# GB CORP DATA ASSESSMENT

## MARINA SAFWAT

**github repo:
https://github.com/WA-on-test2/Sales_Full_Data_Project.git**

# PART A
# DATA QUALITY ISSUES

1. Inconsistency
   - Mixed Upper and Lower cases
   - Mixed Languages e.g. some names are written in Arabic and others in English
   - Emails had some letters in Arabic and had random spaces
   - Formats were so various e.g. Phone no. & Dates even some dates were written in English Except for the month is in Arabic
   - The currency for example some entries had just value and others had the currency stated(Mixed Ar&En)
   - Value Representation was various e.g. Gender&Payment method & Status

2. DUPLICATE ISSUES
   - Columns that must be unique had some duplicates e.g. Customer ID
   - The governorate was repeated in Arabic, in different styles, and we had lots of redundancy that needs to be addressed, because if ignored no meaningful analysis can be done.
   - Duplicate categories with spelling variations & some were in Arabic

# PART A
# DATA QUALITY ISSUES

## 3. MISSING DATA ISSUES (NULLs)

- Some columns were severely missing and cannot be imputed, we need to get the value from the source again for example if it is a system e.g. names, contact info like  email, phone no.

## 4.CALCULATION/FORMULA ERRORS

- Some column were derived from other through formula and was computed wrong

## 5. DATA TYPE ISSUES

- Numbers were stored as text, and that is a major problem in case we needed to perform calculations on them.

## 6. Lack of Standardization

- even the product SKU and the customer ID which suppose to be generated by a system non-human were also various and lacked having a standard style.

## 7.COMPLETENESS ISSUES

- Missing crucial fields.

**WHY this happened?**

I THINK THE ERRORS ARE BECAUSE OF THE LACK OF CONTROL ON THE ENTRIES AND THE HUMAN MISTAKES, SYSTEM INTEGRATION PROBLEMS, LACK OF DATA GOVERNANCE

## PART A

# HOW TO GET EXTRA DATA?

1.Surveys
- (About their experience, what they want to change, what qualifications does the company offer that keeps the customers loyal)

2.Frequency of their purchases
- (Asses their loyalty)

3. Extract Implicit feedback
- from the customers by collecting data for their behavior(Tracking clicks, engagement level, time spent on a page....) and build a feedback matrix that can be used for building a recommendation system.

4.Seasonal Pattern of purchases

5.Collect data about their profession

6. Wishlist or Saved Items to capture their purchase intent and may be put some discounts on those items

7.Birthday Personal celebrations or for example send 20% off on the customer birthday, keep him loyal

# PART A

# CONTROLS TO ADD ON THE SCREEN

1. Mandatory input (if the user doesn't enter a field error message arises)
2. for the phone se regex for example its mandatory to write country's code then 10 numbers only
3. for the email put regex to control the format,only English letters , the position of the @, the domain all of these are controls
4. If someone enter in the gender "انثى" arise error
5. may be even put as much as we can of menus and choices between F, M (Drop down Lists)
6. Reduce the calculations done by human ,Let the system calculates
7. Ask the user for confirmation before submitting

# PART B

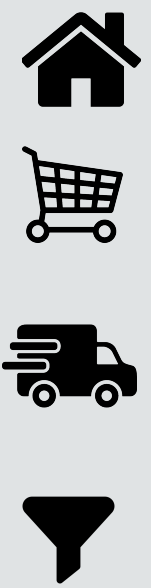# ALL THE INSIGHTS REQUIRED IN THIS PART ARE IN THE DASHBOARD

# SALES
# DASHBOARD

**OVERVIEW**

**SHIPPING**

# Overview

## Total Revenue
9.06M

## Total Orders
136

## Average Order Value
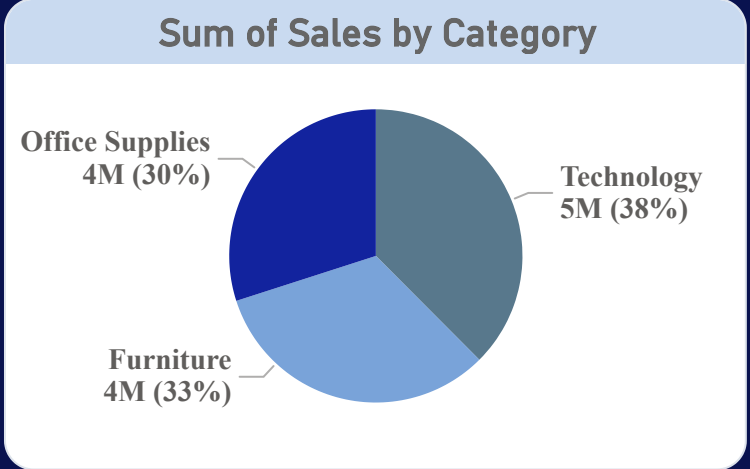66.65K

## Avg Delivery Time
5.88

## Delayed Orders
9.14

## Sales per Gender

4M (49%)

5M (51%)

## Total Revenue by Governorate

IRAQ
JORDAN
ISRAEL
Cairo
KUWAIT
LIBYA
EGYPT
Manama
QATAR
Riyadh
SAUDI ARABIA
U.A.

Microsoft Bing

© 2025 TomTom, © 2025 Microsoft Corporation, © OpenStreetMap

## Sum of Sales by Category

Office Supplies
4M (30%)

Technology
5M (38%)

Furniture
4M (33%)

## Revenue by Months

0.3M

0.2M

0.1M

0.0M

Jan 2024   Feb 2...   Mar 2...   Apr 2...   May 2...   Jun 2024

## Top 10 Products by Sales Revenue

electric kettle - غلاية

men t-shirt - تيشيرت

air fryer 4l - قلاية

data warehousing 101

organic 1 ازيت زيتون

football - كرة قدم

puzzle 1000pcs - أحجية

0M                    1M

## Sales and MaxSales

10M

9.06M

0.00M                    15M

126.43%

### Sidebar
Select all

Qtr 1

Qtr 2

Qtr 3

Qtr 4

# Shipping

## Avg Shipping duration

| Shipper | Duration |
|---------|----------|
| Egypt Post | 9 |
| dhl | 7 |
| aramex | 5 |
| fedex | 2 |

(axis: 0, 5, 10)

## Delivery Performance by Location



© 2025 TomTom, © 2025 Microsoft Corporation, © OpenStreetMap
Microsoft Bing

LIBYA
IRAQ
JORDAN
ISRAEL
Cairo
EGYPT
KUWAIT
Manama
QATAR
Riyadh
SAUDI ARABIA
U.A.

## Delivery Performance

- 12 (9.84%)
- 14 (11.48%)
- 26 (21.31%)
- 26 (21.31%)
- 25 (20.49%)
- 19 (15.57%)

## Orders Distribution Across Categories and Shipping Companies

**ShipperName** ● aramex ● dhl ● Egypt Post ● fedex

| Category | aramex | dhl | Egypt Post | fedex |
|----------|--------|-----|------------|-------|
| electronics | 16 | 6 | 9 | 8 |
| home | 9 | 10 | 7 | 5 |
| fashion | 6 | 3 | 8 | 3 |
| grocery | 3 | 3 | 4 | 4 |
| books | 4 | 4 | 4 | |
| sports | 4 | | 5 | |
| toys | 3 | | | 4 |

(axis: 0, 10, 20, 30, 40)

## Shipping by Category



- Office Supplies 410K (30%)
- Technology 507K (37%)
- Furniture 441K (32%)

### Sidebar

Select all

Qtr 1

Qtr 2

Qtr 3

Qtr 4

# Payment & Risk Analysis

## Unpaid Orders - Revenue at Risk

| CustomerName | Month | Sum of TotalAmount | PaymentStatus | Pro |
|---|---|---|---|---|
| ahmed ali | April | 33,111.54 | unpaid | air |
| ahmed ali | May | 75.00 | unpaid | org |
| ahmed fathi | | 31,041.74 | unpaid | pu |
| ahmed ibrahim | March | 163,504.61 | unpaid | foo |
| ahmed mahmoud | March | 62,163.35 | unpaid | org |
| ahmed nasr | February | 32,584.31 | unpaid | foo |
| **Total** | | **3,903,271.18** | | |

## Underperforming Products

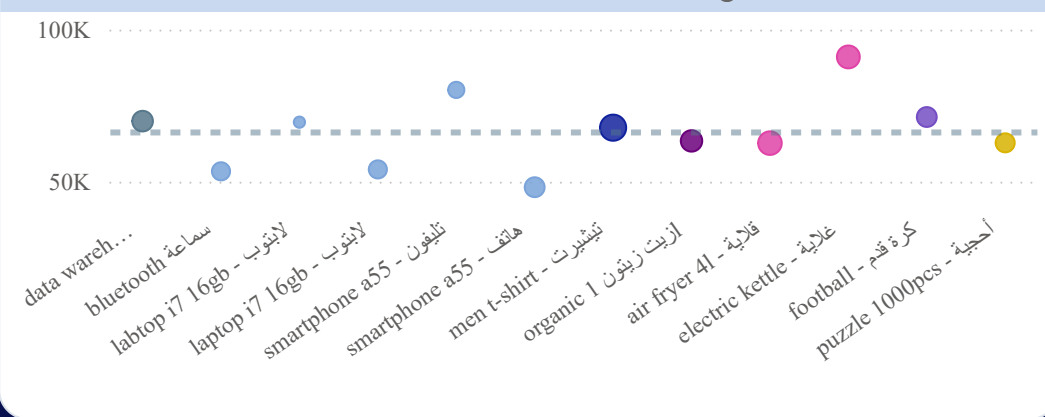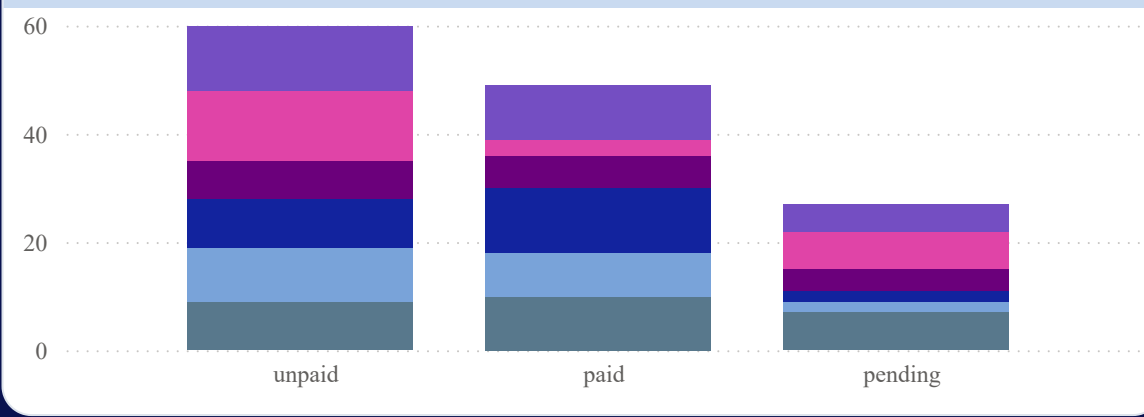| ProductName | Category | Total Orders | TotalRevenue | Average Order Va |
|---|---|---|---|---|
| labtop i7 16gb - لابتوب | electronics | 3 | 208,691.61 | 69,563 |
| bluetooth سماعة | electronics | 9 | 480,671.31 | 53,407 |
| laptop i7 16gb - لابتوب | electronics | 9 | 486,297.22 | 54,033 |
| smartphone a55 - هاتف | electronics | 11 | 530,021.51 | 48,183 |
| smartphone a55 - تليفون | electronics | 7 | 561,258.58 | 80,179 |
| puzzle 1000pcs - أحجية | toys | 10 | 627,715.66 | 62,771 |
| **Total** | | **136** | **9,064,869.06** | **66,653** |

## Unpaid Orders %

**64**

## Revenue at Risk

**5.97M**

## Avg Order Value

**66.65K**

## Products Above/Below Average AOV



## Payment Status Analysis



Select all

Qtr 1

Qtr 2

Qtr 3

Qtr 4

# PART C

## DATA QUALITY

Check and fix data quality issues based on the DAMA dimensions (accuracy, completeness, etc.). Produce three key items: a DQ report, a cleansing guide (step-by-step specs), and a working script or SQL.

- The report is in the next page
- both of the cleaning and script are in the repo
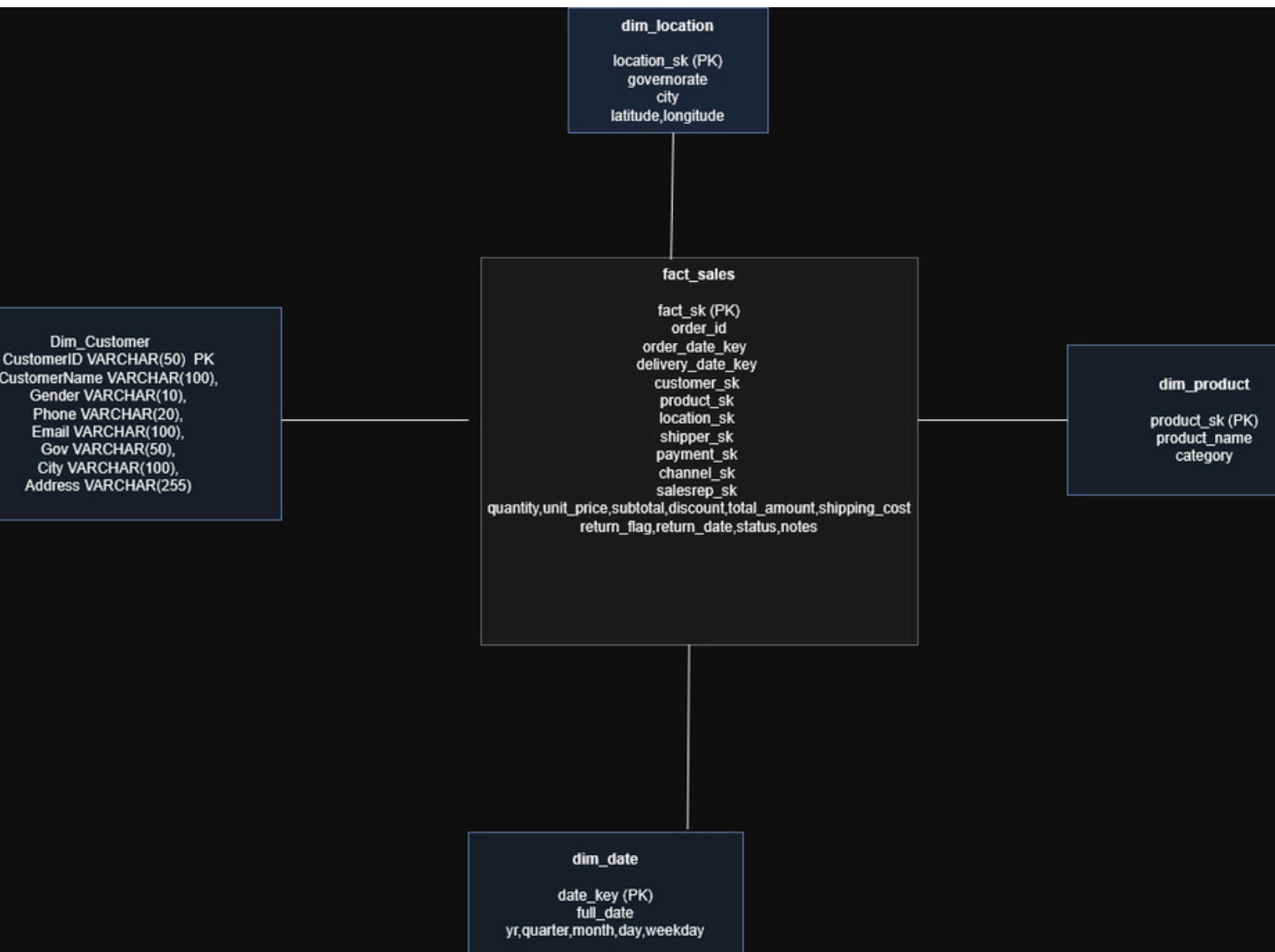
# 1) DATA QUALITY REPORT

| Data Quality Violation | How did it show in the dataset |
|---|---|
| 1.Accuracy | TotalAmount(Wrong calculations)<br>UnitPrice (wrong prices )<br>Quantity<br>Data entries are in mixed languages - Decrease the accuracy of the entries as it losses its value and become not readable |
| 2.Completeness | This dataset had lots of missing values that were crucial and cannot be imputed, how would you impute customerID unless you go back to the system (tried using the customer sheet but it didn't fill all the Nulls in the Sales_Order table regarding the customer ID<br>ProductSKU , Email (Lots of Nulls ), Phone<br>DeliveryDate (Nulls), Category (but could be imputed), Address<br>ProductDescription , ShippingCost |
| 3.Consistency | data was very far away from being consistent in the same column entries vary a lot their is an obvious example in the category column<br><br>variations: "Electronics" vs "Electronic" vs 50%), "الكترونيات"<br><br>Phone (+20, 012, 010"), CustomerName ( capitalization:"), PaymentStatus Nasr City" vs "Nasr", Email , Channel Status |
| 4.Validity | The data isn't even valid<br>Email(most of the email either wrong fromate or filler test.test)  Phone , CustomerID , ProductSKU, DeliveryDate some of them are before OrderDate , Quantity ( negative values), UnitPrice ( zero/negative), Discount (exceeds subtotal) |
| 5.Uniqueness | There are many duplicates even in the Customer ID,some personal info is duplicated, CustomerName ( duplicate variations: "Ahmed Ali" vs "Ahmad Aly"), ProductName ( English/Arabic duplicates), ProductSKU ( duplicate formats), Category, OrderID |
| 6.Timeliness | Data is not even valid for the moment to be timeliness but it doesn't satisfy this constraint and some delivery date is earlier than order date |

# PART D

## 1) THE DATA MODEL AND DESIGN

I UPLOADED IT CLEARER ON THE GITHUB REPO

# PART D

**THE CODES ETL& THE MODELLING IN THE
WAREHOUSE ARE ON GITHUB**

Thank you