# PROJECT - 2

**(BUAN 6320: Database Foundations for Business Analytics)**

## Topic: Retail Data of Turkish Sector

| Abhishek Arya | axa220149@utdallas.edu |
|---|---|
| Manan Mistry | mxm220108@utdallas.edu |
| Tanishka Patel | txp210040@utdallas.edu |
| Het Shah | hrs220002@utdallas.edu |
| Shekhar Vashist | sxv220008@utdallas.edu |

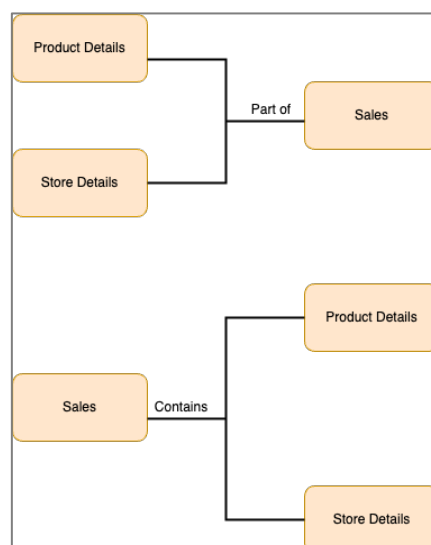**Project made under the guidance of: Farzad Kamalzadeh**

# SECTION 1: DESIGN A DATABASE

Queries 1 and 2 are completed in project 1.

### 3. SCHEMA DESIGN

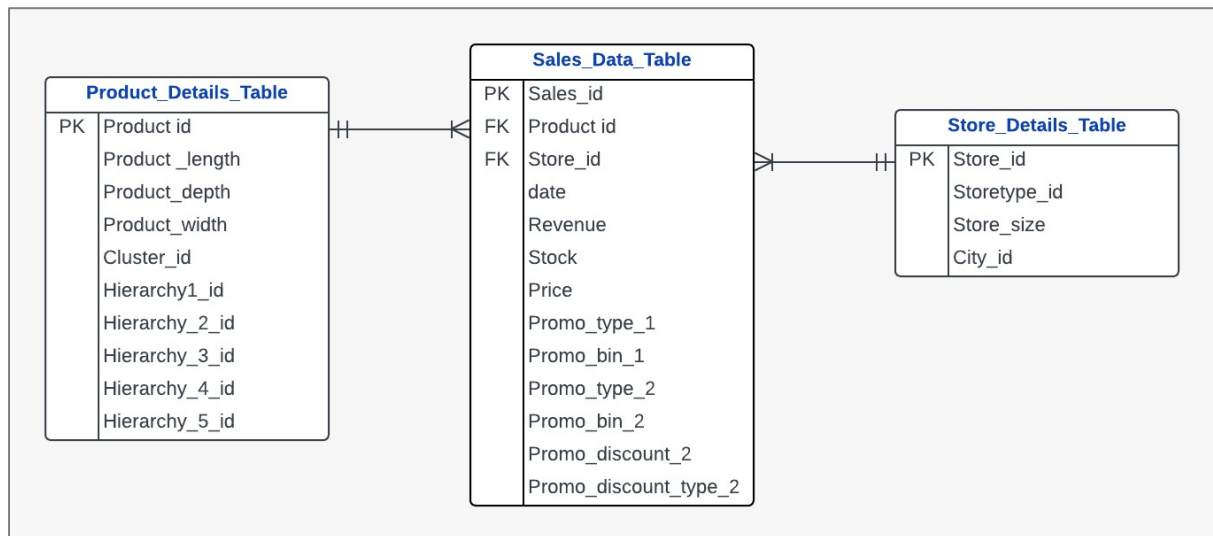a. Find entities, their attributes, their primary keys, and relationships between them.

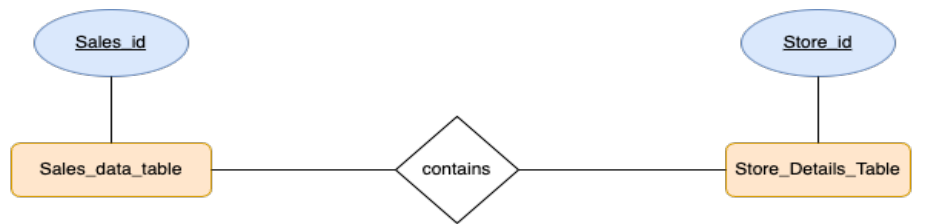| ENTITIES | ATTRIBUTES | PRIMARY KEYS |
|---|---|---|
| Product_Details_Table | Product_id<br>Product_length<br>Product_width<br>Cluster_id<br>Hierarchy1_id<br>Hierarchy2_id<br>Hierarchy3_id<br>Hierarchy4_id<br>Hierarchy5_id | Product_id |
| Store_Details_Table | Store_id<br>Storetype_id<br>Store_size<br>City_id | Store_id |
| Sales_Data_Table | Sales_id<br>Store_id<br>Product_id<br>Date<br>Sales<br>Revenue<br>Stock<br>Price<br>Promo_type_1<br>Promo_bin_1<br>Promo_type_2<br>Promo_bin_2 | Sales_id |

b. Model all the constraints you believe should be there in your schema.

1. Each Sales record must contain only one Product Details Table record and one Store Details Table record.
2. Each Store Details Table can be part of 1 or more Sales records
3. Each Product Details Table can be part of 1 or more Sales records.

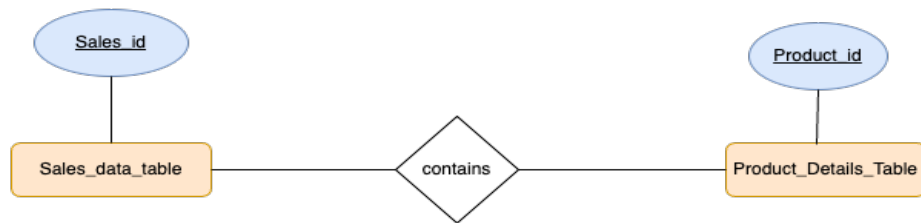c. Draw and ER diagram of your dataset.

d. Translate your ER diagram into relations.



| Sales Data |
|---|
| (sales_id, product_id,<br> store_id,<br>date,<br>sales,<br>revenue,<br>stock,<br>price,<br>promotype_1,<br>promobin_1,<br>promotype_2,<br>promobin_2) |

| Contains |
|---|
| sales_id ,<br><br>store_id |

| Store Details |
|---|
| (Store_id,<br>storetype_id,<br>store_size,<br>city_id) |



| Sales Data |
|---|
| (sales_id, product_id,<br> store_id,<br>date,<br>sales,<br>revenue,<br>stock,<br>price,<br>promotype_1,<br>promobin_1,<br>promotype_2,<br>promobin_2) |

| Contains |
|---|
| sales_id ,<br><br>product_id |

| Product Details |
|---|
| ( product_id,<br>product_length,<br>product_depth,<br>product_width,<br>cluster_id,<br>hierarchy1_id,<br>hierarchy2_id,<br>hierarchy3_id,<br>hierarchy4_id,<br>hierarchy5_id ) |

# 4. SCHEMA NORMALIZATION

## a. Functional Dependencies:

**Sales Table**

**sales_id** → product_id, store_id, date

product_id, store_id, date → sales_id (product id and store id gives a unique sales id)

**sales_id** → product_id, store_id, date, sales, revenue, stock, price, promo_type_1, promo_bin_1, promo_type_2

**Product Table**

**hierarchy5_id** → hierarchy4_id

**hierarchy4_id** → hierarchy3_id

**hierarchy3_id** → hierarchy2_id

**hierarchy2_id** → hierarchy1_id

**product_id** → product_length, product_width, product_depth, cluster_id,

hierarchy1_id, hierarchy2_id, hierarchy3_id, hierarchy4_id, hierarchy5_id

**hierarchy5_id** → hierarchy4_id, hierarchy3_id, hierarchy2_id, hierarchy1_id

**Store Details**

**store_id** → storetype_id, store_size, city_id

## b. Minimal Keys:

Currently, sales_id aside, the primary key for sales {product_id, store_id, date} and that is minimal because subsets of this are not candidate keys. The same goes for product_id is the primary key for product details and store_id is the primary key for store details.

Each minimal key is a unique combination of attributes that can identify a tuple in a relation without any redundancy. Identifying minimal keys is essential for enforcing data integrity and avoiding redundancy in the database schema.

## c. Check if your schema is in BCNF (Boyce-Codd Normal Form)

A relation is in BCNF if, for every non-trivial functional dependency X -> Y, X is a super key. Before checking if the schema is in BCNF (Boyce-Codd Normal Form), let us recap the functional dependencies which we identified for each table: The functional dependencies indicate that the schema is in BCNF, as the left side of each functional dependency represents the primary key of each entity.

### Sales Table

sales_id → product_id, store_id, date

product_id, store_id, date → sales_id

sales_id → product_id, store_id, date, sales, revenue, stock, price, promo_type_1, promo_bin_1, promo_type_2

promo_id -> promo_type_1, promo_bin_1, promo_type_2, promo_bin_2, promo_discount_2, promo_discount_type_2

**(Created separate tables for promotion for sets of values that apply to multiple rows.)**

### Product Details

{product_id, product_length, product_width, product_depth, cluster_id, hierarchy1_id,

hierarchy2_id, hierarchy3_id, hierarchy4_id, hierarchy5_id}

product_id → product_length, product_width, product_depth, cluster_id, hierarchy1_id,

hierarchy2_id, hierarchy3_id, hierarchy4_id, hierarchy5_id → **Follows BCNF**

 hierarchy5_id gives the following ids in the promotion table.

{hierarchy5_id} → {hierarchy4_id, hierarchy3_id, hierarchy2_id, hierarchy1_id} →

{product_id, product_length, product_width, product_depth, cluster_id, hierarchy5_id}, {hierarchy1_id, hierarchy2_id, hierarchy3_id, hierarchy4_id, hierarchy5_id}

**Violates BCNF because it is a transitive dependency; the elements on the right depend on the element hierarchy5_id, which is not a key**.

hierarchy4_id → hierarchy3_id

hierarchy3_id → hierarchy2_id

hierarchy2_id → hierarchy1_id

{product_id, product_length, product_width, product_depth, cluster_id, hierarchy5_id},

{hierarchy1_id, hierarchy2_id}, {hierarchy2_id, hierarchy3_id}, {hierarchy3_id,

hierarchy4_id}, {hierarchy4_id, hierarchy5_id}.

### Store Details:

store_id → storetype_id, store_size, city_id → **Follows BCNF**


## d. The <u>whole schema is not in BCNF</u> because the product details section has transitive dependency. So, we will decompose:

1. Sales table into Promotion and Sales.
2. ProductDetails table into productDetails and Hierarchy table.

Now, we checked if the schema is in BCNF. A schema is in BCNF if for every non-trivial functional dependency, the left-hand side is a super key. A super key is a set of one or more attributes that can uniquely identify each tuple in a relation.

• In our case, we have the following tables in our schema:

  1. Sales table
  2. Product Details table
  3. Store table
  4. Hierarchy Table
  5. Promotion Table

• Each of these tables satisfies the BCNF conditions, meaning that there are no non-trivial functional dependencies that violate the BCNF rules. As a result, there is no need to decompose these tables further.

The schema is now normalized to an acceptable degree, ensuring minimal redundancy, and reducing the risk of update anomalies.

## Decomposition with the minimal keys:

### Sales Table

● sales_id → product_id, store_id, date, sales, revenue, stock, price, promo_id

### Product Details Table

● product_id → product_length, product_width, product_depth, cluster_id, hierarchy5_id

(Product table has a unique value of various products along with its dimensions.)
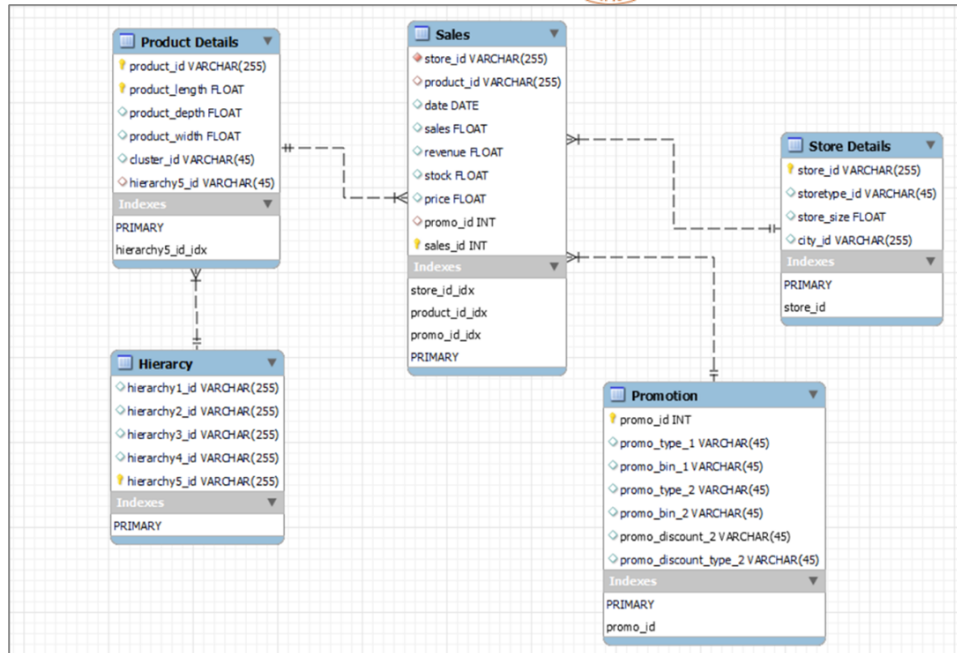
### Store Details Table

● store_id → storetype_id, store_size, city_id (This table provides the store type, its size and respective city.)

### Hierarchy Table

● hierarchy5_id → hierarchy4_id, hierarchy3_id, hierarchy2_id, hierarchy1_id

### Promotion Table

●promo_id → promo_type_1, promo_bin_1, promo_type_2, promo_bin_2, promo_discount_2, promo_discount_type_2 (This table now contains all the details of the promotion table along with its discount details.)

## 5. CREATING THE DATA BASE WITH THE NEW SCHEMA

- Creating a database from the new ER Diagram had some minor issues,
  - While creating product table following issue came:
    "Table storage engine for <TABLE> doesn't have this option on order by query (ERROR 1031)"
    **Resolution:** We resolved this by changing the ROW_FORMAT = FIXED to DEFAULT in the auto created code by the wizard.

  - Later while creating store details table we faced the following error:
    "Error 1064: You have an error in your SQL syntax; check the manual that corresponds to your MySQL server version for the right syntax to use near ') INVISIBLE)
    ENGINE = InnoDB
    KEY_BLOCK_SIZE = 4' at line 10"
    **Resolution:** This was resolved by correcting the index statement by removing INVISIBLE from the statement.

  - Later while creating Promotion table we faced the following error:
    "Error 1064: You have an error in your SQL syntax; check the manual that corresponds to your MySQL server version for the right syntax to use near "promo_id')
    )ENGINE = InnoDB' at line 13"
    **Resolution:** We resolved this by correcting the index statement by removing VISIBLE from the statement.

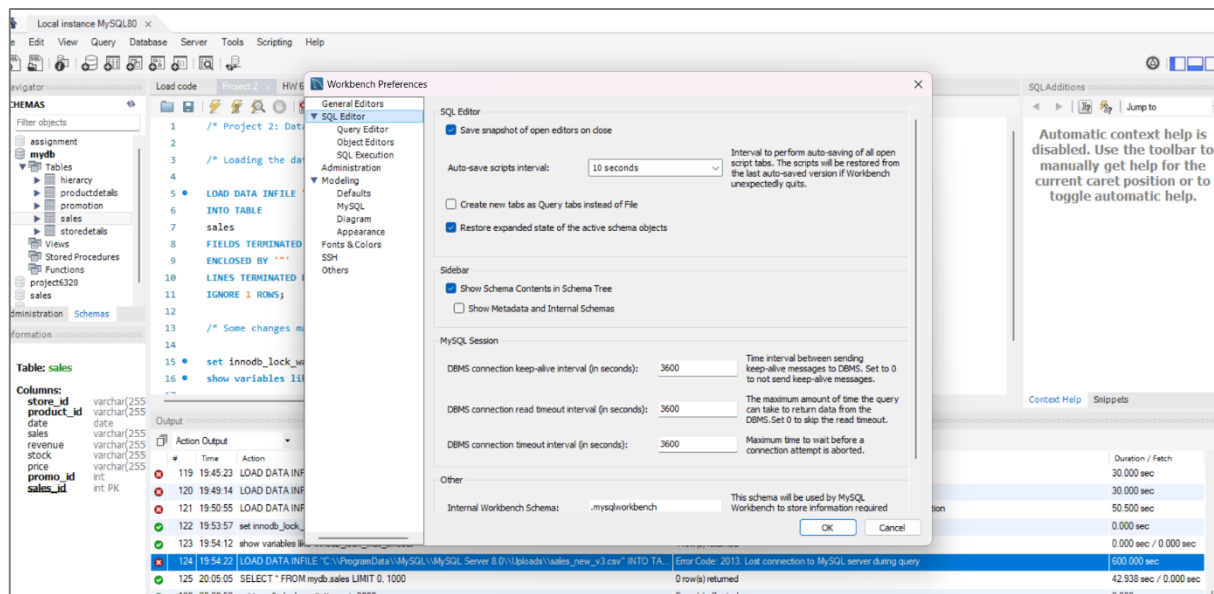## 6. IMPORTING THE DATA INTO THE TABLES

Before importing the data into MySQL workbench we did come data manipulation on R such as splitting the data into tables and added the unique identifiers in some tables.

*OTHER ERRORS THAT WE FACED:*

1. We faced administrative and infile restriction while importing the error.
   **Resolution:** This issue was resolved by changing the location of data to uploads folder in MySQL folder of the system where it was given permission to load the files. Also, we had to use the "\\" in the path section of our code to resolve this issue.

2. "Error Code: 1262. Row 1 was truncated; it contained more data than there were input columns"
   **Resolution:** This was resolved by inserting the data via import wizard

3. "Error Code: 1265. Data truncated for column 'price' at row 1"
   While making changes with data in R, for importing there were NA values and they created problem while importing the data in MySQL workbench.
   **Resolution:** We resolved this issue by changing the numerical NA argument to 0 while exporting the data file from R and then loading the data in MySQL.

4. Error Code: 1452. Cannot add or update a child row: a foreign key constraint fails (`mydb`.`sales`, CONSTRAINT `product_id` FOREIGN KEY (`product_id`) REFERENCES `product details` (`product_id`))
   While importing the sales data we had this issue which describes that the product_id does not exist in the productDetail table and hence it does not meet the foreign key constraints.
   **Resolution:** We resolved this error by checking the productDetails table and saw that SQL was removing some records with few NA entries and hence some records of product_id were missing because of the import by import wizard. We changed the NA in our data set to blank and then again imported the data in the table which eventually resolved the error.

5. Faced following timeout errors while uploading the sales data with load command.
   Error Code: 1205. Lock wait timeout exceeded; try restarting transaction
   Error Code: 2013. Lost connection to MySQL server during query
   **Resolution:** We resolved this by increasing the timeout interval and also by changing the "innodb_lock_wait_timeout" limit.

# SECTION 2: DATA CLEANING AND DATABASE TESTING

1.  ***For each table in your database, check all the columns and the values they contain***

    After updating the new schema we have the following tables and attributes (along with their data types)

| Table Name | Attribute | Attribute type |
|---|---|---|
| StoreDetails | store_id | VARCHAR(255) |
| | storetype_id | VARCHAR(45) |
| | store_size | FLOAT |
| | city_id | VARCHAR(255) |
| | | |
| Promotion | promo_id | VARCHAR(45) |
| | promo_type_1 | VARCHAR(45) |
| | promo_bin_1 | VARCHAR(45) |
| | promo_type_2 | VARCHAR(45) |
| | promo_bin_2 | VARCHAR(45) |
| | promo_discount_2 | VARCHAR(45) |
| | promo_discount_type_2 | VARCHAR(45) |
| | | |
| ProductDetails | product_id | VARCHAR(255) |
| | product_length | FLOAT |
| | product_depth | FLOAT |
| | product_width | FLOAT |

| Table Name | Attribute | Attribute type |
|---|---|---|
| | cluster_id | VARCHAR(45) |
| | hierarchy5_id | VARCHAR(45) |
| | | |
| Hierarchy | hierarchy1_id | VARCHAR(255) |
| | hierarchy2_id | VARCHAR(255) |
| | hierarchy3_id | VARCHAR(255) |
| | hierarchy4_id | VARCHAR(255) |
| | hierarchy5_id | VARCHAR(255) |
| | | |
| Sales | store_id | VARCHAR(255) |
| | product_id | VARCHAR(255) |
| | date | DATE |
| | sales | FLOAT |
| | revenue | FLOAT |
| | stock | FLOAT |
| | price | FLOAT |
| | promo_id | INT |
| | sales_id | INT |

## 2. Check them against the information you found in step 3 of project 1

While uploading the sales data in our table in MySQL workbench we had to do some modification with NA values and replace it with 0 for importing.The Statistics for the numerical attributes are approximately the same what mentioned in step 3 of project 1. The values of these statistical measures are mentioned in the query log attached below.

| Table | Variable | Mean | Standard Deviation |
|---|---|---|---|
| Sales | Sales | 0.41 | 14.21 |
| Sales | Revenue | 1.94 | 38.66 |
| Sales | Stock | 14.91 | 36.32 |
| Sales | Price | 15.13 | 31.93 |

## 3. Checking for data discrepancies:

a. **Discrepancy:** In productDetails table we have some records showing product length, depth, and width as 0, which we are assuming are blanks because while importing the data the NAs were converted as zeros.
**Resolution:** We will replace these zeros with the average values of the attributes, and we will be taking the average values with respect to non zeros

b. **Discrepancy:** In promotion table while importing the data we imported the NAs in product_discount_2 column
**Resolution:** We changed these values to NULL by using update command

### 4. Make sure all the values of these columns are from the same type (all numeric)

After reviewing all the character variables in all the data tables, everything is correct, and all the columns are from the same data type as per our schema.


**QUERY LOGS FOR IMPORT AND DATA CHECK / APPENDIX:**

```
/* Project 2: Data Base Foundation for Business Analytics */

/* Loading the data in sales table*/

LOAD DATA INFILE "C:\\ProgramData\\MySQL\\MySQL Server
8.0\\Uploads\\sales_new_v5.csv"
INTO TABLE
sales
FIELDS TERMINATED BY ','
ENCLOSED BY ''
LINES TERMINATED BY '\n'
IGNORE 1 ROWS;

/* Some changes made temporarily to change the timeout time to upload the data */

set innodb_lock_wait_timeout=3600;
show variables like 'innodb_lock_wait_timeout';

/*
===================================================================
====================================*/

/* Checking for data discrepancies*/

/* Checking for data discrepancies in productdetails table*/

SELECT *
FROM productdetails;

SELECT *
FROM productdetails
WHERE product_length = 0;

UPDATE productdetails
SET product_length = (
SELECT AVG(product_length)
FROM productdetails WHERE product_length != 0)
WHERE product_length = 0;


SELECT AVG(product_length)
```

```
FROM productdetails
WHERE product_length <> 0;

UPDATE productdetails
SET product_length = 7.243
WHERE product_length = 0;

SELECT AVG(product_width)
FROM productdetails
WHERE product_width <> 0;

UPDATE productdetails
SET product_width = 13.474
WHERE product_width = 0;

SELECT AVG(product_depth)
FROM productdetails
WHERE product_depth <> 0;

UPDATE productdetails
SET product_depth = 18.478
WHERE product_depth = 0;

/* Checking for data discrepancies in hierarchy table: No numerical variables in this table*/
/* Checking for leading and trailing spces in hierarchy table*/

SELECT *
FROM hierarcy
WHERE hierarchy1_id LIKE ' %' or hierarchy1_id LIKE '% ';

SELECT *
FROM hierarcy
WHERE hierarchy2_id LIKE ' %' or hierarchy2_id LIKE '% ';

SELECT *
FROM hierarcy
WHERE hierarchy3_id LIKE ' %' or hierarchy3_id LIKE '% ';

SELECT *
FROM hierarcy
WHERE hierarchy4_id LIKE ' %' or hierarchy4_id LIKE '% ';

SELECT *
FROM hierarcy
WHERE hierarchy5_id LIKE ' %' or hierarchy5_id LIKE '% ';


/* Checking for data discrepancies in promotion table*/

SELECT *
```

```
FROM promotion;

UPDATE promotion
SET promo_discount_2 = NULL
WHERE promo_discount_2 = "NA";


/* Checking for data discrepancies in storedetails table*/

SELECT *
FROM storedetails;

SELECT *
FROM storedetails
WHERE store_size IS NULL;

/* storedetail table seems fine as there are no missing values in the numeric data column i.e.
store_size */


/* Chcking the discrepancies in sales table */

SELECT AVG(sales)
FROM sales;

/* mean = 0.41 */

SELECT stddev(sales)
FROM sales;

/* Standard Dev = 14.21 */

SELECT AVG(revenue)
FROM sales;
/* mean = 1.94 */
SELECT stddev(revenue)
FROM sales;
/* Standard Dev = 38.66 */
SELECT AVG(stock)
FROM sales;
/* mean = 14.91 */
SELECT stddev(stock)
FROM sales;
/* Standard Dev = 36.32 */
SELECT AVG(price)
FROM sales;
/* mean = 15.13 */
SELECT stddev(price)
FROM sales;
/* Standard Dev = 31.93 */
```

## DATA CHECK FOR CHARACTER VARIABLES

   *1.  Check them against the information you found in step 3 of project 1.*
For the character columns, a few of the columns in the sales table contained large number of null values and removed those majorly from the promotion column. As shown in section 2-part 1 table, All the promotion columns are of character table.

```
2 ●     SELECT * FROM mydb.sales WHERE promo_id IS NULL;
3 ●     UPDATE mydb.sales SET promo_id=0 WHERE sales IS NULL;
4 ●     SELECT * FROM mydb.sales WHERE promo_id IS NULL;
```

| store_id | product_id | date | sales | revenue | stock | price | promo_id | sales_id |
|---|---|---|---|---|---|---|---|---|
| NULL | NULL | NULL | NULL | NULL | NULL | NULL | NULL | NULL |

   *2.  Checking for data discrepancies:*
For removing discrepancies, separated out the unused products from the product_details table and deleted incomplete or missing data from the sales table. The following query was performed in  order to remove discrepancies.

   a.  **Discrepancy:** In promo_bin_1 column contains missing values in the sales table which are of character type.
      **Resolution:** For removing the white space and null values in the promo_bin_1 column we are separating the whole column and making a new table for it. So, by doing these, we can remove the anomalies in the sales table

   b.  **Discrepancy:** cluster_id column contains missing values in the product_hierarchy table which are of character type.
      **Resolution:** For removing the white space in the cluster_id column, we are querying in the particular cluster_id column to remove all the null values. So, by doing these, we can remove the anomalies in the product_hierarchy table.

   3.  *Make sure all the values of these columns are from the same type (all character)*
      After reviewing all the numerical variables in all the data tables, everything is correct, and all the columns are from the same data type as per our schema.

**QUERY LOGS FOR CHARATER VARIABLES DATA CHECK / APPENDIX:**

CREATE TABLE unused_products as SELECT productdetails.*
FROM productdetails
LEFT JOIN sales
ON productdetails.product_id = sales.product_id
where sales.product_id IS NULL;

/* to delete the null values */
DELETE FROM sales where (sales.sales IS NULL) OR (sales.revenue is NULL) OR (sales.stock IS
NULL)

## APPLYING JOINS TO CHECK WHETHER JOINS WORK PROPERLY IN OUR NEW SCHEMA:

**Join 1:**



**Join 2:**



**Join 3:**

## QUERIES TO CHECK FOR CONSTRAINTS:

Constraints are also working as shown in the below queries.

### Query 1



Error Code: 1452. Cannot add or update a child row: a foreign key constraint fails (`mydb`.`sales`, CONSTRAINT `product_id` FOREIGN KEY (`product_id`) REFERENCES `productdetails` (`product_id`))

### Query 2

Error Code: 1452. Cannot add or update a child row: a foreign key constraint fails (`mydb`.`productdetails`, CONSTRAINT `hierarchy5_id` FOREIGN KEY (`hierarchy5_id`) REFERENCES `hierarcy` (`hierarchy5_id`))

**Query 3**

```
202
203 •   DELETE FROM hierarcy
204     WHERE  hierarchy5_id = "H0315080028";
205
```

Context Help  Snippets

Output

Action Output         ▾

| # | Time | Action | Message | Duration / Fetch |
|---|------|--------|---------|------------------|
| ⊗ 1 | 21:11:23 | DELETE FROM hierarcy WHERE hierarchy5_id = "H0315080028" | Error Code: 1451. Cannot delete or update a parent row: a foreign key constraint fails ('mydb'.'productdetails', CON... | 0.000 sec |

Error Code: 1451. Cannot delete or update a parent row: a foreign key constraint fails (`mydb`.`productdetails`, CONSTRAINT `hierarchy5_id` FOREIGN KEY (`hierarchy5_id`) REFERENCES `hierarcy` (`hierarchy5_id`))

*** End of Document***