# CLEANING AND TRANSFORMING THE DATASET USING POWER QUERY

## Introduction

Source File: malltransactions_mockdataset.csv

Description of the file

- A mocked data of well-known malls in the capital region of the Philippines

- Created from mockaroo.com

Dimensions

- transaction_id
- date
- mall
- city
- gender
- age
- product category
- price
- quantity
- discount
- payment

Business Requirements

- All Ages less than 21 or 60 and above will have 20% discount
- For undeclared ages, replaced it with average age per gender and per mall.
- For undeclared gender, replace it with "Undeclared"
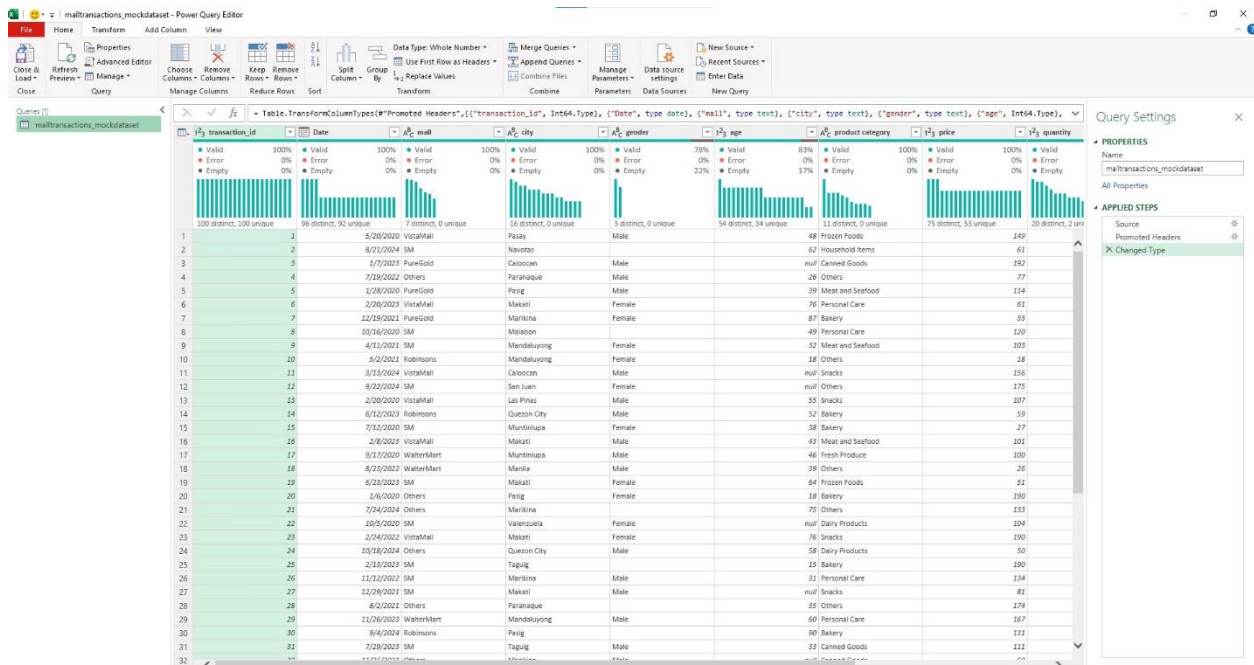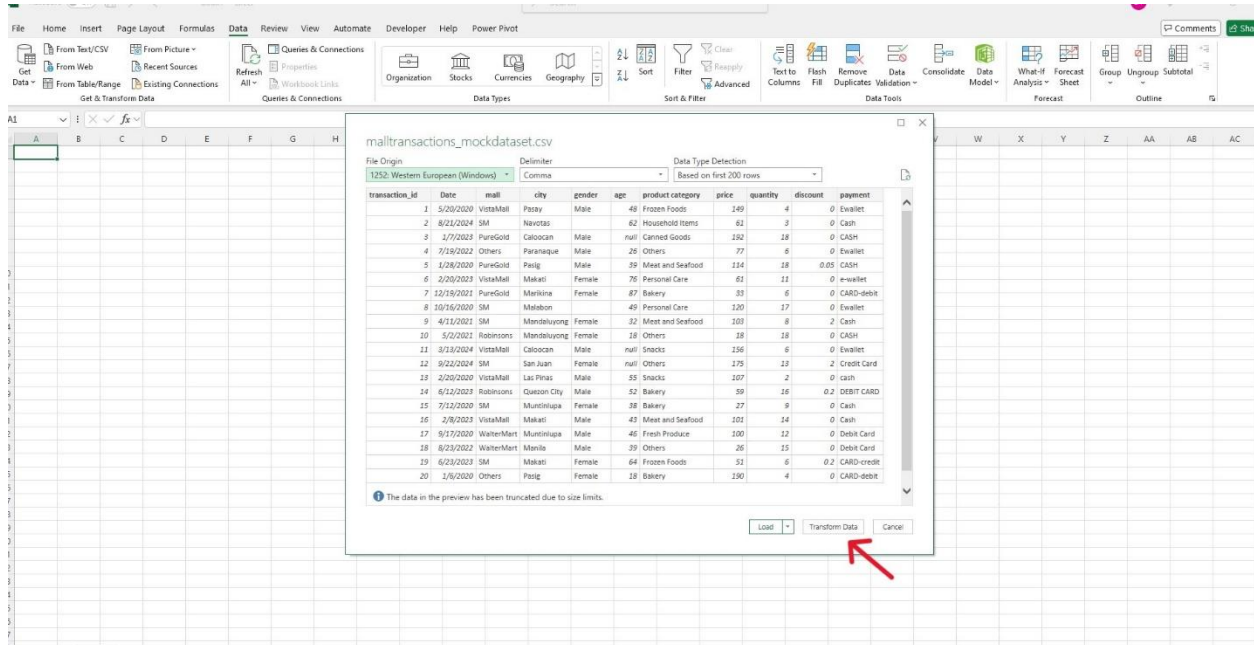
Business Notes

- There are data entry mistakes on discounts and payment

## Steps In Cleaning and Transforming the Dataset

1. Open file in Power Query
2. Do an Exploratory Data Analysis of the file
   a. Check for error/null values on each column
   b. Check for misspelled / case inconsistency in the values
   c. Check for Outliers
3. Clean all necessary columns
   a. Correct misspelled / case inconsistency values
   b. Filter or Replace Outliers with standard values
   c. Replace or Remove error/null values
   d. Check if all columns have the right data type and format
4. Save and Close file to be used on Reports/Creating Dashboard

# Step 1: Open file in Power Query

- I use Excel to use the built-in Power Query tool.
- I clicked the Transform Data for the next Step.





# Step 2: Do an Exploratory Data Analysis (EDA)

***On View Ribbon, check the following:***

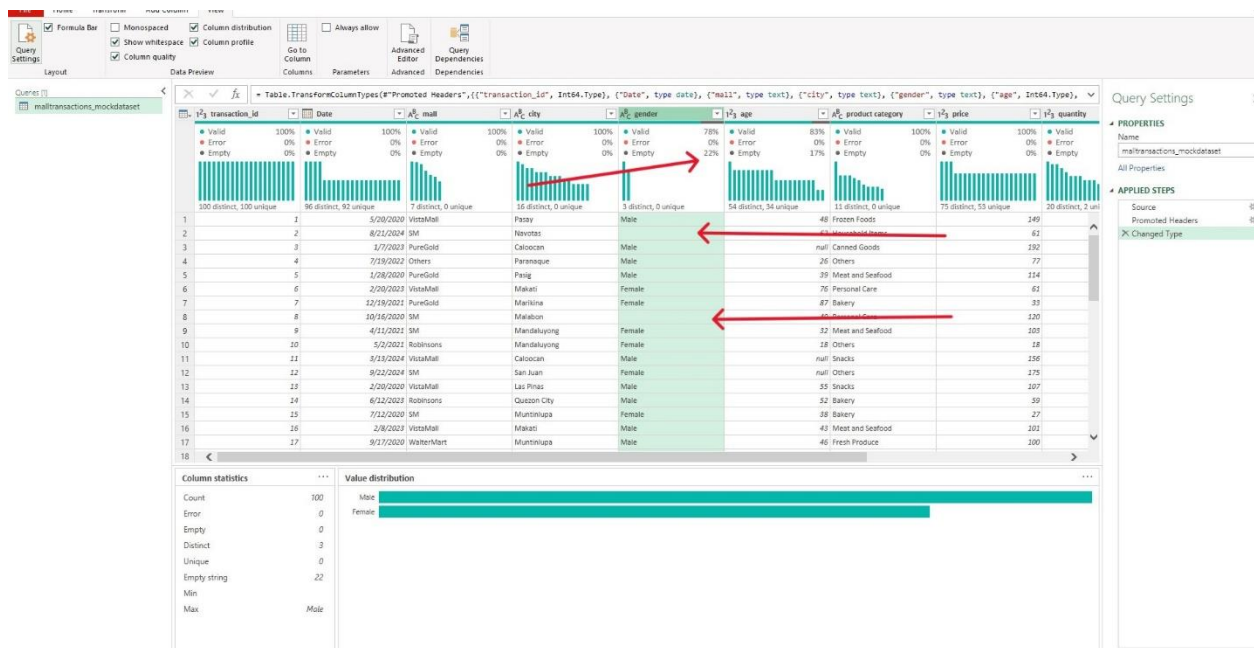- Column Profile

- Column Quality
- Column Distribution



This is used for EDA, checking all the needed things for cleaning and transforming.

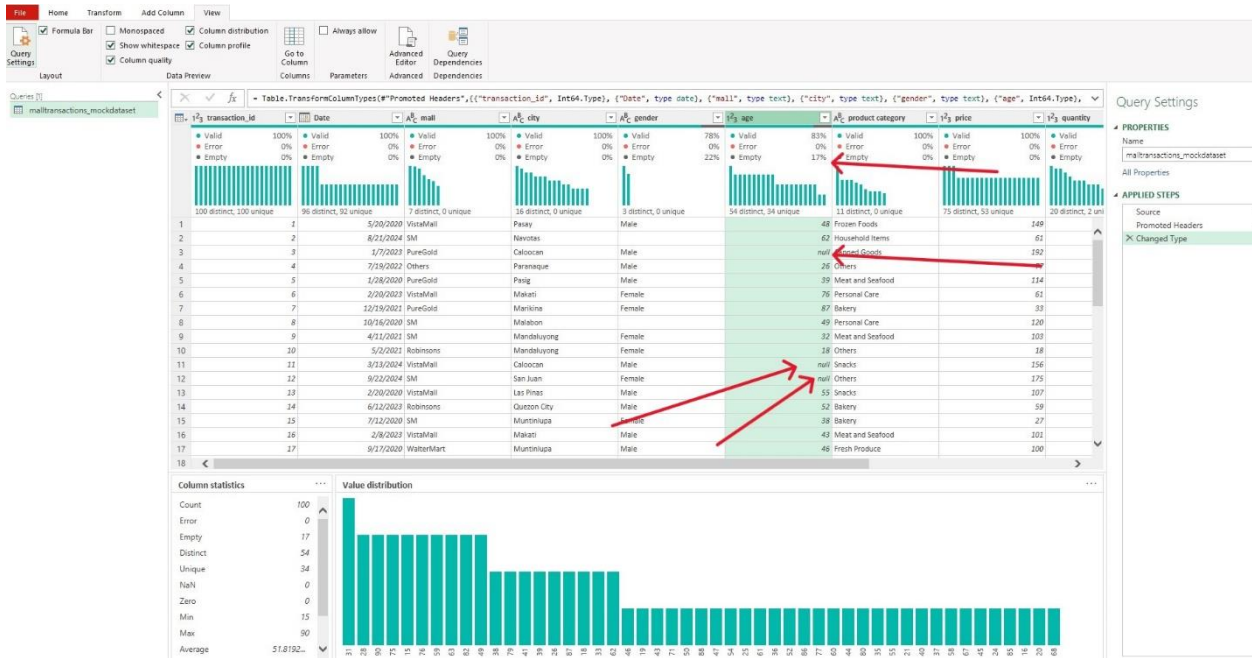## Check all column profiles, column quality, column distribution of each column

### *Gender column*

- Gender column have EMPTY cells which is 22%.
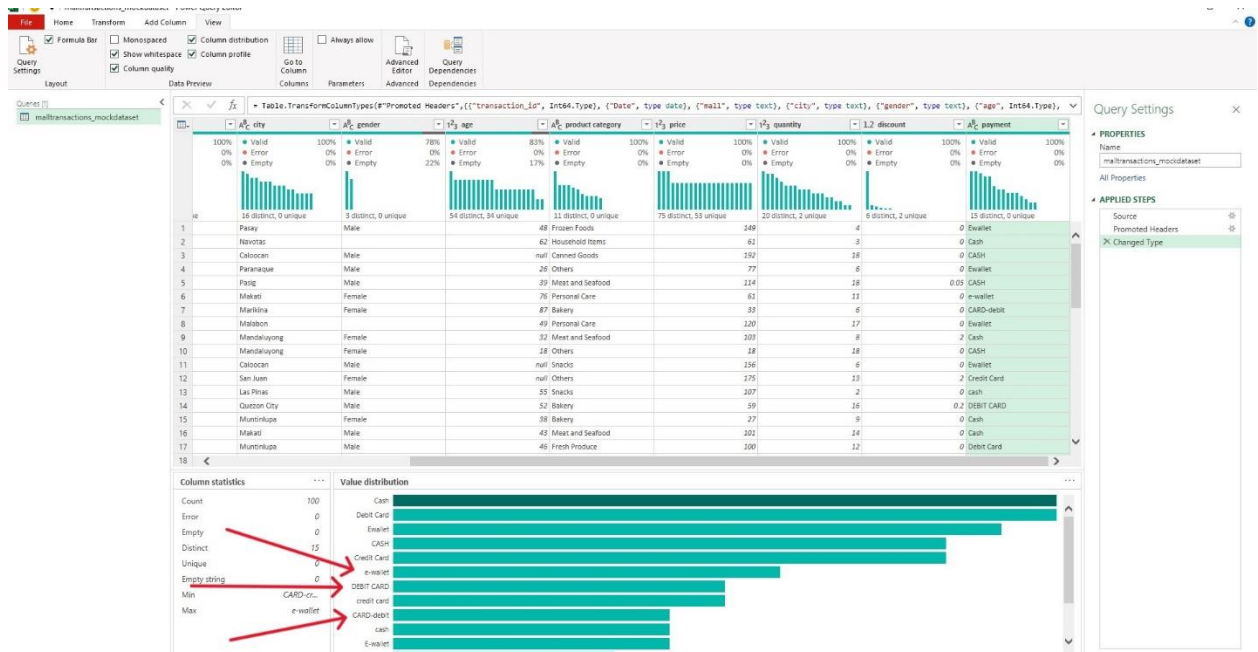


### *Age column*

- Age column has null values and is about 17%.

### Discount column

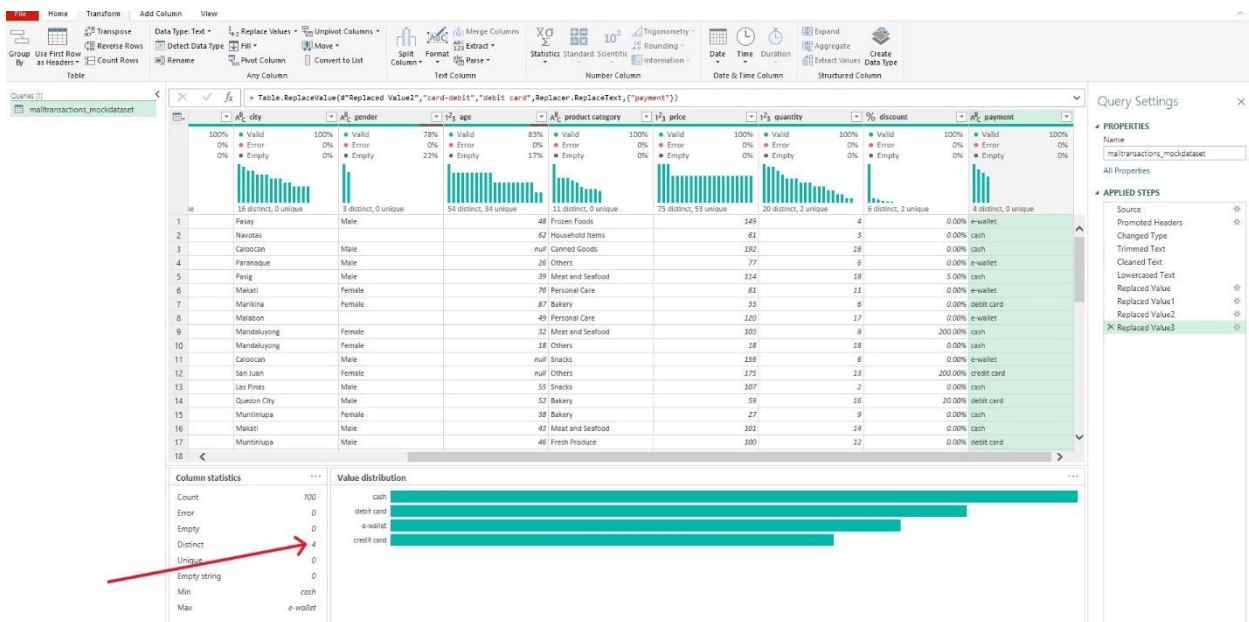- Discount column has negative values and an Outlier. Data should be in percentage format.

### *Payment column*

- Payment column has misspelled / case inconsistency in its values



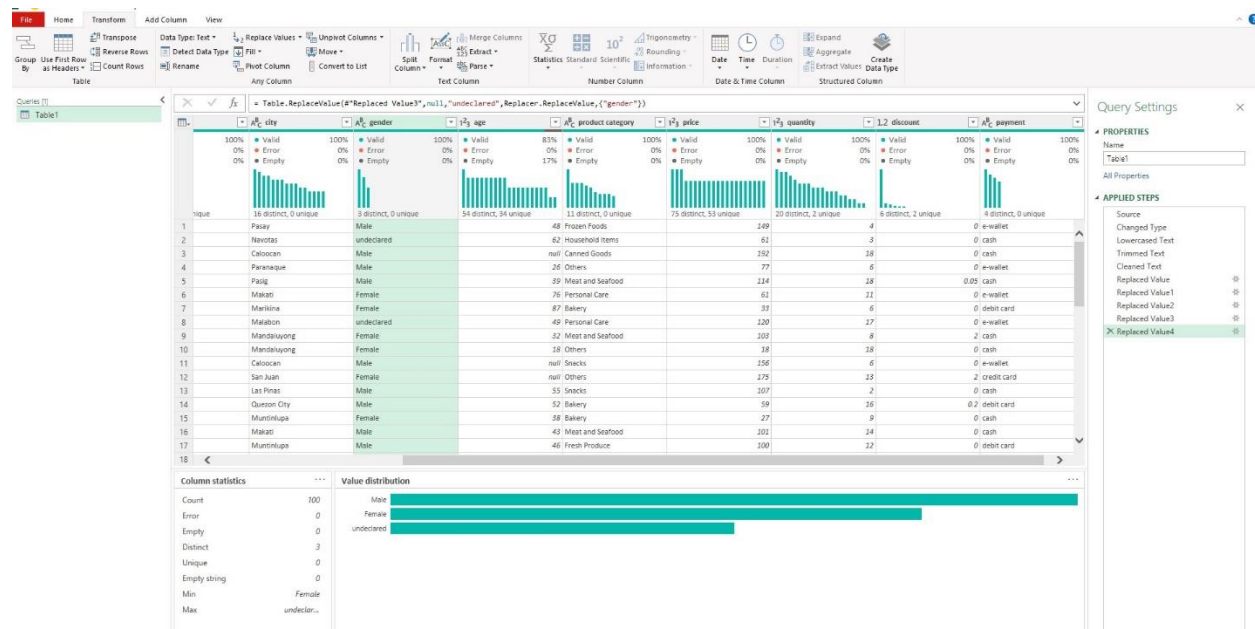## Step 3: Clean all necessary columns

### *Payment Column*

- For Payment Column, I trim, clean, and lower case the values of the column
- Replace ewallet to e-wallet, card-debit to debit card, coins to cash, card-credit to credit card for standardization
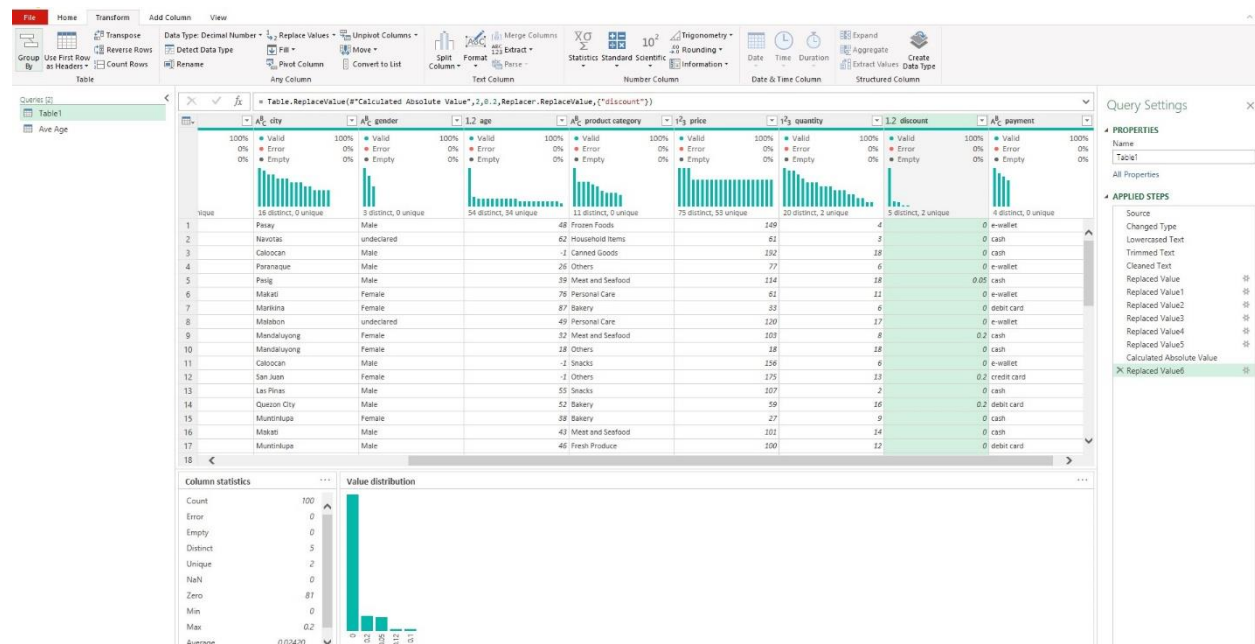- Only 4 distinct values remain

### Gender Column

- For Gender Column, I replaced blank values to undeclared values



### Discount Column

- For Discount Column, I remove all negative values by using the absolute value feature
- Replace the Outlier 2 to .2



- Created a new column "discount_new" that will change values based on the age (All Ages less than 21 or 60 and above will have 20% discount)

- Temporarily, I replace all null values with -1 on the Age column to avoid Errors after manipulation



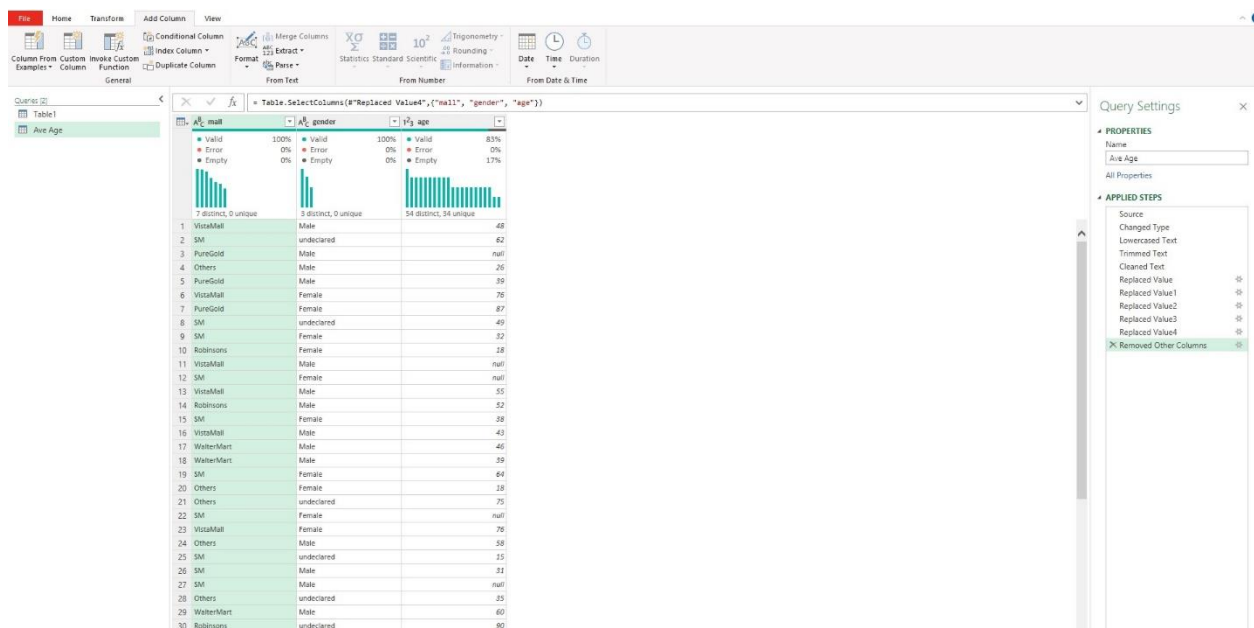- I use conditional column for manipulation for the discount_new column



- Replace the -1 to Null again and delete the old column.
- Replace the format to percentage

### Age Column

- For Age Column, I need to get the average age per gender per mall.
  - I duplicate the query, remove other columns except Age, Mall, and Gender



  - I use the Group BY feature to get the average age

## Group By

Specify the columns to group by and one or more outputs.

○ Basic  ● Advanced

| mall | ▼ |
| gender | ▼ |

Add grouping

| New column name | Operation | Column |
| --- | --- | --- |
| Ave_Age | Average ▼ | age ▼ |

Add aggregation

OK    Cancel

○ Replace the *null* to the average age between Female and Male of WalterMart
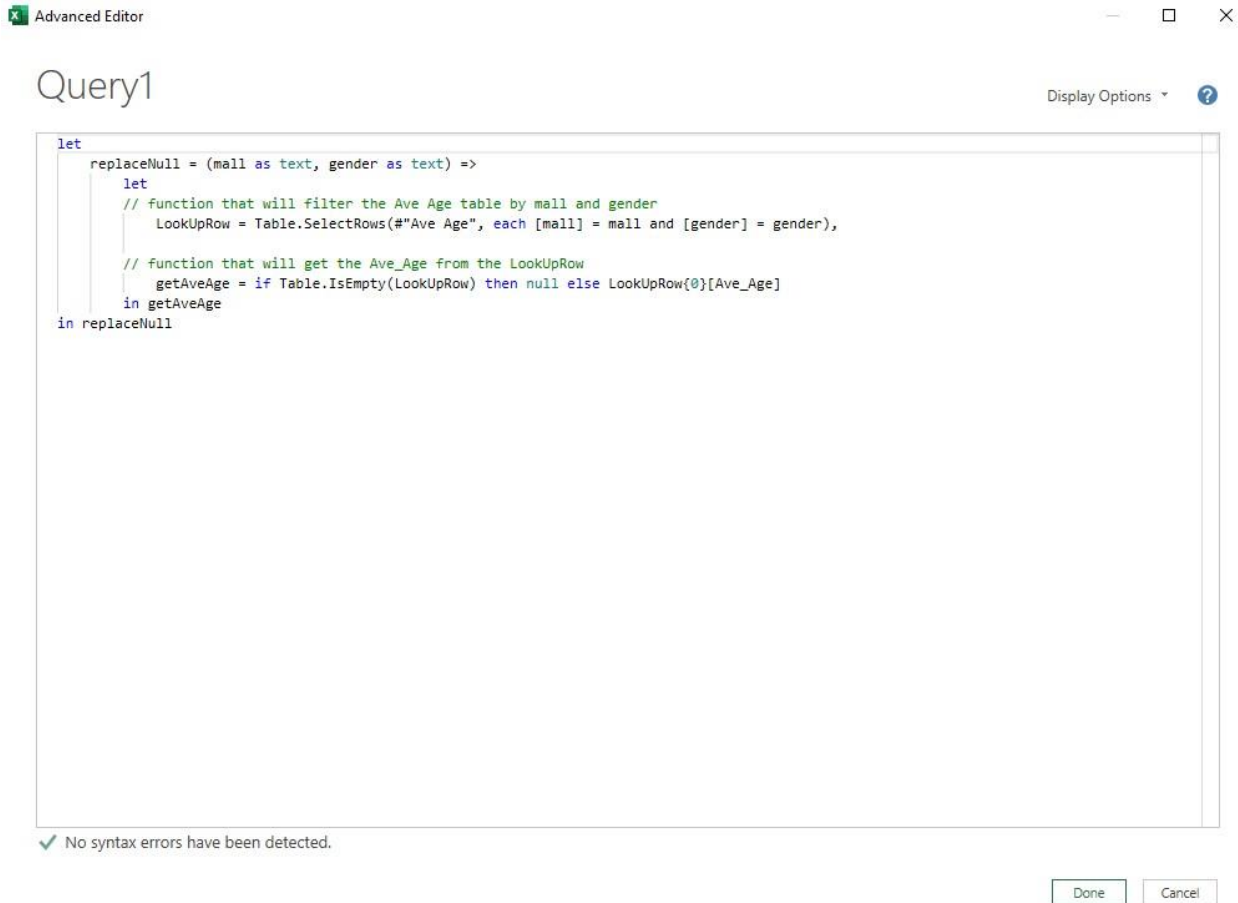
o    This will be the reference table for replacing the null values.

- I'll create a new query and create a custom function that will replace the null values based on the mall and the gender
  - On Advanced Editor, Put this code



  - Rename the function avgAge_function



  - On the main table, add column by custom column and put this code:

## Custom Column

Add a column that is computed from the other columns.

New column name

age_clean

Custom column formula ⓘ

```
= if [age] = null then avgAge_function( [mall], [gender] )
   else [age]
```

Available columns

transaction_id
Date
mall
city
gender
age
product category

<< Insert

Learn about Power Query formulas

✔ No syntax errors have been detected.    OK    Cancel

○ Remove the old age column and change format to whole Number, rename the column to age

## Step 4: Save and Close file to be used on Reports/Creating Dashboard

• Close and Load the file. Save the file for Reports/ Creating Dashboard