisions regarding products with the most and least profits for the regional branches.

Going forward, the Operational report will be utilized by the regional operational teams for future monitoring and reference. Similarly, the Executive report will assist the executive team in current and future decision-making.

## Assumptions
- Monetary values are represented in US dollars.
- The customer's 'Region' was associated with the 'Region' of the product
- All rows with duplicate order IDs were dropped from the database.
- Negative profits were retained in the database and were interpreted as a loss.
- Negative profits were due to discounts and products being sold at a price lower than their original cost or big discounts due to loyalty programs.
- Products that have a negative profit on discount were either purchased through store loyalty points or during a clearance sale to clear out inventory.
- Discounts were expressed as the percentage of the sales that were reduced (expressed as a ratio)
- The 'Order ID' matching between 'Orders' and 'Returns' sheets determines which specific products were returned. Consequently, when a product is returned, we can assume that all the products from that order were returned and are available for resale.
- Each region aims to hit its business goal of a 2% increase in gross profit per year.
- Numbers formatted in the colour "Red" correspond to poor performance (e.g. negative profits or region did not meet the target KPI)

## Proposed Information to Report

### Operational Report

The Operational report will include the following fields:

- **Region** - The list of regions (West, East, Central, South) where products are sold
- **Segment** - The various customer segments (Consumer, Corporate, Home Office)
- **Category** - The categories of products (Office Supplies, Furniture, Technology)

   *Calculated Fields*

- **(#) Orders** - The number of unique order IDs
- **Quantity (#)** - The sum of the total number of products sold per product category
- **Sales (USD $)** - The revenue generated from products across various categories
- **Profits (USD $)** - The sum of the profits generated from products across various categories
- **Order Return Rate (%)**
  - The ratio of products returned to the quantity sold per product category
  - **Formula**: Number of orders returned / (#) Orders.
- **No. of products with discounts** - The total number of products that were offered discounts per product category.
- **Average Transaction Value (USD $)**
  - the average amount a customer spends on a single order.
  - **Formula**: Sales (USD $) / (#) Orders.

**Executive Report**

The Executive report will include the following fields:

- **Region** - The list of regions (West, East, Central, South) where products are sold

  *Calculated Fields*

- **2019 Gross Profits (USD $)**
  - **Quarter 1 (Q1)** - The total profits per region for Q1 of 2019
  - **Quarter 2 (Q2)** - The total profits per region for Q2 of 2019
  - **Q2 vs Q1 (%)** - The ratio of Q1 to Q2 gross profits for 2019

- **KPI: Q2 (%) To Target (%)**
  - The proportion of the quarter's target gross profit accumulated.
  - **Formula**: Q2 Gross Profit / Q2 Target Gross Profit
- **Q2 Gross Profits**
  - **2018 (USD $)** - The total profits per region for Q2 of 2018
  - **2019 vs 2018 (%) -** The ratio of 2018 to 2019 gross profits for Q2

- **KPI: 2019 (%) To Target**
  - The proportion of the year's target gross profit accumulated
  - **Formula:** 2019 Gross Profit / 2019 Target Gross Profit

## 4.1.4 Operational Report Template

### Monthly Regional Sales Team Operations Report

*04-01-2019 - 04-30-2019*

| Region | Segment | Category | (#) Orders | Order Return Rate (%) | (#) Products with discount | Quantity | Sales (USD $) | Profits (USD $) | Average Transaction Value (USD $) |
|---|---|---|---|---|---|---|---|---|---|
| West | Consumer | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Corporate | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Home Office | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Total | | | | | | | | |
| East | Consumer | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Corporate | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |

| Region | Segment | Category | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Home Office | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | **Total** | | | | | | | | |
| **Central** | Consumer | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Corporate | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Home Office | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | **Total** | | | | | | | | |
| **South** | Consumer | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |
| | Subtotal | | | | | | | | |
| | Corporate | Office Supplies | | | | | | | |
| | | Furniture | | | | | | | |
| | | Technology | | | | | | | |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Subtotal | | | | | | | | |
| Home Office | Office Supplies | | | | | | | | |
| | Furniture | | | | | | | | |
| | Technology | | | | | | | | |
| | Subtotal | | | | | | | | |
| **Total** | | | | | | | | | |
| | *Grand Total* | | | | | | | | |

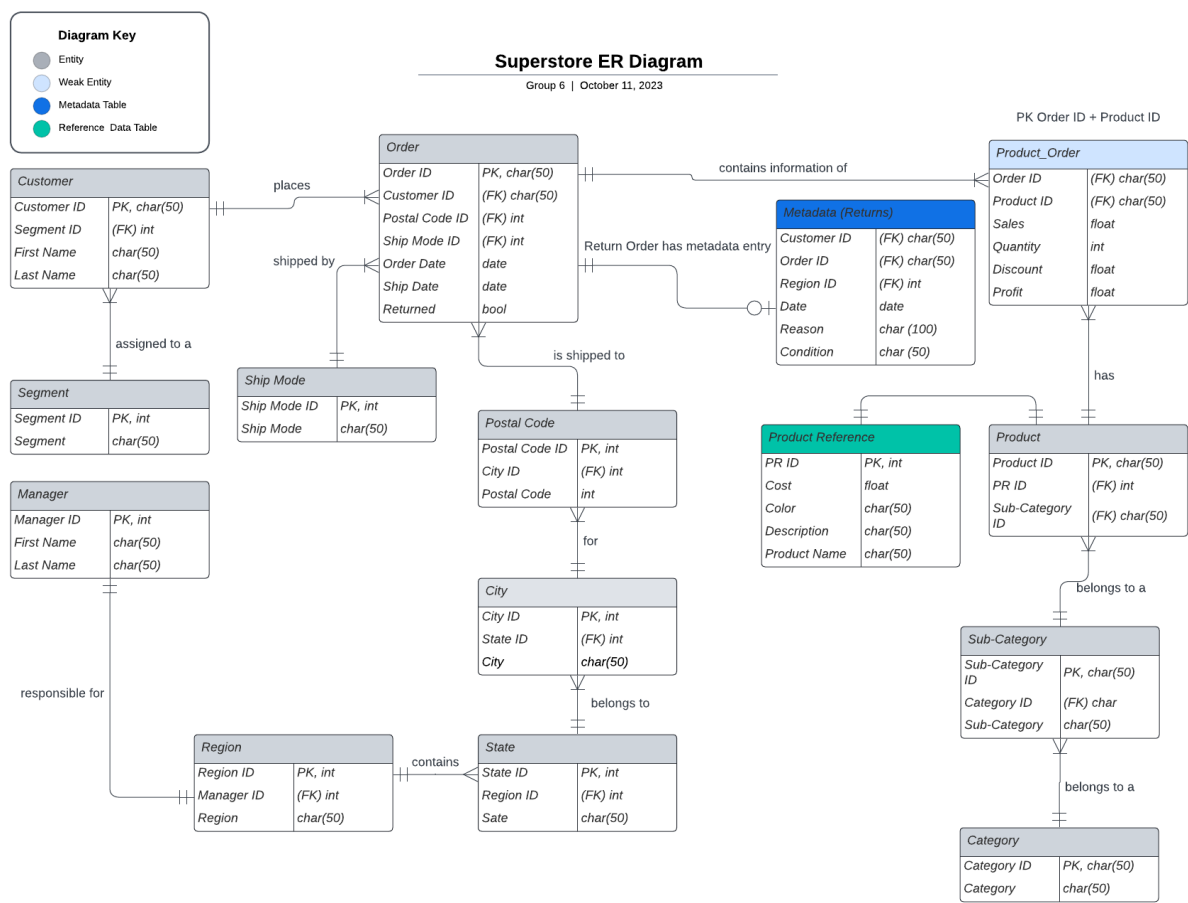## 4.1.5 Executive Report Template

### Quarterly Executive Report

*Gross Profits for Q1 and Q2 (2019) vs Q2 of 2018 and 2019*

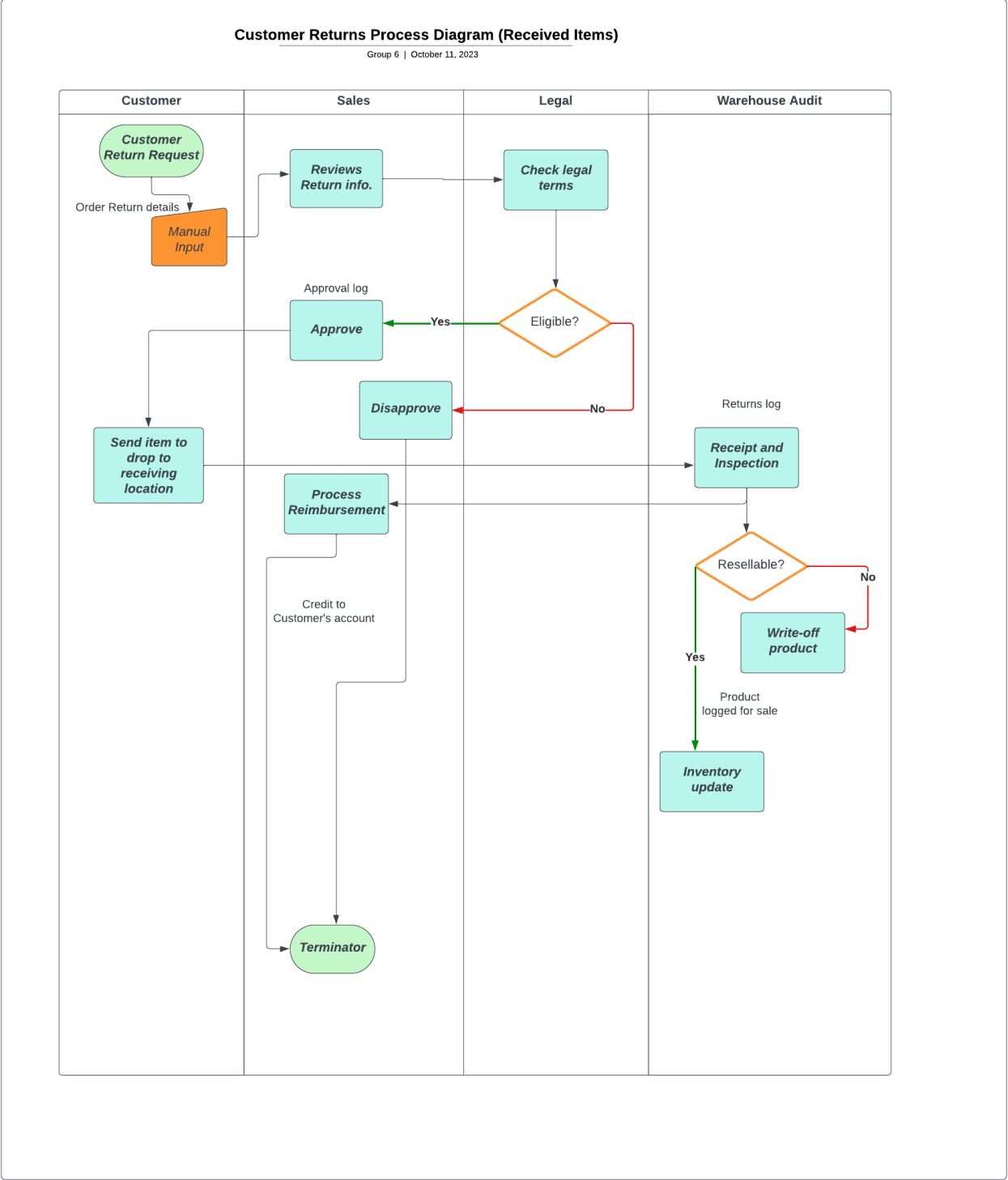| Region | 2019 Gross Profits USD ($) | | | KPI: Q2 (%) To Target | Q2 Gross Profits USD ($) | | KPI: 2019 (%) To Target |
|---|---|---|---|---|---|---|---|
| | **Q1** | **Q2** | **Q2 vs Q1 (%)** | | **2018** | **2019 vs 2018 (%)** | |
| West | | | | | | | |
| East | | | | | | | |
| Central | | | | | | | |
| South | | | | | | | |
| *Grand Total* | | | | | | | |

Link to the spreadsheet with operational report template and executive report template: ⊞ Report Templates

# 4.2 Lab Exercise 2 - submission B
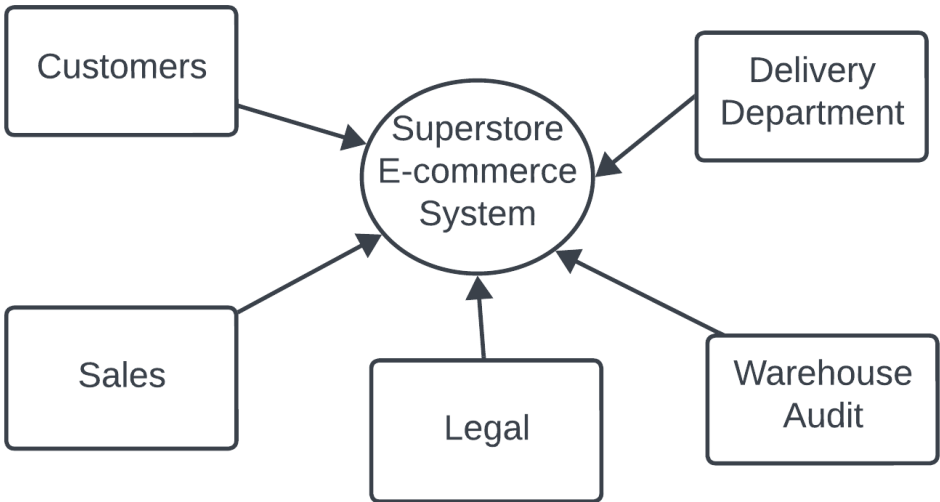
## 4.2.1 Logical ERD Diagram

## 4.2.2 Dataflow



**Customer Returns Process Diagram (Received Items)**

Group 6 | October 11, 2023

| Customer | Sales | Legal | Warehouse Audit |
|---|---|---|---|

Customer Return Request

Order Return details

Manual Input

Reviews Return info.

Check legal terms

Approval log

Approve

Eligible?

Yes

No

Disapprove

Send item to drop to receiving location

Returns log

Receipt and Inspection

Process Reimbursement

Resellable?

No

Write-off product

Credit to Customer's account

Yes

Product logged for sale

Inventory update

Terminator

20

# Customer Ordering Process Diagram

Group 6 | October 11, 2023

**Level 0 - Superstore DF Diagram**

Group 6  |  October 11, 2023

Customers

Delivery Department

Superstore E-commerce System

Sales

Legal

Warehouse Audit

LEVEL 0



**Level 1 - Customer Ordering DF Diagram**

Group 6  |  October 11, 2023

E-commerce DB

Product Information

Query Product Information

Find Products

Product information

Product

Order information

Order modification information

Change Order

Product Preferences e.g. quantity, style

Cart Information

Ordering Process

Shopping Cart

Shipping mode, Address, etc.

Interarctions

Payment method and validations

Shipping Details

Cart Information

Customer

Shipping Information

Confirm

Payment Details

Payment information

Checkout Item

Make Payments

LEVEL 1

22

## 4.2.3 Database Schema



## 4.2.4 Data Cleaning

The following are all the changes that we have implemented to obtain a final clean version of the superstore dataset:

1. Dropped all duplicate rows present within the dataset.

2. We switched the date format to YYYY-MM-DD to match the default format of MySQL workbench.

3. Added a separate column for the category code that was taken from the first part of the product ID

4. Added a separate column for the subcategory code that was taken from the second part of the product ID

5. We made another column called costs to make a single product to have the cost of the product by using the formula, $Cost = (Sales - Profit)/Quantity$

6. Since we have some products that have the same name but different Product IDs, and we want to make Product ID unique, we simply appended an A or B for each product that has duplicates:

7. Switched all Product IDs with Product Name: "Sauder Forest Hills Library, Woodland Oak Finish" from ID "FUR-BO-10002213" to "FUR-BO-10002213A".

8. Switched all Product IDs with Product Name: "DMI Eclipse Executive Suite Bookcases" from ID "FUR-BO-10002213" to "FUR-BO-10002213B".

9.  Switched all Product IDs with Product Name: "Global Value Mid-Back Manager's Chair, Gray" from ID "FUR-CH-10001146" to "FUR-CH-10001146A".
10. Switched all Product IDs with Product Name: "Global Task Chair, Black" from ID "FUR-CH-10001146" to "FUR-CH-10001146B".

11. Switched all Product IDs with Product Name: "DAX Wood Document Frame" from ID "FUR-FU-10001473" to "FUR-FU-10001473A".
12. Switched all Product IDs with the Product Name: "Eldon Executive Woodline II Desk Accessories, Mahogany" from ID "FUR-FU-10001473" to "FUR-FU-10001473B".

13. Switched all Product IDs with Product Name: "Executive Impressions 13" Chairman Wall Clock" from ID "FUR-FU-10004017" to "FUR-FU-10004017A".
14. Switched all Product IDs with Product Name: "Tenex Contemporary Contur Chairmats for Low and Medium Pile Carpet, Computer, 39" x 49"" from ID "FUR-FU-10004017" to "FUR-FU-10004017B".

15. Switched all Product IDs with Product Name: "Eldon 200 Class Desk Accessories, Black" from ID "FUR-FU-10004091" to "FUR-FU-10004091A".
16. Switched all Product IDs with "Howard Miller 13" Diameter Goldtone Round Wall Clock" from ID "FUR-FU-10004091" to "FUR-FU-10004091B".

17. Switched all Product IDs with Product Name: "Eldon Image Series Desk Accessories, Burgundy" from ID "FUR-FU-10004270" to "FUR-FU-10004270A".
18. Switched all Product IDs with Product Name: "Executive Impressions 13" Clairmont Wall Clock" from ID "FUR-FU-10004270" to "FUR-FU-10004270B".

19. Switched all Product IDs with Product Name: "DAX Solid Wood Frames" from ID "FUR-FU-10004848" to "FUR-FU-10004848A".
20. Switched all Product IDs with Product Name: "Howard Miller 13-3/4" Diameter Brushed Chrome Round Wall Clock" from ID "FUR-FU-10004848" to "FUR-FU-10004848B".

21. Switched all Product IDs with Product Name: "Eldon 500 Class Desk Accessories" from ID "FUR-FU-10004864" to "FUR-FU-10004864A".
22. Switched all Product IDs with Product Name: "Howard Miller 14-1/2" Diameter Chrome Round Wall Clock" from ID "FUR-FU-10004864" to "FUR-FU-10004864B".

23. Switched all Product IDs with Product Name: "Belkin 7 Outlet SurgeMaster II" from ID "OFF-AP-10000576" to "OFF-AP-10000576A".
24. Switched all Product IDs with Product Name: "Belkin 325VA UPS Surge Protector, 6" from ID "OFF-AP-10000576" to "OFF-AP-10000576B".

25. Switched all Product IDs with Product Name: "Avery Hi-Liter Comfort Grip Fluorescent Highlighter, Yellow Ink" from ID "OFF-AR-10001149" to "OFF-AR-10001149A".
26. Switched all Product IDs with Product Name: "Sanford Colorific Colored Pencils, 12/Box" from ID "OFF-AR-10001149" to "OFF-AR-10001149B".

27. Switched all Product IDs with Product Name: "Ibico Recycled Linen-Style Covers" from ID "OFF-BI-10002026" to "OFF-BI-10002026A".
28. Switched all Product IDs with Product Name: "Avery Arch Ring Binders" from ID "OFF-BI-10002026" to "OFF-BI-10002026B".

29. Switched all Product IDs with Product Name: "GBC Binding covers" from ID "OFF-BI-10004632" to "OFF-BI-10004632A".
30. Switched all Product IDs with Product Name: "Ibico Hi-Tech Manual Binding System" from ID "OFF-BI-10004632" to "OFF-BI-10004632B".

31. Switched all Product IDs with Product Name: "Avery Binding System Hidden Tab Executive Style Index Sets" from ID "OFF-BI-10004654" to "OFF-BI-10004654A".
32. Switched all Product IDs with Product Name: "VariCap6 Expandable Binder" from ID "OFF-BI-10004654" to "OFF-BI-10004654B".

33. Switched all Product IDs with Product Name: "White Dual Perf Computer Printout Paper, 2700 Sheets, 1 Part, Heavyweight, 20 lbs., 14 7/8 x 11" from ID "OFF-PA-10000357" to "OFF-PA-10000357A".
34. Switched all Product IDs with the Product Name: "Xerox 1888" from ID "OFF-PA-10000357" to "OFF-PA-10000357B".

35. Switched all Product IDs with the Product Name: "Xerox 1952" from ID "OFF-PA-10000477" to "OFF-PA-10000477A".
36. Switched all Product IDs with the Product Name: "Xerox 22" from ID "OFF-PA-10000477" to "OFF-PA-10000477B".

37. Switched all Product IDs with Product Name: 'Adams Phone Message Book, Professional, 400 Message Capacity, 5 3/6" x 11' from ID "OFF-PA-10000659" to "OFF-PA-10000659A".
38. Switched all Product IDs with Product Name: "TOPS Carbonless Receipt Book, Four 2-3/4 x 7-1/4 Money Receipts per Page" from ID "OFF-PA-10000659" to "OFF-PA-10000659B".

39. Switched all Product IDs with the Product Name: "Xerox 2" from ID "OFF-PA-10001166" to "OFF-PA-10001166A".
40. Switched all Product IDs with the Product Name: "Xerox 1932" from ID "OFF-PA-10001166" to "OFF-PA-10001166B".

41. Switched all Product IDs with the Product Name: "Xerox 1881" from ID "OFF-PA-10001970" to "OFF-PA-10001970A".
42. Switched all Product IDs with the Product Name: "Xerox 1908" from ID "OFF-PA-10001970" to "OFF-PA-10001970B".

43. Switched all Product IDs with the Product Name: "RSVP Cards & Envelopes, Blank White, 8-1/2" X 11", 24 Cards/25 Envelopes/Set" from ID "OFF-PA-10001970" to "OFF-PA-10001970A".
44. Switched all Product IDs with the Product Name: "Xerox 1966" from ID "OFF-PA-10001970" to "OFF-PA-10001970B".

45. Switched all Product IDs with the Product Name: "Adams Telephone Message Book W/Dividers/Space For Phone Numbers, 5 1/4"X8 1/2", 200/Messages" from ID "OFF-PA-10002377" to "OFF-PA-10002377A".
46. Switched all Product IDs with the Product Name: "Xerox 1916", 200/Messages" from ID "OFF-PA-10002377" to "OFF-PA-10002377B".

47. Switched all Product IDs with the Product Name: "Standard Line "While You Were Out" Hardbound Telephone Message Book" from ID "OFF-PA-10003022" to "OFF-PA-10003022A".
48. Switched all Product IDs with the Product Name: "Xerox 1992" from ID "OFF-PA-10003022" to "OFF-PA-10003022B".

49. Switched all Product IDs with Product Name: "Fellowes Personal Hanging Folder Files, Navy" from ID "OFF-ST-10001228" to "OFF-ST-10001228A".
50. Switched all Product IDs with the Product Name: "Personal File Boxes with Fold-Down Carry Handle" from ID "OFF-ST-10001228" to "OFF-ST-10001228B".

51. Switched all Product IDs with Product Name: "Acco Perma 3000 Stacking Storage Drawers" from ID "OFF-ST-10004950" to "OFF-ST-10004950A".
52. Switched all Product IDs with the Product Name: "Tenex Personal Filing Tote With Secure Closure Lid, Black/Frost" from ID "OFF-ST-10004950" to "OFF-ST-10004950B".

53. Switched all Product IDs with the Product Name: "Logitech G19 Programmable Gaming Keyboard" from ID "TEC-AC-10002049" to "TEC-AC-10002049A".
54. Switched all Product IDs with the Product Name: "Plantronics Savi W720 Multi-Device Wireless Headset System" from ID "TEC-AC-10002049" to "TEC-AC-10002049B".

55. Switched all Product IDs with the Product Name: "Maxell 4.7GB DVD-RW 3/Pack" from ID "TEC-AC-10002550" to "TEC-AC-10002550A".
56. Switched all Product IDs with Product Name: "Memorex 25GB 6X Branded Blu-Ray Recordable Disc, 30/Pack" from ID "TEC-AC-10002550" to "TEC-AC-10002550B".

57. Switched all Product IDs with Product Name: "Imation 16GB Mini TravelDrive USB 2.0 Flash Drive" from ID "TEC-AC-10003832" to "TEC-AC-10003832A".
58. Switched all Product IDs with the Product Name: "Logitech P710e Mobile Speakerphone" from ID "TEC-AC-10003832" to "TEC-AC-10003832B".

59. Switched all Product IDs with the Product Name: "Okidata MB491 Multifunction Printer" from ID "TEC-MA-10001148" to "TEC-MA-10001148A".
60. Switched all Product IDs with Product Name: "Swingline SM12-08 MicroCut Jam Free Shredder" from ID "TEC-MA-10001148" to "TEC-MA-10001148B".

61. Switched all Product IDs with the Product Name: "Cisco Unified IP Phone 7945G VoIP phone" from ID "TEC-PH-10001530" to "TEC-PH-10001530A".
62. Switched all Product IDs with the Product Name: "Plantronics Voyager Pro Legend" from ID "TEC-PH-10001530" to "TEC-PH-10001530B".

63. Switched all Product IDs with Product Name: "ClearOne CHATAttach 160 - speaker phone" from ID "TEC-PH-10001795" to "TEC-PH-10001795A".
64. Switched all Product IDs with Product Name: "RCA H5401RE1 DECT 6.0 4-Line Cordless Handset With Caller ID/Call Waiting" from ID "TEC-PH-10001795" to "TEC-PH-10001795B".

65. Switched all Product IDs with Product Name: "Aastra 6757i CT Wireless VoIP phone" from ID "TEC-PH-10002200" to "TEC-PH-10002200A".
66. Switched all Product IDs with the Product Name: "Samsung Galaxy Note 2" from ID "TEC-PH-10002200" to "TEC-PH-10002200B".

67. Switched all Product IDs with the Product Name: "Panasonic KX T7731-B Digital phone" from ID "TEC-PH-10002310" to "TEC-PH-10002310A".
68. Switched all Product IDs with Product Name: "Plantronics Calisto P620-M USB Wireless Speakerphone System" from ID "TEC-PH-10002310" to "TEC-PH-10002310B".

69. Switched all Product IDs with Product Name: "AT&T CL2909" from ID "TEC-PH-10004531" to "TEC-PH-10004531A".
70. Switched all Product IDs with the Product Name: "OtterBox Commuter Series Case - iPhone 5 & 5s" from ID "TEC-PH-10004531" to "TEC-PH-10004531B".

## 4.2.5 Extract, Transform and Load (ETL)

*Assumption*: All primary keys are created by MySQL Workbench

*Entity*: Segment
*Extract*: Unique values from the "Segment" Column from superstore dataset
*Transform*: Export data as a CSV File
*Load*: Load the csv file into MySQL workbench

*Entity*: Customer
*Extract*: Unique values from "Customer ID", "Customer Name", "Segment" Columns from superstore dataset

***Transform***:
1. Split the "Customer Name" into two columns: First Name and Last
2. Drop "Customer Name" Column
3. Export data as a CSV File

***Load***: Load the csv file into MySQL workbench


***Entity***: Ship Mode
***Extract***: Unique values from "Ship" Column from superstore dataset
***Transform***: Export data as a CSV File
***Load***: Load the csv file into MySQL workbench


***Entity***: Manager
***Extract***: Unique values from "Regional Manager" Column from superstore dataset
***Transform***:
1. Split the "Regional Manager" into two columns: First Name and Last
2. Drop "Regional Manager" Column
3. Export Data as a CSV file

***Load***: Load the csv file into MySQL workbench


***Entity***: Region
***Extract***: Unique values from "Region" and "Manger Name" Column from superstore dataset
***Transform***:
1. Replace "Manager Name" with its corresponding Manager ID
2. Export Data as a CSV file

***Load***: Load the csv file into MySQL workbench


***Entity***: State
***Extract***: Unique values from "Region" and "State" Column from superstore dataset
***Transform***:
3. Replace Region with its corresponding Region ID
4. Export Data as a CSV file

***Load***: Load the csv file into MySQL workbench


***Entity***: City
***Extract***: Find all unique pairs for "State" and "City" Columns from superstore dataset
***Transform***:
1. Replace City with its corresponding City ID
2. Export Data as a CSV file

***Load***: Load the csv file into MySQL workbench


***Entity***: Category
***Extract***: Find all unique values for "Category" Columns from superstore dataset
***Transform***: Export Data as a CSV file
***Load***: Load the csv file into MySQL workbench


***Entity***: Sub Category
***Extract***: Find all unique values for "Sub-Category" and "Category" Columns from superstore dataset
***Transform***:
1. Replace "Category" with its corresponding Category ID
2. Export Data as a CSV file

***Load***: Load the csv file into MySQL workbench

***Entity***: Postal Code
***Extract***: Find all unique values for "Postal Code" and "City" Columns from superstore dataset
***Transform***:
1.  Replace "City" with its corresponding City ID
2.  Export Data as a CSV file

***Load***:  Load the csv file into MySQL workbench


***Entity***: Product Reference
***Extract***: Find all unique values for the "Product Name" and "Cost" Columns from superstore dataset
***Transform***:
1.  Created a list of colours for the products
2.  Created an empty "Colour" Column
3.  Randomly select a colour from the colour list for each product
4.  Created an empty Description Column and populate it with sample descriptions

***Load***:  Load the csv file into MySQL workbench


***Entity***: Product
***Extract***: Find all unique values for "Product ID",  "Product Name", "Subcategory" and  "Cost" Columns from superstore dataset
***Transform***:
1.  We transformed Product ID to be in the following format: Category ID - SubCategory ID - Product ID
2.  Replaced "Subcategory" with its corresponding Subcategory ID
3.  Created a new column "PR ID" and populated it with the corresponding product reference ID values from the Product Reference Table

***Load***:  Load the csv file into MySQL workbench


***Entity***: Order
***Extract***: Find all unique values for "Order ID", "Customer ID", "Postal Code", "Ship Mode" "OrderDate" and "ShipDate" Columns from superstore dataset
***Transform***:
1.  Replaced "Postal Code" with its corresponding Postal Code ID
2.  Replaced "Ship Mode" with its corresponding Ship Mode ID
3.  Created a new boolean column "Returned"
4.  Populated the return table with boolean values "True" if the Order ID for that record is present in our list of returned Orders and "False" if it is not present

***Load***:  Load the csv file into MySQL workbench


***Entity***: Product Order
***Extract***: Find all unique values for "Order ID", "Product ID", "Sales", "Quality", "Discount"  and "Profit" Columns from superstore dataset
***Transform***:  Export the data as CSV
***Load***:  Load the csv file into MySQL workbench


***Entity***: Metadata (Returns)
***Extract***: Find all unique values for "Order ID", "Customer ID", "Region"  Column from superstore dataset
***Transform***:
1.  Created a list with return reasons
2.  Created an empty column "Reason"
3.  Populated the reason column with a random selection from the list of reasons
4.  Created an empty column "Date" and populated it with today's date as a sample return date
5.  Created an empty column "Condition" and populated it with 'Good' as a sample condition
6.  Replaced Region with its corresponding Region ID

***Load***:  Load the csv file into MySQL workbench

## 4.3 Lab Exercise 3 - submission C

### 4.3.1 Introduction

In this exercise, the Data Analytics department (Group 6) was tasked with creating Operational and Executive reports for communicating useful insights regarding regional sales-related information to the Regional Sales Managers and Chief Operating Officer (COO) of the **Future Shoppers Superstore**. Future Shoppers Superstore is an e-commerce brand that offers a range of furniture, office supplies, and technology. They market products to three customer segments: Consumers, Home Offices, and Corporate companies.

The following report includes the reporting context and shows the final reports that were created with the aid of SQL for querying the superstore's database, which resulted in the production of two CSV files (Operational and Executive information). These CSV files were further formatted in Google Sheets to produce a more professional and organized appearance.

### 4.3.2 Reporting Context

We created Operational and Executive reports for the Regional Managers and COO of **Future Shoppers Superstore** who are preparing for the new quarter's sales (Q2) of 2020.

To support the process, we examined the profits and the product trends for April 2019 as well as comparisons for Q1 and Q2 of 2018 and 2019. These reports will aid the administrative officers in making decisions regarding products with the most and least profits for the regional branches.

Going forward, the Operational report will be utilized by the regional operational teams for future monitoring and reference. Similarly, the Executive report will assist the executive team in current and future decision-making.

### Assumptions

- Monetary values are represented in US dollars.
- The customer's 'Region' was associated with the 'Region' of the product
- All rows with duplicate order IDs were dropped from the database.
- Negative profits were retained in the database and were interpreted as a loss.
- Negative profits were due to discounts and products being sold at a price lower than their original cost or big discounts due to loyalty programs.
- Products that have a negative profit on discount were either purchased through store loyalty points or during a clearance sale to clear out inventory.
- Discounts were expressed as the percentage of the sales that were reduced (expressed as a ratio)
- The 'Order ID' matching between 'Orders' and 'Returns' sheets determines which specific products were returned. Consequently, when a product is returned, we can assume that all the products from that order were returned and are available for resale.
- Each region aims to hit its business goal of a 2% increase in gross profit per year.
- Numbers formatted in the colour "Red" correspond to poor performance (e.g. negative profits or region did not meet the target KPI)

### 4.3.3 Fields Utilized

**Operational Report**

The Operational report included the following fields:

- **Region** - The list of regions (West, East, Central, South) where products are sold
- **Segment** - The various customer segments (Consumer, Corporate, Home Office)
- **Category** - The categories of products (Office Supplies, Furniture, Technology)

*Calculated Fields*

- **(#) Orders** - The number of unique order IDs
- **(#) Quantity** - The sum of total number of products sold per product category
- **Sales (USD $)** - The revenue generated from products across various categories
- **Profits (USD $)** - The sum of the profits generated from products across various categories
- **Order Return Rate (%)**
  - The ratio of products returned to the quantity sold per product category
  - **Formula**: Number of orders returned / (#) Total Orders.
- **(#) Discounted Products** - The total number of products that were offered discounts per product category.
- **Average Transaction Value (USD $)**
  - the average amount a customer spends on a single order.
  - **Formula**: Sales (USD $) / (#) Orders.

**Executive Report**

The Executive report included the following fields:
- **Region** - The list of regions (West, East, Central, South) where products are sold

*Calculated Fields*

- **2019 Gross Profits** (**USD $**)
  - **Quarter 1 (Q1)** - The total profits per region for Q1 of 2019
  - **Quarter 2 (Q2)** - The total profits per region for Q2 of 2019
  - **Q2 vs Q1 (%)** - The ratio of Q2 to Q1 gross profits for 2019
    - If the previous Quarter had negative profits, this will be blank as a negative percent does not make much sense in terms of profit

- **KPI: Q2 (%) To Target (%) - (*2% increase*)**
  - The proportion of the quarter's target gross profit accumulated.
  - **Formula**:  Q2 Gross Profit / Q2 Target Gross Profit
  - Values in red means that has not reached the target for Q2
- **Q2 Gross Profits**
  - **2018 (USD $)** - The total profits per region for Q2 of 2018
  - **2019 vs 2018 (%)** - The ratio of 2019 to 2018 gross profits for Q2

- **KPI: 2019 (%) To Target (*2% increase*)**
  - The proportion of the year's target gross profit accumulated
  - **Formula:** 2019 Gross Profit / 2019 Target Gross Profit
  - Values marked red have not reached 50% to target