## Project Title

**Supermart Profit Prediction & Retail Business Analytics System**

---

# Business Problem

Retail businesses face challenges in predicting profit due to fluctuating sales, discount strategies, seasonal demand, and product-level variations. Business teams need a reliable system to **analyze historical performance** and **predict future profitability** before making pricing or promotional decisions.

---

# Project Objective

The objectives of this project were:

- To analyze historical supermarket sales data
- To understand the relationship between sales, discounts, and profit
- To build a machine learning model for profit prediction
- To deploy an interactive application for real-time business decision support
- To demonstrate the use of **Excel, SQL, Python, and Machine Learning** in a single end-to-end project

---

# Dataset Overview

The dataset contains historical grocery sales transactions with attributes such as:

- Order details (Order ID, Customer Name)
- Product information (Category, Sub-category)
- Location details (City, Region, State)
- Financial metrics (Sales, Discount, Profit)
- Time-related information (Order Date)

This dataset represents a realistic retail environment suitable for business analytics and predictive modeling.

---

## TOOLS & TECHNOLOGIES USED

- **Excel** – Initial data understanding & validation
- **SQL** – Structured querying & business-level analysis
- **Python** – Data cleaning, EDA, feature engineering
- **Machine Learning (Scikit-learn)** – Model development
- **Streamlit** – Model deployment & dashboard creation

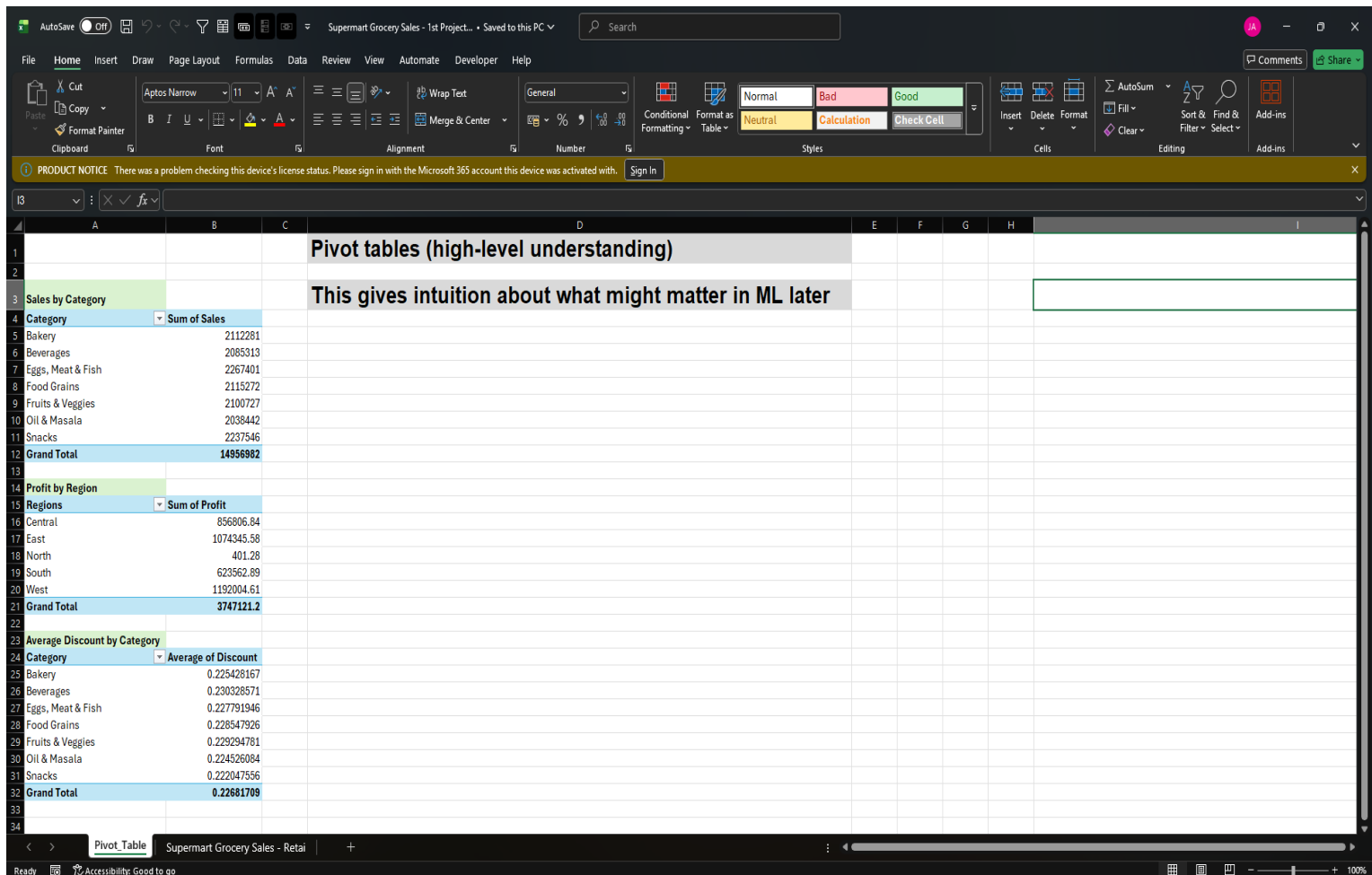# ☐ Step 1: Excel Analysis (WHY & HOW)

## ☐ Why Excel was used

Excel was used as the **first-level analysis tool** to quickly understand the data before moving to advanced tools.

## ☐ What was done in Excel

- Reviewed raw data structure and column meanings
- Checked for missing values and incorrect data types
- Identified outliers in Sales, Discount, and Profit
- Created basic pivot tables:
  - Category-wise sales and profit
  - Region-wise performance
- Used charts to visually understand:
  - Sales vs Profit trends
  - Discount impact on profit

## ☐ Outcome:
Excel helped validate data quality and provided initial business insights that guided further analysis in SQL and Python.

# ☐ Step 2: SQL Analysis (WHY & HOW)

## ☐ Why SQL was used

SQL was used to perform **structured, business-driven analysis** similar to how data is analyzed in real organizations using databases.

## ☐ What was done using SQL

- Queried total sales and profit by:
    - Category
    - Sub-category
    - Region
- Identified top-performing and low-performing products
- Analyzed discount impact using grouped aggregations
- Validated business KPIs such as:
    - Average discount per region
    - Profit contribution by category

## ☐ Example business questions answered via SQL:

- Which category generates the highest profit?
- Which region is most sensitive to discounts?
- How does profit vary across different product segments?

## ☐ Outcome:
SQL helped transform raw data into **business-ready insights** and ensured the analysis aligned with real-world reporting practices.

---

# ☐ Step 3: Data Cleaning & Preprocessing (Python)

Python was used for deeper data preparation:

- Cleaned and standardized column names
- Converted numerical and date fields to proper formats
- Removed irrelevant columns for modeling
- Handled encoding of categorical variables
- Prepared a clean dataset suitable for machine learning

---

# ☐ Step 4: Exploratory Data Analysis (EDA)

EDA was performed using Python libraries:

- Analyzed distributions of sales, profit, and discount
- Studied relationships such as:
    - Discount vs Profit
    - Category vs Profit

- Identified seasonal patterns using year and month

This step confirmed many insights initially observed in Excel and SQL.

---

# ☐ Step 5: Feature Engineering

Additional features were created to enhance model learning:

- Year, Month, Quarter
- Weekend indicator
- One-hot encoding for categorical variables
- Numeric-only feature selection for ML compatibility

---

# ☐ Step 6: Machine Learning Modeling

## Baseline Model: Linear Regression

- Used to establish a baseline performance
- Provided interpretability of linear relationships

## Advanced Model: Random Forest Regressor

- Captured non-linear patterns
- Improved handling of complex feature interactions
- Provided better business-level insights

Models were evaluated using:

- Mean Absolute Error (MAE)
- R² Score

---

# ☐ Step 7: Model Deployment using Streamlit

The trained Random Forest model was deployed using Streamlit Cloud:

- Users can input order details
- Real-time profit predictions are generated
- Business warnings and insights are displayed
- The app converts ML output into actionable insights

This step transformed the analysis into a **production-style business application**.

---

# ☐ Business Impact & Insights

- High discounts negatively impact profitability
- Certain categories consistently perform better
- Profit prediction supports pricing and promotion decisions
- Scenario testing reduces financial risk

---

# ☐ Final Conclusion

This project demonstrates a complete retail analytics solution by combining Excel, SQL, Python, Machine Learning, and Streamlit deployment. Excel and SQL were used for initial data validation and business-level analysis, while Python enabled advanced analytics and modeling. The final deployed application allows stakeholders to predict profit in real time and make informed business decisions. This project reflects real-world data science workflows and highlights the importance of combining technical skills with business understanding.

# SQl Results :

MySQL Workbench — Local instance MySQL80

```sql
1   -- Basic validation
2   -- Why:
3   -- Confirms one row = one order.
4   select * from grocery_sales
5   SELECT COUNT(*) FROM grocery_sales;
6   SELECT COUNT(DISTINCT Order_ID) FROM grocery_sales;
7
8   -- Business aggregations
9   -- Why:
10  -- Confirms which features may influence the target variable.
11  SELECT Category,
12      SUM(Sales) AS total_sales,
13      SUM(Profit) AS total_profit
14  FROM grocery_sales
15  GROUP BY Category;
16
17  SELECT Region,
```

Result:

| COUNT(DISTINCT Order_ID) |
| --- |
| 9994 |

Output:

| # | Time | Action | Message | Duration / Fetch |
| --- | --- | --- | --- | --- |
| 26 | 12:42:54 | select * from grocery_sales LIMIT 0, 5000 | 5000 row(s) returned | 0.000 sec / 0.015 sec |
| 27 | 12:43:54 | SELECT Region,   AVG(Discount) AS avg_discount,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Region LIMIT 0, 5... | 5 row(s) returned | 0.031 sec / 0.000 sec |
| 28 | 12:46:05 | SELECT  CASE   WHEN Discount < 0.15 THEN 'Low'   WHEN Discount < 0.30 THEN 'Medium'   ELSE 'High'  END AS discount_buc... | 3 row(s) returned | 0.016 sec / 0.000 sec |
| 29 | 12:47:42 | select * from grocery_sales LIMIT 0, 5000 | 5000 row(s) returned | 0.000 sec / 0.015 sec |
| 30 | 12:49:22 | SELECT COUNT(*) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.000 sec / 0.000 sec |
| 31 | 12:49:32 | SELECT COUNT(DISTINCT Order_ID) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.015 sec / 0.000 sec |

---

Result Grid:

| Category | total_sales | total_profit |
| --- | --- | --- |
| Oil & Masala | 2038442 | 497895.2899999996 |
| Beverages | 2085313 | 525605.7600000001 |
| Food Grains | 2115272 | 529162.6400000006 |
| Fruits & Veggies | 2100727 | 530400.3800000008 |
| Bakery | 2112281 | 528521.06 |
| Snacks | 2237546 | 568178.8499999999 |
| Eggs, Meat & Fish | 2267401 | 567357.2200000002 |

Output:

| # | Time | Action | Message | Duration / Fetch |
| --- | --- | --- | --- | --- |
| 27 | 12:43:54 | SELECT Region,   AVG(Discount) AS avg_discount,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Region LIMIT 0, 5... | 5 row(s) returned | 0.031 sec / 0.000 sec |
| 28 | 12:46:05 | SELECT  CASE   WHEN Discount < 0.15 THEN 'Low'   WHEN Discount < 0.30 THEN 'Medium'   ELSE 'High'  END AS discount_buc... | 3 row(s) returned | 0.016 sec / 0.000 sec |
| 29 | 12:47:42 | select * from grocery_sales LIMIT 0, 5000 | 5000 row(s) returned | 0.000 sec / 0.015 sec |
| 30 | 12:49:22 | SELECT COUNT(*) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.000 sec / 0.000 sec |
| 31 | 12:49:32 | SELECT COUNT(DISTINCT Order_ID) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.015 sec / 0.000 sec |
| 32 | 12:49:44 | SELECT Category,   SUM(Sales) AS total_sales,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Category LIMIT 0, 5000 | 7 row(s) returned | 0.031 sec / 0.000 sec |

---

```sql
8   -- Business aggregations
9   -- Why:
10  -- Confirms which features may influence the target variable.
11  SELECT Category,
12      SUM(Sales) AS total_sales,
13      SUM(Profit) AS total_profit
14  FROM grocery_sales
15  GROUP BY Category;
16
17  SELECT Region,
18      AVG(Discount) AS avg_discount,
19      SUM(Profit) AS total_profit
20  FROM grocery_sales
21  GROUP BY Region;
22
23  -- Discount vs profit logic
24  -- Why:
```

Result Grid:

| Region | avg_discount | total_profit |
| --- | --- | --- |
| North | 0.12 | 401.28 |
| South | 0.226775797523162255 | 623562.8899999996 |
| West | 0.224726618607556608 | 119200.4.60999999982 |
| Central | 0.228725785220.4054 | 856806.8400000016 |
| East | 0.227672050516179823 | 1074345.5799999982 |

Output:

| # | Time | Action | Message | Duration / Fetch |
| --- | --- | --- | --- | --- |
| 28 | 12:46:05 | SELECT  CASE   WHEN Discount < 0.15 THEN 'Low'   WHEN Discount < 0.30 THEN 'Medium'   ELSE 'High'  END AS discount_buc... | 3 row(s) returned | 0.016 sec / 0.000 sec |
| 29 | 12:47:42 | select * from grocery_sales LIMIT 0, 5000 | 5000 row(s) returned | 0.000 sec / 0.015 sec |
| 30 | 12:49:22 | SELECT COUNT(*) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.000 sec / 0.000 sec |
| 31 | 12:49:32 | SELECT COUNT(DISTINCT Order_ID) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.015 sec / 0.000 sec |
| 32 | 12:49:44 | SELECT Category,   SUM(Sales) AS total_sales,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Category LIMIT 0, 5000 | 7 row(s) returned | 0.031 sec / 0.000 sec |
| 33 | 12:49:56 | SELECT Region,   AVG(Discount) AS avg_discount,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Region LIMIT 0, 5... | 5 row(s) returned | 0.031 sec / 0.000 sec |

---

```sql
19      SUM(Profit) AS total_profit
20  FROM grocery_sales
21  GROUP BY Region;
22
23  -- Discount vs profit logic
24  -- Why:
25  -- Builds business intuition before ML.
26
27  SELECT
28      CASE
29          WHEN Discount < 0.15 THEN 'Low'
30          WHEN Discount < 0.30 THEN 'Medium'
31          ELSE 'High'
32      END AS discount_bucket,
33      AVG(Profit) AS avg_profit
34  FROM grocery_sales
35  GROUP BY discount_bucket;
```

Result Grid:

| discount_bucket | avg_profit |
| --- | --- |
| Low | 374.9035185185 1847 |
| Medium | 376.5306747939 5514 |
| High | 370.9870137103 683 |

Output:

| # | Time | Action | Message | Duration / Fetch |
| --- | --- | --- | --- | --- |
| 29 | 12:47:42 | select * from grocery_sales LIMIT 0, 5000 | 5000 row(s) returned | 0.000 sec / 0.015 sec |
| 30 | 12:49:22 | SELECT COUNT(*) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.000 sec / 0.000 sec |
| 31 | 12:49:32 | SELECT COUNT(DISTINCT Order_ID) FROM grocery_sales LIMIT 0, 5000 | 1 row(s) returned | 0.015 sec / 0.000 sec |
| 32 | 12:49:44 | SELECT Category,   SUM(Sales) AS total_sales,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Category LIMIT 0, 5000 | 7 row(s) returned | 0.031 sec / 0.000 sec |
| 33 | 12:49:56 | SELECT Region,   AVG(Discount) AS avg_discount,   SUM(Profit) AS total_profit FROM grocery_sales GROUP BY Region LIMIT 0, 5... | 5 row(s) returned | 0.031 sec / 0.000 sec |
| 34 | 12:50:09 | SELECT  CASE   WHEN Discount < 0.15 THEN 'Low'   WHEN Discount < 0.30 THEN 'Medium'   ELSE 'High'  END AS discount_buc... | 3 row(s) returned | 0.015 sec / 0.000 sec |