

Sales Analysis of Retail Stores

Evgenii Kozin

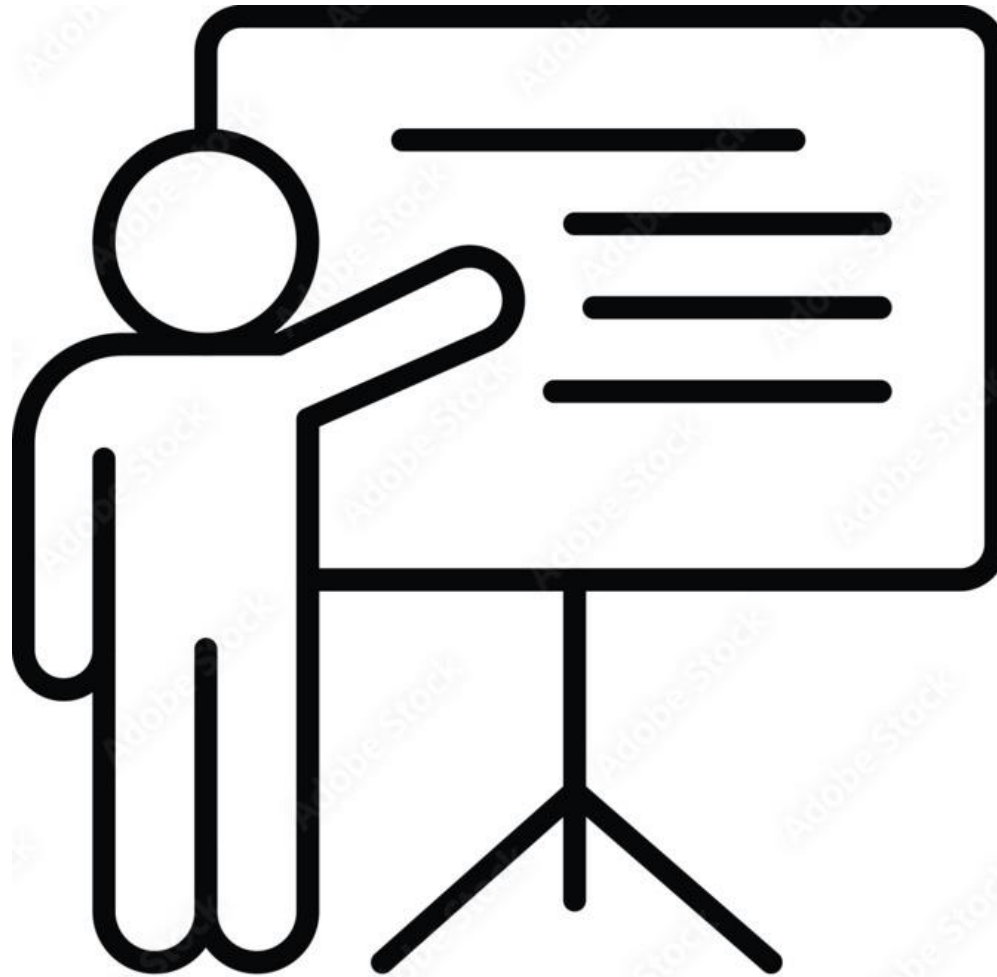


Table of Contents

- Introduction
- Methodology
- Results
- Discussion
- Conclusion
- Appendix



Introduction



Introduction

- The aim of this project is to conduct a comprehensive analysis of the retail chain's sales performance across different regions and periods. The focus is on identifying the main factors influencing sales, including store locations, product categories, and marketing activities. The insights gained will guide strategic decisions and help drive business growth.
- To stay competitive and enhance results, the company is dedicated to leveraging data-driven insights to refine its sales strategies and increase operational efficiency.





Methodology

Data Details



When I downloaded the data first I went back & forth to understand it's format,structure,size of the dataset before jumping into the task



The dataset consist of the following files

sales.csv, product_hierarchy.csv, store_cities.csv, store_names.csv, city_names.csv, and product_names.csv.



Where sales.csv is the target file to analyse and derive insights.sales.csv consist of the date of purchase, productid information,storeid information , types of promo-code, revenue, stock & price information



product_hierarchy.csv and product_names.csv consists of product category, product names, product id which is a unique identity for each product.



store_names.csv & store_cities.csv consists of storeid info and storetypeid info to cityid to establish the relationship between variables.



city_names.csv consists of cityid info and city names

Module 1: Data Cleaning and Preparation Using Excel

The initial and most important step in obtaining meaningful insights is to remove unnecessary or irrelevant data from the dataset. For this task, I used Excel, which is an excellent tool for data analysis.

Step 1: Detecting missing values is critical for any data-driven business, as unaddressed gaps can lead to incorrect conclusions. Using the filter feature, I located columns containing null values and removed those records from the dataset.

Step 2: Applied Excel's *Remove Duplicates* function to identify and eliminate duplicate entries, ensuring only unique records remained for further analysis. This step safeguards data integrity.

Step 3: Standardized the data formatting throughout the worksheet to keep it consistent and presentable, such as aligning date formats and setting numerical values to two decimal places.

Step 4: Removed extra spaces within data fields, converted text to uppercase, and inserted spaces where needed to maintain uniformity. This was achieved using the TRIM() function in Excel.

These steps were applied across all datasets to ensure clean, reliable data for analysis.

Module 2: Data Querying and Analysis Using PostgreSQL

- Moving on to PostgreSQL, several steps were followed to load the data into the server:
- **Step 1:** Created a database and tables, then populated them using the BI-dump.sql file, which contains datasets such as sales.csv and product.csv.
- **Step 2:** In the command-line interface, a new database named *BI-Capstone* was created and connected by running the command `\connect capstone`. The tables and data from BI-dump.sql were restored using the command `\include BI-dump.sql`.
- **Step 3:** Verified that the data loaded correctly using *pgAdmin*, the graphical interface for PostgreSQL. For each table, ran the query `SELECT * FROM table_name LIMIT 5;` to confirm the presence and accuracy of the data.
- **Step 4:** Performed sales analysis using features such as *Cubes* and *Rollups* to calculate total sales, as well as breakdowns by product, city, and store.

Module 3: Data Visualization and Statistical Analysis

- **Regression Analysis: A Step-by-Step Guide**
- The core purpose of this statistical analysis is to interpret regression results to understand how sales are influenced by various factors such as **date**, **revenue**, **stock**, and **price**. This process provides a detailed report on sales performance across these key variables.
- **The Process**
- **Define Variables:** Begin by identifying **sales** as the dependent variable. Analyze how other variables correlate with sales to find any strong relationships.
- **Run the Analysis:** Use a regression analysis tool, such as the one in Excel, to generate a summary of the output that will explain the relationships between the variables.
- **Use a Sample Dataset:** To save time and computational resources, analyze a **sample** of about 20,000 rows from the dataset instead of the entire dataset.
- **Interpret Coefficients:** Examine the coefficients to determine the expected change in sales for a one-unit increase in each independent variable.
- **Assess Significance:** Evaluate the **p-values** (significance levels) to confirm if the influence of the independent variables is statistically significant.
- **Evaluate Model Fit:** Analyze the **R-squared** value to understand the proportion of the variance in sales that can be explained by the independent variables.

Module 4: Data Visualization and Dashboards Using Tableau

The next task is to develop an **interactive dashboard** that effectively visualizes the entire story through charts. This process involves a structured approach to data handling, visualization, and final presentation.

Step-by-Step Process

- **Data Integration and Preparation:** The first step is to **import data** into a business intelligence tool like **Tableau**. Once the data is loaded, you must establish the correct relationships between the variables to ensure the output is accurate. You can also create **calculated fields** at this stage to simplify complex calculations and aggregations for later use in visualizations.
- **Visualization Development:** Next, you will create individual sheets to visualize the data. This involves building various **charts**, such as line charts, bar charts, and bubble charts, to represent the data in a clear and compelling manner. Each chart should be designed to tell a specific part of the story.
- **Dashboard Assembly:** After creating the individual sheets and charts, they are combined into a single **dashboard**. The goal is to arrange these sheets in a logical flow to present a cohesive and insightful narrative for the audience.
- **Enhancing Interactivity:** The final step is to refine the dashboard's user experience. This is achieved by utilizing Tableau's features like **filters** and **highlighters**. These features allow stakeholders to interact with the data directly; for example, clicking on an item in one chart will **highlight** or **filter** the corresponding data points in all other related charts, making the analysis intuitive and easy to understand.



Results

Module 1

Lesson 1: Data Cleaning and Preparation

Data cleaning and preparation are fundamental to any successful data project, whether it's for analysis, machine learning, or generating business insights. These processes are crucial because they ensure the data used is **accurate**, **consistent**, and in a **usable** format. Properly prepared data serves as a reliable foundation for all subsequent steps. This includes building accurate predictive models and making informed, data-driven decisions. Without this critical phase, the quality and reliability of any conclusions or outputs would be compromised, potentially leading to flawed results.

product_id	store_id	date	sales	revenue	stock	price	promo_type_1	promo_bin_1	promo_type_2	promo_bin_2	promo_discount_2	promo_discount_type_2	promo_discount_type_3	store_names
P0055	S0001	01/02/2017	0	0	14	3.5	PR05	verylow	PR03					Electro World (National Chain)
P0068	S0001	01/02/2017	0	0	6	5.1	PR10	verylow	PR03					Electro World (National Chain)
P0080	S0001	01/02/2017	0	0	10	7.9	PR06	veryhigh	PR03					Electro World (National Chain)
P0083	S0001	01/02/2017	0	0	2	6.9	PR05	moderate	PR03					Electro World (National Chain)
P0144	S0001	01/02/2017	0	0	4	4.95	PR10	low	PR03					Electro World (National Chain)
P0175	S0001	01/02/2017	0	0	5	5.5	PR05	high	PR03					Electro World (National Chain)
P0212	S0001	01/02/2017	0	0	8	9.95	PR05	low	PR03					Electro World (National Chain)
P0218	S0001	01/02/2017	0	0	6	59.9	PR05	moderate	PR03					Electro World (National Chain)
P0229	S0001	01/02/2017	0	0	2	6.85	PR05	verylow	PR03					Electro World (National Chain)
P0237	S0001	01/02/2017	0	0	33	59.9	PR05	high	PR03					Electro World (National Chain)
P0241	S0001	01/02/2017	0	0	4	12.95	PR05	verylow	PR03					Electro World (National Chain)
P0267	S0001	01/02/2017	0	0	17	4	PR10	high	PR03					Electro World (National Chain)
P0277	S0001	01/02/2017	0	0	2	3.45	PR10	verylow	PR03					Electro World (National Chain)
P0282	S0001	01/02/2017	0	0	9	1.9	PR06	high	PR03					Electro World (National Chain)
P0287	S0001	01/02/2017	1	5.56	8	7.5	PR05	verylow	PR03					Electro World (National Chain)
P0312	S0001	01/02/2017	0	0	16	12.9	PR05	high	PR03					Electro World (National Chain)
P0327	S0001	01/02/2017	1	1.81	12	1.95	PR07	verylow	PR03					Electro World (National Chain)
P0332	S0001	01/02/2017	0	0	9	26.9	PR07	verylow	PR03					Electro World (National Chain)
P0345	S0001	01/02/2017	0	0	13	3.15	PR06	verylow	PR03					Electro World (National Chain)
P0348	S0001	01/02/2017	2	3.11	12	2.1	PR05	verylow	PR03					Electro World (National Chain)
P0389	S0001	01/02/2017	0	0	4	8.85	PR05	low	PR03					Electro World (National Chain)
P0399	S0001	01/02/2017	0	0	12	3.9	PR05	moderate	PR03					Electro World (National Chain)
P0413	S0001	01/02/2017	2	3.12	4	1.99	PR03	verylow	PR03					Electro World (National Chain)
P0427	S0001	01/02/2017	0	0	16	2.45	PR06	low	PR03					Electro World (National Chain)
P0446	S0001	01/02/2017	0	0	5	5.95	PR05	verylow	PR03					Electro World (National Chain)
P0458	S0001	01/02/2017	0	0	122	4.25	PR05	moderate	PR03					Electro World (National Chain)
P0500	S0001	01/02/2017	0	0	94	5.95	PR10	verylow	PR03					Electro World (National Chain)
P0506	S0001	01/02/2017	0	0	11	5.95	PR05	verylow	PR03					Electro World (National Chain)
P0605	S0001	01/02/2017	0	0	1	59.9	PR10	high	PR03					Electro World (National Chain)
P0680	S0001	01/02/2017	0	0	19	28.9	PR10	low	PR03					Electro World (National Chain)
P0694	S0001	01/02/2017	0	0	8	5.99	PR10	high	PR03					Electro World (National Chain)
P0705	S0001	01/02/2017	0	0	14	7.95	PR12	verylow	PR03					Electro World (National Chain)
P0055	S0001	01/03/2017	0	0	14	3.5	PR05	verylow	PR03					Electro World (National Chain)
P0068	S0001	01/03/2017	0	0	6	5.1	PR10	verylow	PR03					Electro World (National Chain)
P0080	S0001	01/03/2017	0	0	10	7.9	PR06	veryhigh	PR03					Electro World (National Chain)
P0083	S0001	01/03/2017	0	0	8	6.9	PR05	moderate	PR03					Electro World (National Chain)
P0144	S0001	01/03/2017	0	0	4	4.95	PR10	low	PR03					Electro World (National Chain)
P0175	S0001	01/03/2017	0	0	5	5.5	PR05	high	PR03					Electro World (National Chain)
P0212	S0001	01/03/2017	0	0	8	9.95	PR05	low	PR03					Electro World (National Chain)
P0229	S0001	01/03/2017	0	0	2	6.85	PR05	verylow	PR03					Electro World (National Chain)

Module 1, Lesson 2: Data Analysis Using Pivot Tables

219	SolarOvenwave	117	352.14	2526
220	SolarTVdrive	7	70.98	449
221	Soundbar with Dolby Atmos	1	5.3	760
222	Stackable Washer and Dryer with Steam Refresh	0	0	276
223	Steam Mops	7	93.46	744
224	Steam Washer with Sanitize Cycle	20	106.45	558
225	String Trimmers	2	14.7	172
226	TechBlendercast	17	46.42	745
227	TechDryerdrive	13	209.14	1800
228	TechDryerlux	25	79.82	1767
229	TechFridgematic	184	511.26	5785
230	TechFridgepulse	60	124.68	4602
231	TechGrillcast	99	222.78	2186
232	TechHeaterflow	2	16.38	1359
233	TechHeaterhub	2,014	40.83	0
234	TechMixerflow	44	459.99	1333
235	TechMixergen	18	180	183
236	TechOvengen	49	130.67	7473
237	TechTVgen	99	168.3	1943
238	TechTVhub	160	146.86	3879
239	TechTVlux	341	315.92	5473
240	TechVacuumflow	93	698.85	930
241	Toaster Oven	66	118.09	5236
242	Upright Vacuums	5	52.97	825
243	Upright Vacuums with Cyclonic Technology and Pet Hair Remover	6	63.61	818
244	Vacuum Sealers with Roll Cutter and Marinate Mode	47	315.48	762
245	Wine Cooler	2	57.46	849
246	Grand Total	13147.60	47367.13	466304.33
247	AVERAGE SALES	54.33		

1				
2				
3	City	Sum of revenue	Sum of sales	Sum of stock
4	Edinburgh	18146.5	5537.989	133535.99
5	Helsinki	17880.98	3674.995	153521.285
6	London	5080.29	1900.615	45375.05
7	Saint Petersburg	4545.47	1247	71470
8	Vienna	1713.89	787	62402
9	Grand Total	47367.13	13147.599	466304.325
10	AVERAGE SALES		2629.5198	

1				
2				
3	Store_name	Sum of sales	Sum of revenue	Sum of stock
4	Currys (National Chain)	389	1055.21	31473
5	Darty	1900.615	5080.29	45375.05
6	DIGI	5537.989	18146.5	133535.99
7	Electro World (National Chain)	3674.995	17880.98	153521.285
8	Elettrodomestici Rossi	1247	4545.47	71470
9	Euronics Lisboa (National Chain)	398	658.68	30929
10	Grand Total	13147.599	47367.13	466304.325
11	AVERAGE_SALES	2191.27		

Module 2

Lesson 1:

Data Querying Using PostgreSQL

The screenshot displays the DBeaver PostgreSQL interface. The left sidebar shows the database structure with 'bicapstone' selected. The main query editor contains the SQL query: `SELECT * FROM product_names LIMIT 5;`. The right sidebar shows the 'Database sessions' tab with a table of session statistics. The bottom pane shows the 'Messages' tab with a table of query results.

Database sessions

Database sessions	Total	Active	Idle	Transactions per second

Messages

store_id	date	sales	revenue	stock
S0002	2017-02-01	0	0	8
S0012	2017-02-01	1	5.3	0
S0013	2017-02-01	2	10.59	0
S0023	2017-02-01	0	0	6
S0025	2017-02-01	0	0	1

Module 2, Lesson 2: Data Analysis Using PostgreSQL

bicapstone/postgres@PostgreSQL

Query Query History

```
1 SELECT
2   p.hierarchy1_id,
3   p.hierarchy2_id,
4   SUM(sl.sales) AS total_sales
5 FROM
6   sales sl
7 JOIN
8   product_hierarchy p
9 ON
10  sl.product_id = p.product_id
11 GROUP BY rollup
12   (p.hierarchy1_id,
13    p.hierarchy2_id);
14
```

Data Output Messages Notifications

	hierarchy1_id character varying	hierarchy2_id character varying	total_sales double precision
1	[null]	[null]	9856.880999999998
2	H03	H0317	2
3	H00	H0001	617
4	H03	H0316	0
5	H03	H0314	190
6	H01	H0107	374

```
3   TO_CHAR(sl.date, 'YYYY-MM') AS fmt_date,
4   SUM(sl.sales) AS total_sales
5 FROM
6   sales sl
7 GROUP BY rollup
8   (sl.store_id,
9    fmt_date)
10
```

Data Output Messages Notifications

	store_id character varying	fmt_date text	total_sales double precision
1	[null]	[null]	9856.880999999998
2	S0052	2017-02	23.675
3	S0120	2017-02	19
4	S0135	2017-02	32.305
5	S0024	2017-02	85
6	S0127	2017-03	9
7	S0107	2017-03	27
8	S0072	2017-02	80
9	S0102	2017-03	23
10	S0082	2017-03	42
11	S0015	2017-02	24

Tools Help

Dependencies X Dependents X Processes X bicapstone/postgres@PostgreSQL* x public.o

bicapstone/postgres@PostgreSQL

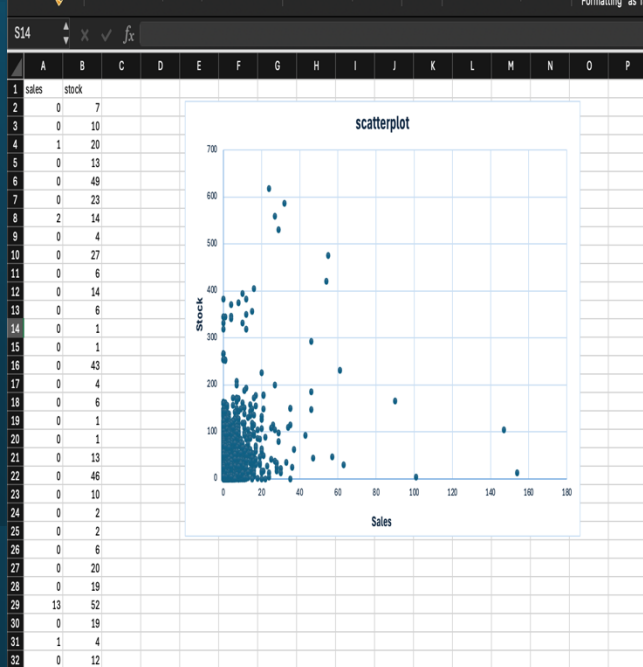
Query Query History

```
1 SELECT
2   c.city_id,
3   cn.city_name,
4   TO_CHAR(s.date, 'YYYY-MM') AS month_sale,
5   SUM(s.sales) AS total_sales_sum
6 FROM
7   sales s
8
9 JOIN
10  store_cities c
11 ON
12  s.store_id = c.store_id
13 JOIN
14  city_names cn
15 ON
16  c.city_id=cn.city_id
17 GROUP BY
18   CUBE (c.city_id, cn.city_name, TO_CHAR(s.date, 'YYYY-MM'));
19
```

Data Output Messages Notifications

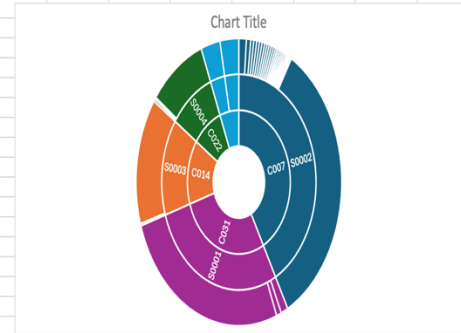
	city_id character varying	city_name character varying (255)	month_sale text	total_sales_sum double precision
1	C002	Berlin	2017-02	100
2	C002	Berlin	2017-03	59
3	C002	Berlin	[null]	159
4	C002	[null]	[null]	159
5	C003	Barcelona	2017-02	32.84
6	C003	Barcelona	2017-03	14
7	C003	Barcelona	[null]	46.84
8	C003	[null]	[null]	46.84
9	C004	Budapest	2017-02	223.49

Module 3, Lesson 1: Data Visualization Using Excel

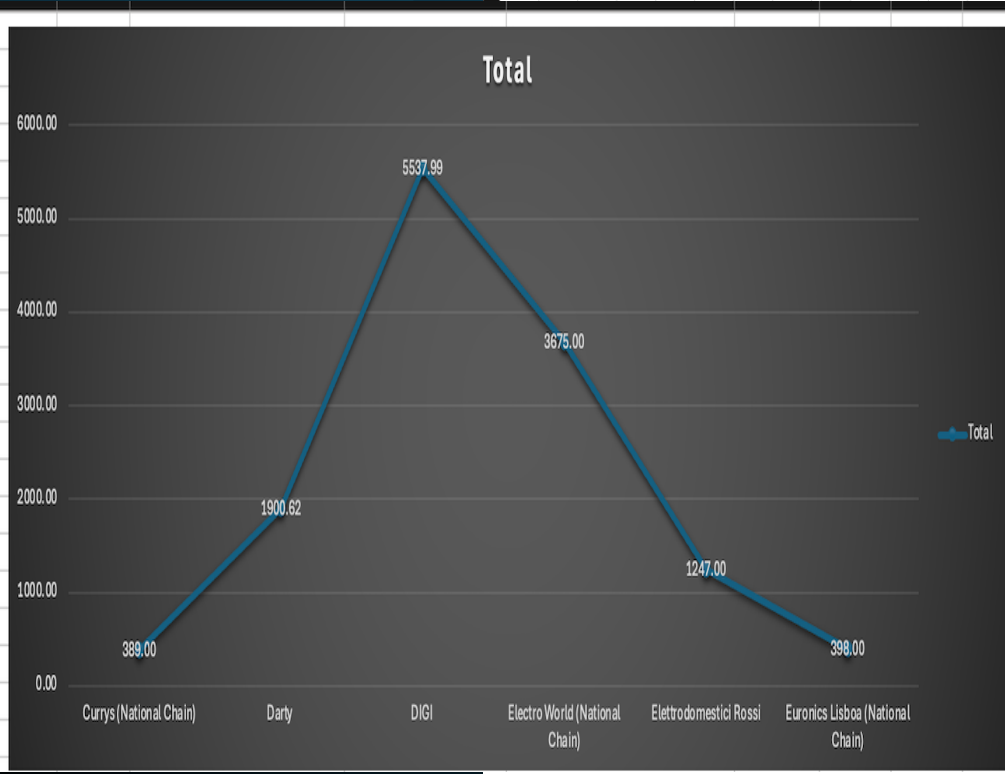


Sheet 1 - store_cities.c Sheet1-store_names. Sheet 1 - product_h sales

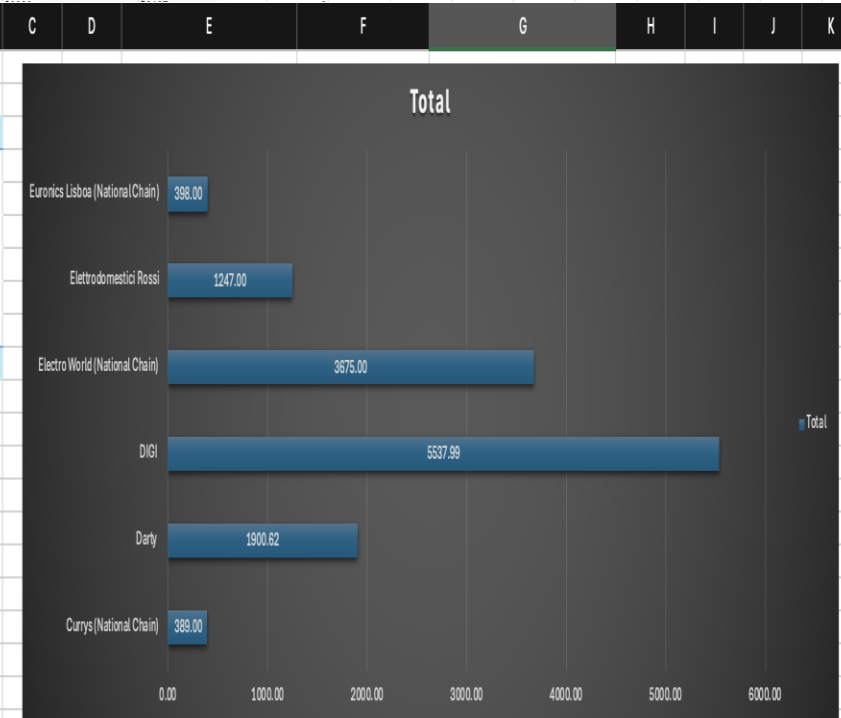
2	C007	S0002	P0001	0
3	C007	S0002	P0005	0
4	C007	S0002	P0011	0
5	C007	S0002	P0015	1
6	C007	S0002	P0017	0
7	C007	S0002	P0018	1
8	C007	S0002	P0024	0
9	C007	S0002	P0035	2
10	C007	S0002	P0046	0
11	C007	S0002	P0051	7
12	C007	S0002	P0054	0
13	C007	S0002	P0055	0
14	C007	S0002	P0057	0
15	C007	S0002	P0060	0
16	C007	S0002	P0067	0
17	C007	S0002	P0070	1
18	C007	S0002	P0079	0
19	C007	S0002	P0083	0
20	C007	S0002	P0090	0
21	C007	S0002	P0092	2
22	C007	S0002	P0099	0
23	C007	S0002	P0102	0
24	C007	S0002	P0103	5
25	C007	S0002	P0109	0
26	C007	S0002	P0110	1
27	C007	S0002	P0114	0
28	C007	S0002	P0116	0
29	C007	S0002	P0125	0
30	C007	S0002	P0129	0
31	C007	S0002	P0131	0
32	C007	S0002	P0134	1



1		
2		
3	RowLabels	Sum of sales
4	Currys (National Chain)	389.00
5	Darty	1900.62
6	DIGI	5537.99
7	Electro World (National Chain)	3675.00
8	Elettrodomestici Rossi	1247.00
9	Euronics Lisboa (National Chain)	398.00
10	Grand Total	13147.60
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		



1		
2		
3	RowLabels	Sum of sales
4	Currys (National Chain)	389.00
5	Darty	1900.62
6	DIGI	5537.99
7	Electro World (National Chain)	3675.00
8	Elettrodomestici Rossi	1247.00
9	Euronics Lisboa (National Chain)	398.00
10	Grand Total	13147.60
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		
22		



Module 3, Lesson 2: Statistical Analysis

SUMMARY OUTPUT

Regression Statistics

Multiple R	0.73269124
R Square	0.53683645
Adjusted R S	0.53674379
Standard Err	1.90960891
Observations	19999

ANOVA

	df	SS	MS	F	Significance F				
Regression	4	84507.6798	21126.92	5793.5842	0				
Residual	19994	72910.2439	3.64660618						
Total	19998	157417.924							

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	-0.5933672	44.846394	-0.0132311	0.98944355	-88.496006	87.3092711	-88.496006	87.3092711
date	1.88E-05	0.00104885	0.0179283	0.98569623	-0.002037	0.00207463	-0.002037	0.00207463
revenue	0.04606812	0.00032757	140.636675	0	0.04542606	0.04671018	0.04542606	0.04671018
stock	0.02416673	0.00049781	48.5463819	0	0.02319098	0.02514247	0.02319098	0.02514247
price	-0.0148461	0.00107736	-13.780087	5.28E-43	-0.0169579	-0.0127344	-0.0169579	-0.0127344

• R-Square (0.5368):

- About 53.68% of the variability in the dependent variable (likely "sales") is explained by the independent variables (date, revenue, stock, price). This suggests a decent model fit but indicates that other factors not included in the model might explain the remaining 46.32%.

○ Interpretation of Predictors:

1. Intercept (-0.5934):

- Represents the predicted value of the dependent variable when all independent variables are zero. This has no practical interpretation here due to the high p-value (0.9894), indicating it is not statistically significant.

2. Date (Coefficient = 1.8804E-05, p-value = 0.9857):

- This predictor is not statistically significant (p-value > 0.05). It likely has no meaningful impact on the dependent variable in this model.

3. Revenue (Coefficient = 0.0461, p-value = 0):

- This is a highly significant predictor (p-value < 0.05).
- For every unit increase in revenue, the dependent variable (e.g., sales) increases by **0.0461** units, holding all other variables constant.

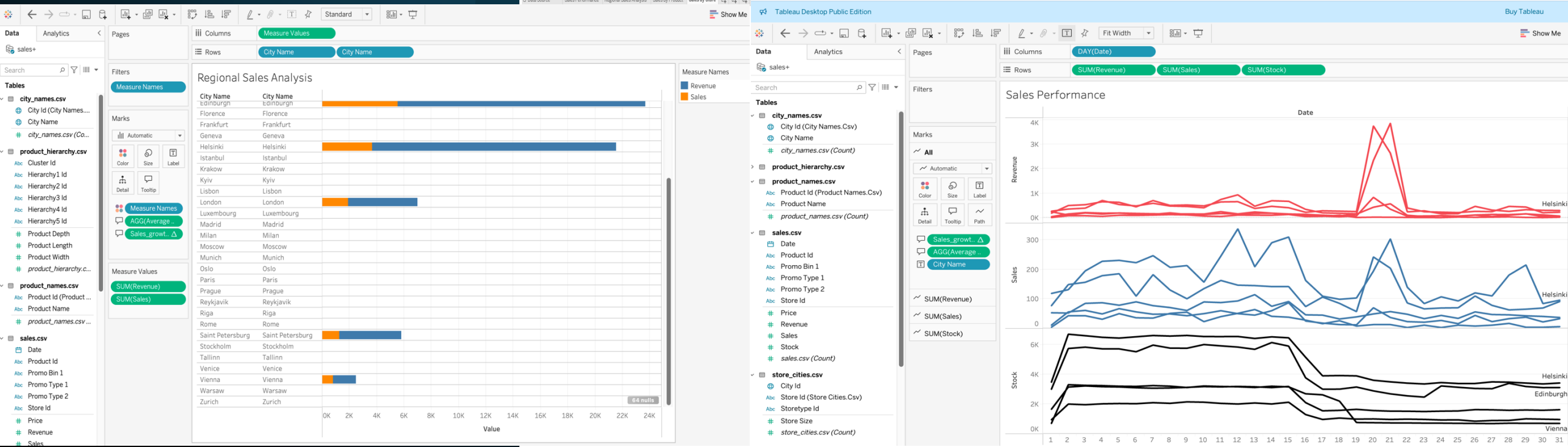
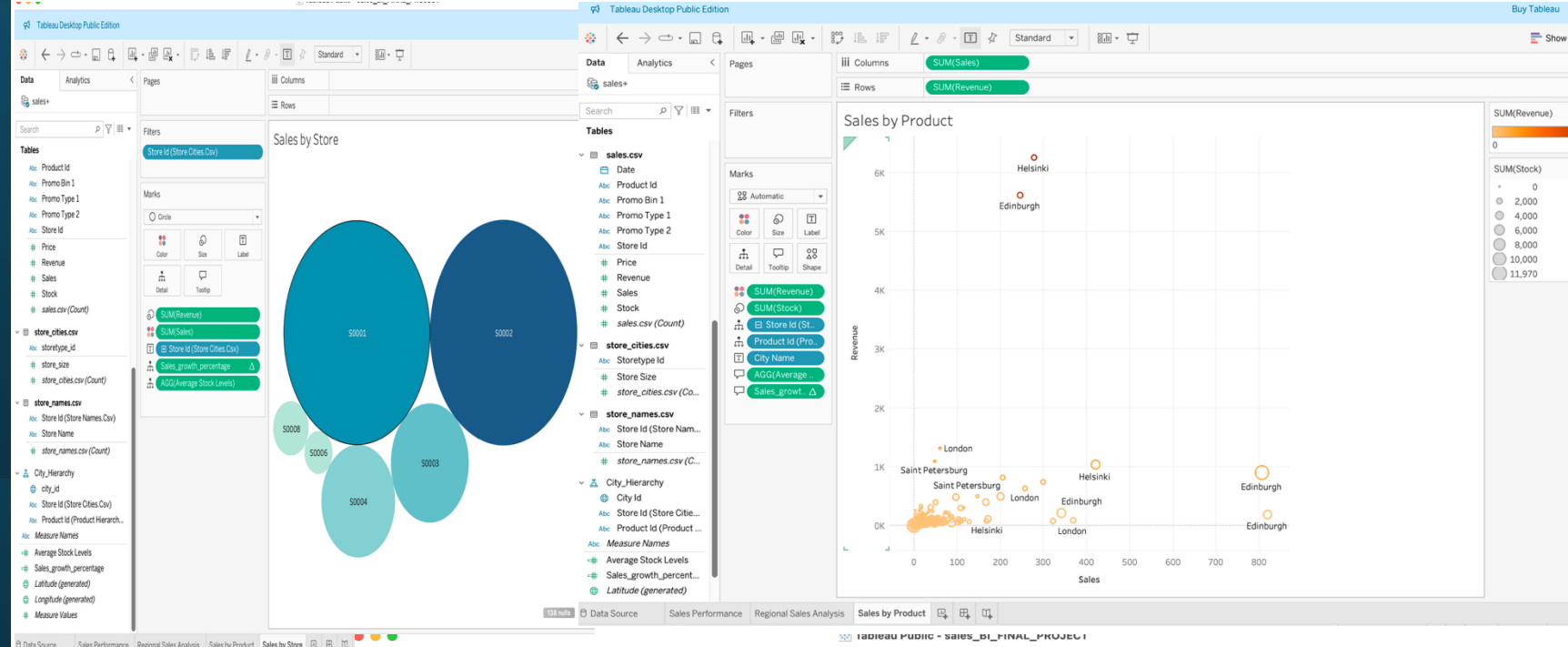
4. Stock (Coefficient = 0.0242, p-value = 0):

- Another highly significant predictor.
- For every unit increase in stock, the dependent variable increases by **0.0242** units, holding other variables constant.

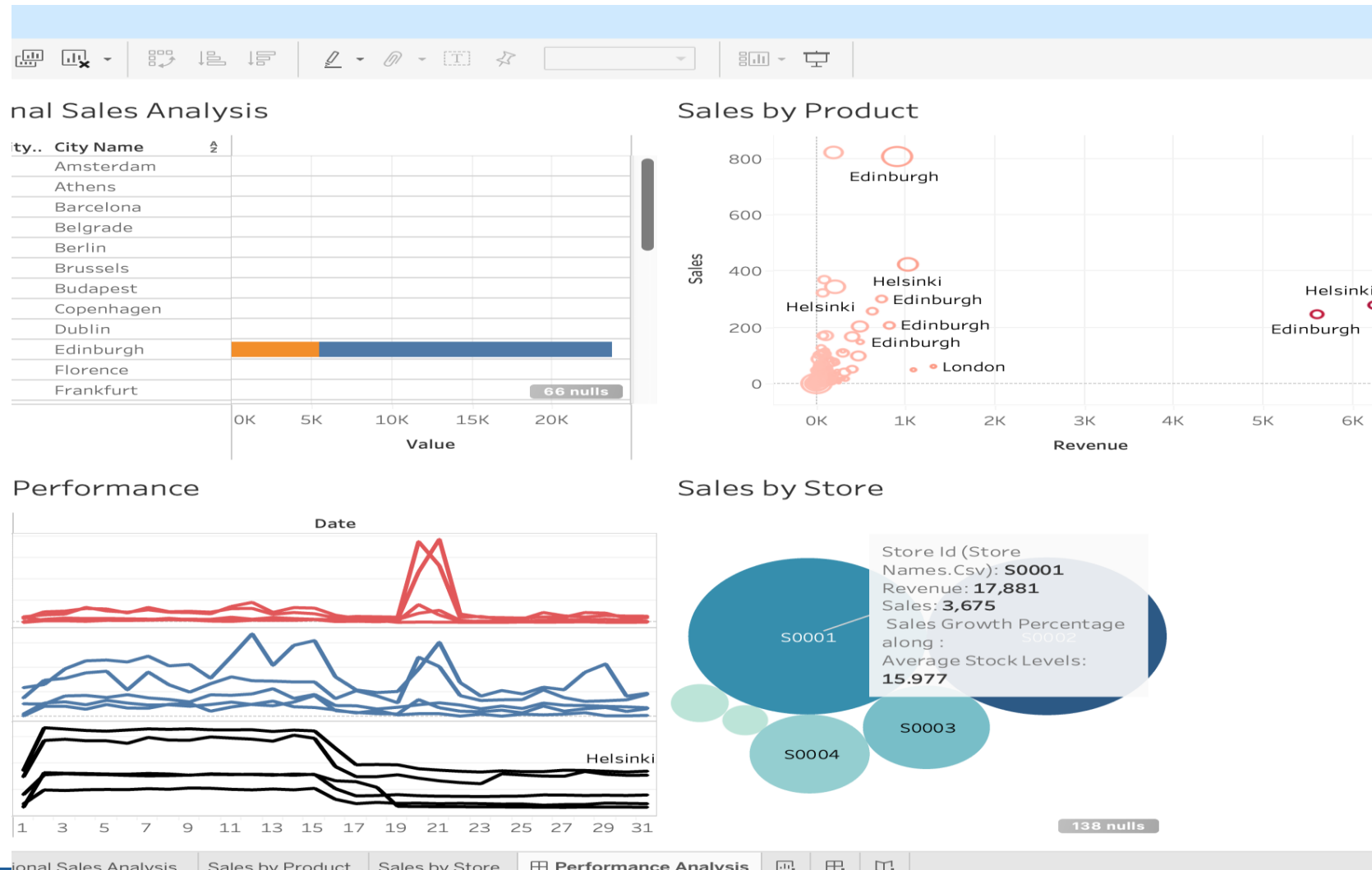
5. Price (Coefficient = -0.0148, p-value = 5.28E-43):

- This is also highly significant.
- For every unit increase in price, the dependent variable decreases by **0.0148** units, holding other variables constant. This suggests that as prices increase, sales decrease, which aligns with typical economic behavior.

Module 4, Lesson 1: Basic Tableau Visualizations



Module 4, Lesson 2: Advanced Visualizations Using Tableau





Discussion

Insights and Recommendations

Revenue is the Most Important Factor:

The analysis shows that revenue has a strong and direct impact on sales. As revenue increases, sales also increase significantly.

Implication: The retail chain should focus on strategies to boost revenue, such as promotions, targeted marketing, and expanding product offerings.

Having more stock available contributes positively to sales. When stock levels increase, sales also go up.

Implication: Ensuring that popular products are always in stock can help maintain and boost sales.

Higher prices tend to reduce sales. Even small price increases can lead to a noticeable drop in sales.

Implication: The retail chain should consider competitive pricing strategies and discount offers to attract more customers.



Conclusion

Summary

Analysis of Sales Performance Across Cities

Based on the provided key findings, here is a rewritten summary and conclusion, emphasizing clarity and actionable insights.

Key Findings

The analysis reveals distinct performance patterns across different cities:

- **Helsinki and Edinburgh** show strong revenue despite moderate sales volumes. This suggests that these locations are either selling **high-value products** or have adopted successful **premium pricing strategies**.
- In contrast, **London and Saint Petersburg** generate lower revenue relative to their high sales volume. This indicates a potential reliance on **lower-priced products** or smaller transaction values.
- **Stock management** appears to be a key area for improvement. Some cities maintain high stock levels without a corresponding increase in revenue, which points to potential inefficiencies in inventory and logistics.

Conclusion and Recommendations

The company's strategy should focus on two primary areas for optimization:

- **Pricing Strategy:** Review and adjust the pricing model in cities like **London and Saint Petersburg** to increase revenue without compromising sales volume. The success of Helsinki and Edinburgh can serve as a benchmark for this effort.
- **Stock Management:** Improve inventory allocation by linking **stock levels to revenue potential**. This will help prevent overstocking in low-revenue cities and ensure that high-revenue locations have adequate supply to meet demand.