

# ***GROUP B***

Contributing Team Members:

- Adhish Pillai
- Aswin Muthusamy
- Ganga Hariharan
- Omkar Krishnapurkar
- Yashwanth Duddupuddi

# **Contents**

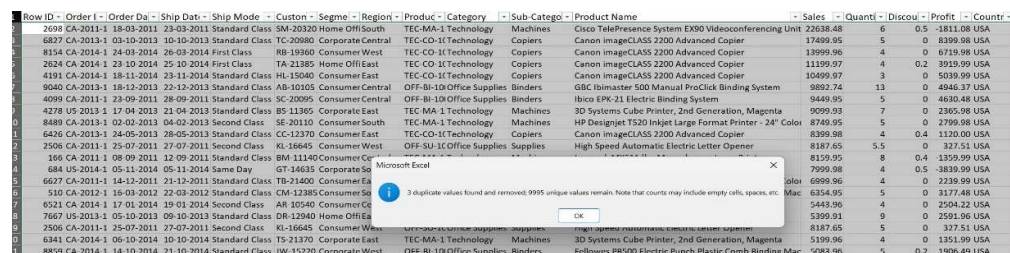
1. *Introduction*
2. *Description of Analysis Done*
3. *Visualizations and Insights*
4. *Team Collaboration*
5. *Conclusion*

## 1. Introduction

9,998 records from a comprehensive US sales dataset will be examined to identify patterns, fix errors, and evaluate the performance of client segments and regions. This will help data-driven initiatives aimed at boosting total sales and accomplishing corporate objectives.

## 2. Description of Analysis done

We took a few crucial actions during the data preparation process to guarantee the dataset's integrity and suitability for analysis. To preserve uniqueness and correctness, we first eliminated three duplicate entries.



Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer Segment	Region	Product Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Country
2698	CA-2011-1	18-03-2011	23-03-2011	Standard Class	SM-20320	Home Office/South	TEC-MA-1 Technology	Machines	Cisco TelePresence System EX90 Videoconferencing Unit	22638.48	6	0.5	-1811.08	USA
6827	CA-2013-1	03-10-2013	10-10-2013	Standard Class	TC-20980	Corporate Central	TEC-CO-1 Technology	Copiers	Canon imageCLASS 2200 Advanced Copier	17409.95	5	0	8399.98	USA
8154	CA-2014-1	24-03-2014	26-03-2014	First Class	RB-19360	Consumer West	TEC-CO-1 Technology	Copiers	Canon imageCLASS 2200 Advanced Copier	13999.96	4	0	6719.98	USA
2628	CA-2014-1	23-10-2014	25-10-2014	First Class	TA-21380	Home Office/East	TEC-CO-1 Technology	Copiers	Canon imageCLASS 2200 Advanced Copier	11199.97	4	0.2	3919.99	USA
4191	CA-2014-1	18-11-2014	23-11-2014	Standard Class	HL-15040	Consumer East	TEC-CO-1 Technology	Copiers	Canon imageCLASS 2200 Advanced Copier	10499.97	3	0	5039.99	USA
9040	CA-2013-1	18-12-2013	22-12-2013	Standard Class	AB-10105	Consumer Central	OFF-BI-1 Office Supplies	Binders	GBC biMaster 500 Manual ProClick Binding System	9892.74	13	0	4946.37	USA
4099	CA-2011-1	28-09-2011	28-09-2011	Standard Class	SC-20995	Consumer Central	OFF-BI-1 Office Supplies	Binders	Ibico EPK-21 Electric Binding System	9449.95	5	0	4630.48	USA
4278	US-2013-1	17-04-2013	21-04-2013	Standard Class	BS-11365	Corporate East	TEC-MA-1 Technology	Machines	3D Systems Cube Printer, 2nd Generation, Magenta	9099.93	7	0	2365.98	USA
8489	CA-2013-1	02-02-2013	04-02-2013	Second Class	SE-20110	Consumer South	TEC-MA-1 Technology	Machines	HP DesignJet T520 Inkjet Large Format Printer - 24" Color	8749.95	5	0	2799.98	USA
6426	CA-2013-1	24-05-2013	28-05-2013	Standard Class	CC-12370	Consumer East	TEC-CO-1 Technology	Copiers	Canon imageCLASS 2200 Advanced Copier	8399.98	4	0.4	1120.00	USA
2508	CA-2011-1	25-07-2011	27-07-2011	Second Class	KL-16645	Consumer West	OFF-BI-1 Office Supplies	Supplies	High Speed Automatic Electric Letter Opener	8187.65	5.5	0	327.51	USA
166	US-2011-1	08-09-2011	12-09-2011	Standard Class	BM-11140	Consumer Central	Microsoft Excel			8159.95	8	0.4	1359.99	USA
684	US-2014-1	05-11-2014	05-11-2014	Same Day	GT-14635	Corporate South				7999.98	4	0.5	3839.99	USA
6627	CA-2011-1	14-12-2011	21-12-2011	Standard Class	TB-21400	Consumer East				6999.96	4	0	2239.99	USA
510	CA-2012-1	16-03-2012	22-03-2012	Standard Class	CM-12385	Consumer South				6354.95	5	0	3177.48	USA
6521	CA-2014-1	17-01-2014	19-01-2014	Second Class	AR-10540	Consumer Central				5449.96	4	0	2504.22	USA
7667	US-2013-1	05-10-2013	09-10-2013	Standard Class	DR-12940	Home Office/East				5399.91	9	0	2591.96	USA
2506	CA-2011-1	25-07-2011	27-07-2011	Second Class	KL-16645	Consumer West				8187.65	5	0	327.51	USA
6341	CA-2014-1	06-10-2014	10-10-2014	Standard Class	TS-21370	Corporate East	TEC-MA-1 Technology	Machines	3D Systems Cube Printer, 2nd Generation, Magenta	9399.96	4	0	1351.99	USA
8859	CA-2014-1	14-10-2014	21-10-2014	Standard Class	JW-10220	Corporate West	OFF-BI-1 Office Supplies	Binders	Fellowes PB500 Electric Punch Plastic Comb Binding Mac	5083.96	5	0.2	1806.49	USA

Figure a: Removing Duplicates in Excel

We found a discrepancy where one of the observations has "Technologi" instead of "Technology," for category.

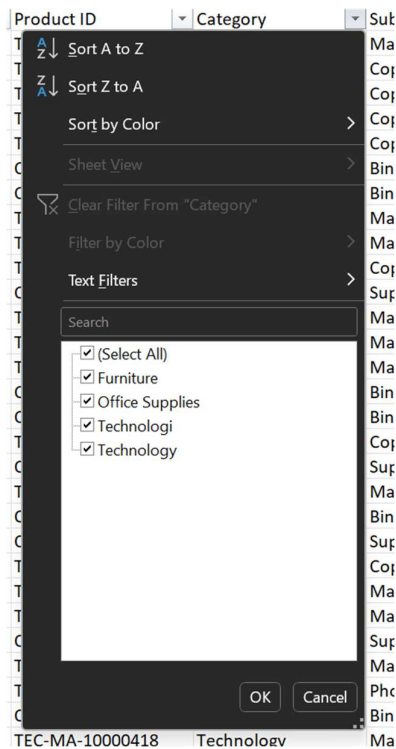


Figure b.1: Filter showing discrepancy

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category
CA-2011-145541	14-12-2011	21-12-2011	Standard Class	TB-21400	Consumer East		TEC-MA-10001127	Technology	

Figure b.2: Record with 'Technologi'

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category
CA-2011-145541	14-12-2011	21-12-2011	Standard Class	TB-21400	Consumer East		TEC-MA-10001127	Technology	

Figure b.3: Record changed from 'Technologi' to 'Technology'

Since there was not enough data and reference values to implement a formula to all product IDs, we have decided to manually handle profits with blank data.

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Price Per Product	Country
US-2012-126977	17-09-2012	23-09-2012	Standard Class	PF-19120	Consumer East		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	4228.7	6	0.2		880.98	USA
CA-2011-116246	12-09-2011	17-09-2011	Standard Class	LW-17215	Consumer East		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	3785.29	6	0.1		700.98	USA

Figure c: Product IDs with blank values.

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Price Per Product	Country
CA-2014-118892	18-08-2014	23-08-2014	Standard Class	TP-21415	Consumer East		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	4416.17	9	0.3	-630.88	700.98	USA
CA-2011-116246	12-09-2011	17-09-2011	Standard Class	LW-17215	Consumer East		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	3785.29	6	0.1		700.98	USA
CA-2013-122903	28-05-2013	30-05-2013	Second Class	LA-16780	Corporate Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	3504.9	5	0	700.98	700.98	USA
CA-2014-102204	02-05-2014	07-05-2014	Standard Class	CI-12010	Consumer South		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2803.92	5	0.2	0	700.98	USA
CA-2012-164777	27-01-2012	29-01-2012	Second Class	SC-20305	Consumer West		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2803.92	5	0.2	0	700.98	USA
CA-2012-139731	15-10-2012	15-10-2012	Same Day	JE-15745	Consumer Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2453.43	5	0.3	-350.49	700.98	USA
CA-2013-136406	16-04-2013	18-04-2013	Second Class	BD-11320	Consumer West		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	1121.57	2	0.2	0	700.98	USA
CA-2011-143903	20-07-2011	24-07-2011	Standard Class	KM-16375	Home Office Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	981.37	2	0.3	-140.2	700.98	USA

Figure d: For Product ID: FUR-CH-10002024

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Price Per Product	Country
CA-2014-118892	18-08-2014	23-08-2014	Standard Class	TP-21415	Consumer East		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	4416.17	9	0.3	-630.88	700.98	USA
CA-2011-116246	12-09-2011	17-09-2011	Standard Class	LW-17215	Consumer East		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	3785.29	6	0.1	420.59	700.98	USA
CA-2013-122903	28-05-2013	30-05-2013	Second Class	LA-16780	Corporate Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	3504.9	5	0	700.98	700.98	USA
CA-2014-102204	02-05-2014	07-05-2014	Standard Class	CI-12010	Consumer South		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2803.92	5	0.2	0	700.98	USA
CA-2012-164777	27-01-2012	29-01-2012	Second Class	SC-20305	Consumer West		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2803.92	5	0.2	0	700.98	USA
CA-2012-139731	15-10-2012	15-10-2012	Same Day	JE-15745	Consumer Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	2453.43	5	0.3	-350.49	700.98	USA
CA-2013-136406	16-04-2013	18-04-2013	Second Class	BD-11320	Consumer West		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	1121.57	2	0.2	0	700.98	USA
CA-2011-143903	20-07-2011	24-07-2011	Standard Class	KM-16375	Home Office Central		FUR-CH-10002024	Furniture	Chairs	HON 5400 Series Task Chairs	981.37	2	0.3	-140.2	700.98	USA

Figure e: For Product ID: FUR-CH-10002024 Profit Calculated

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Price Per Product	Country
CA-2012-117086	08-11-2012	12-11-2012	Standard Class	QJ-19255	Corporate East		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	4404.9	5	0	1013.13	880.98	USA
US-2012-126977	17-09-2012	23-09-2012	Standard Class	PF-19120	Consumer East		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	4228.7	6	0.2	158.58	880.98	USA
US-2012-150630	17-09-2012	21-09-2012	Standard Class	TB-21520	Consumer East		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	3083.43	7	0.5	-1665.05	880.98	USA
CA-2013-108987	09-09-2013	11-09-2013	Second Class	AG-10675	Consumer Central		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	2396.27	4	0.32	-317.15	880.98	USA
US-2014-109316	09-06-2014	11-06-2014	Second Class	MG-17680	Home Office West		FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	1497.67	2	0.15	140.96	880.98	USA

Figure f: Figure e: For Product ID: FUR-BO-10004834

Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Segment	Region	Product ID	Category	Sub-Category	Product Name	Sales	Quantity	Discount	Profit	Price Per Product	Country
CA-2012-117086	08-11-2012	12-11-2012	Standard Class	QJ-19255	Corporate East	FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	4404.9	5	0	1013.13	880.98	USA	
US-2012-126977	17-09-2012	23-09-2012	Standard Class	PF-19120	Consumer East	FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	4228.7	6	0.2		880.98	USA	
US-2012-150630	17-09-2012	21-09-2012	Standard Class	TB-21520	Consumer East	FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	3083.43	7	0.5	-1665.05	880.98	USA	
CA-2013-108987	09-09-2013	11-09-2013	Second Class	AG-10675	Consumer Central	FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	2396.27	4	0.32	-317.15	880.98	USA	
US-2014-109316	09-06-2014	11-06-2014	Second Class	MG-17680	Home Office West	FUR-BO-10004834	Furniture	Bookcases	Riverside Palais Royal Lawye	1497.67	2	0.15	140.96	880.98	USA	

Figure g: Figure e: For Product ID: FUR-BO-10004834 Profit Calculated

Then we used **Tableau Prep Builder** to efficiently complete the cleaning process.

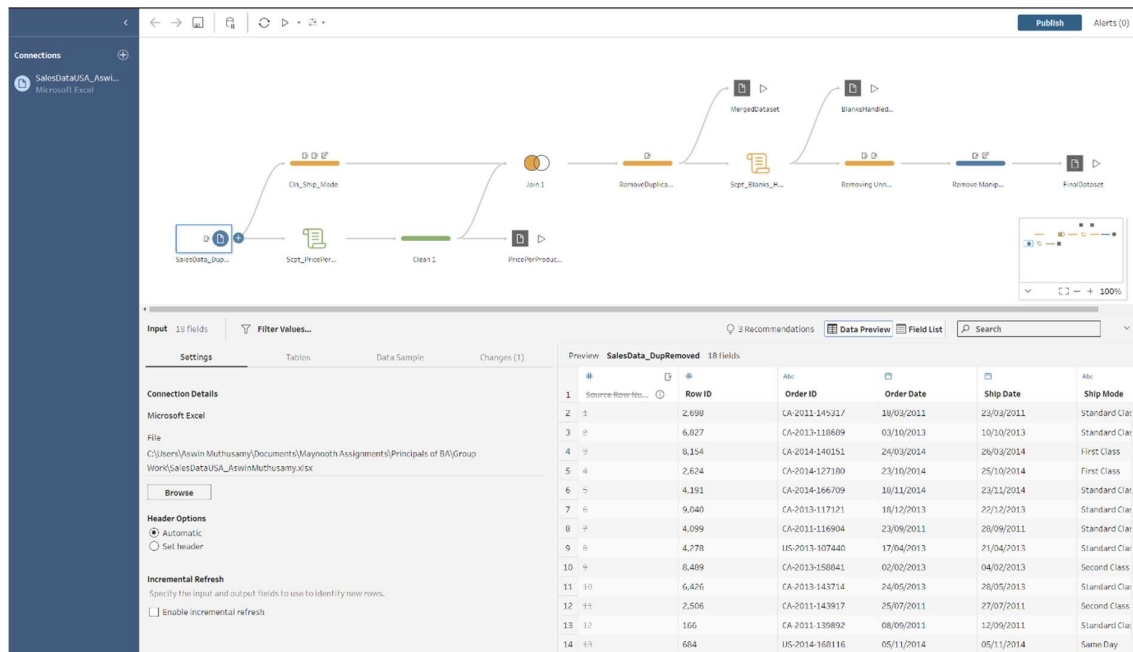


Figure h: Complete Tableau Prep Flow

We developed a Python script that produced a distinct dataset to determine the pricing per product for every unique product ID. This dataset was then combined with the original using a Left Outer Join based on "Product ID." We were able to create a new "Price per Product" column as a result. While in original dataset we are adding a dateDifference calculated field to correct the ship mode based on the below,

dateDifference	Ship Mode
= 0	Same Day
= 1	First Class
= 2	Second Class
>= 3	Standard Class

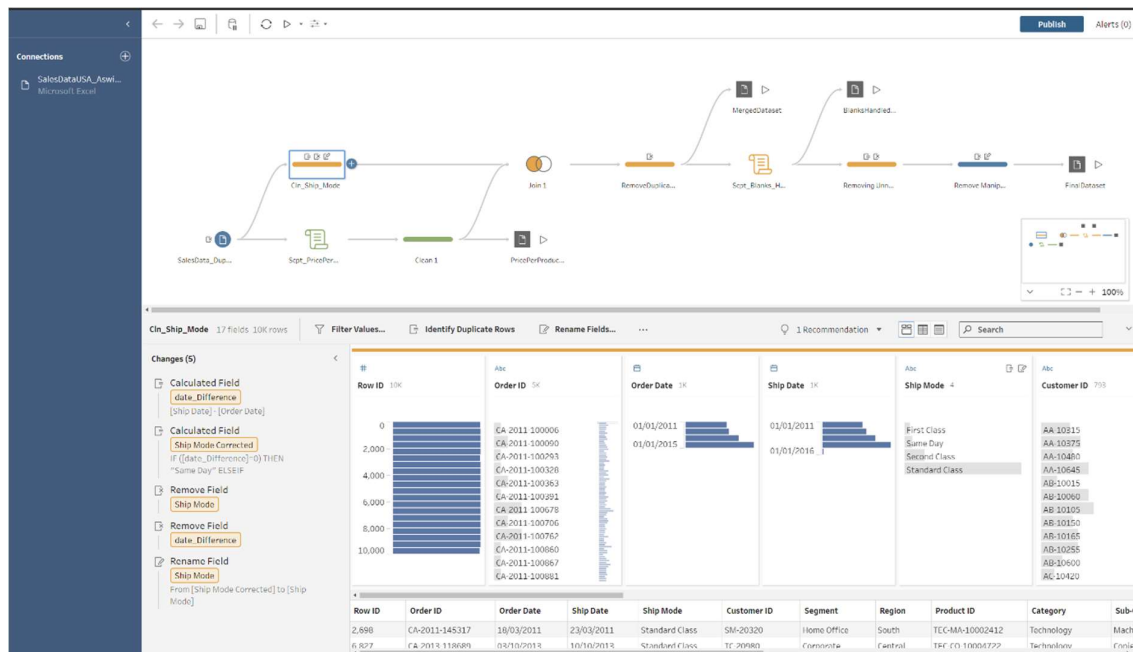


Figure i: Cleaning Ship Mode

```

generatePricePerProduct.py 5 • handlingBlanks.py 9+
Python_Scripts_For_DataCleaning > generatePricePerProduct.py > get_output_schema
1 import pandas as pd
2
3 def generatePrice(SalesData):
4     df_SalesData = pd.DataFrame(SalesData)
5     df_nan_removed = df_SalesData[(df_SalesData['Quantity'] * 10) % 10 == 0]
6     df_nansale_removed = df_nan_removed[(df_nan_removed['Sales'] >= 0)]
7     df_price_per_product_code = df_nansale_removed.drop_duplicates(subset=['Product ID'], keep='first')
8     df_price_per_product_code['Price Per Product'] = ((df_price_per_product_code['Sales'] / df_price_per_product_code['Quantity']) * 100)
9     / (100 - (df_price_per_product_code['Discount'] * 100))
10
11
12
13
14
15     df_price_per_product_code.drop(['Row ID', 'Order ID', 'Order Date',
16     'Ship Date', 'Ship Mode', 'Customer ID', 'Segment', 'Region',
17     'Category', 'Sub-Category', 'Product Name', 'Sales', 'Quantity',
18     'Discount', 'Profit', 'Country'], axis=1, inplace=True)
19     return(df_price_per_product_code)
20
21 # df_price_per_product_code.drop(['Row ID', 'Order ID', 'Order Date',
22 # 'Ship Date', 'Ship Mode', 'Customer ID', 'Segment', 'Region',
23 # 'Category', 'Sub-Category', 'Product Name', 'Profit', 'Country'], axis=1, inplace=True)
24 # return(df_price_per_product_code)
25
26
27 def get_output_schema():
28     return pd.DataFrame({
29         'Product ID': prep_string(),
30         'Price Per Product': prep_decimal()
31     })

```

Figure j: Script that Generates the new 'Price Per Product' dataset

To eliminate blanks and manage decimal values in the "Quantity" and "Sales" columns, we additionally included another Python script. This script will handle blank and decimal values for "Quantity" and the blank values for 'Sales'. We

also removed 'Product IDs' with no 'Sales' or reference values to Calculate 'Price Per Product'.

```
generatePricePerProduct.py 5 • handlingBlanks.py 9+ •
Python_Scripts_For_DataCleaning > handlingBlanks.py > handleBlanks > df_blank_sales_and_qty_handled
1 import pandas as pd
2
3 def handleBlanks(MergedSalesData):
4     #Handling Blanks in Sales
5     df_merged_salesData = pd.DataFrame(MergedSalesData)
6     df_merged_salesData.loc[df_merged_salesData['Sales'].isna(), 'Sales'] =
7     ((df_merged_salesData['Price Per Product'] * df_merged_salesData['Quantity']) * (100 - (df_merged_salesData['Discount'] * 100))) / 100
8     df_blankSales_handled = df_merged_salesData
9
10    #Handling Decimal Values in Quantity
11    df_blankSales_handled.loc[(df_blankSales_handled['Quantity'] * 10) % 10 != 0, 'Quantity'] = 0
12
13    #Handling Blanks and Decimal Values in Quantity
14    df_blankSales_handled.loc[df_blankSales_handled['Quantity'].isna() | df_blankSales_handled['Quantity'] == 0, 'Quantity'] =
15    (df_blankSales_handled['Sales'] * 100) / (df_blankSales_handled['Price Per Product'] * (100 - (df_blankSales_handled['Discount'] * 100)))
16
17    df_blank_sales_and_qty_handled = df_blankSales_handled
18
19    #Eliminating Data with No Sales value and no reference value to calculate Price per product
20    df_final = df_blank_sales_and_qty_handled.dropna(subset='Price Per Product')
21
22    return(df_final)
23
24
25 def get_output_schema():
26     return pd.DataFrame({
27         'Row ID' : prep_int(),
28         'Order ID' : prep_string(),
29         'Order Date' : prep_date(),
30         'Ship Date' : prep_date(),
31         'Ship Mode' : prep_string(),
32         'Customer ID' : prep_string(),
33         'Segment' : prep_string(),
34         'Region' : prep_string(),
35         'Product ID' : prep_string(),
36         'Category' : prep_string(),
37         'Sub-Category' : prep_string(),
38         'Product Name' : prep_string(),
39         'Sales' : prep_decimal(),
40         'Quantity' : prep_decimal(),
41         'Discount' : prep_decimal(),
42         'Profit' : prep_decimal(),
43         'Price Per Product' : prep_decimal(),
44         'Country' : prep_string()
45     });
```

Figure k: Script that handles blanks in 'Sales' and 'Blanks'

Finally, we constructed calculated fields to round values of "Sales," "Quantity," "Profit," and "Price per Product" columns to two decimal places. The Excel worksheet was then renamed to "Refined Dataset" which will now allow us to extract insightful information from the data.

### **3. Visualizations and Insights:**

*In the visualization section, we created a Power Bi Dashboard. Below are the insights and type of visualizations used-*

- **Sales per Month over the years:** *From 2011 to 2014, the "Sales per Month over the Years" figure shows a steady increase trend. Notable sales peaks usually occur in late-year months like October and December, reflecting seasonal demand trends. Sales in the early months of the year, especially January, are frequently lower, indicating both seasonality in consumer behavior and overall company expansion during the studied time.*
- **Monthly Sales Trend:** *With sales peaking in March, September, and October and reaching a maximum of \$0.31 million in October, the "Monthly Sales Trend" chart clearly shows seasonality. January, February, and May are the early and mid-year months with lower sales volumes. These recurrent peaks point to regular demand cycles, which can guide year-round inventory and marketing plans.*
- **Total Sales by Region:** *Sales in the Central, East, South, and West regions are shown in a stacked bar chart. This graph illustrates the regions that contribute the most to sales, demonstrating that some regions—such as the West—have the largest sales. Sales are strongest in the West, which may be the company's primary market for future expansion.*
- **Average Shipping Time for each month:** *Shipping times, which range from 3.7 to 4.2 days, are comparatively stable, according to the "Average Shipping Time for Each Month" figure. Because of effective*



logistics, January and February have the greatest averages (4.2 days), while November has the lowest (3.81 days). April and September peaks could point to locations for delivery optimization due to strong demand or operational difficulties.

- **Total no. of orders based on shipping time:** According to the "Total Number of Orders Based on Shipping Time" figure, the majority of orders are delivered within two to four days, with Standard Class and 3-day delivery being the most common. This implies that consumers value dependable, reasonably priced delivery over speed. Promoting First Class and Same Day delivery could draw clients looking for quicker solutions, while highlighting Standard and Second-Class options may increase customer satisfaction.
- **Correlation between Avg. Shipping speed and total sales:** Faster shipping is positively correlated with better sales, according to the "Correlation Between Average Shipping Speed and Total Sales" graphic. Reduced shipping durations are associated with higher sales, indicating that effective delivery can boost client satisfaction and spur revenue expansion, underscoring the need of optimizing delivery procedures.
- **Profit & Sales by Segment and Region:** With technology accounting for 48.96% of sales, the office category is exceptionally profitable. Office supplies and furniture come in second and third. By directing optimal discount tactics and focused investments in high-performing regions and segments, minimal discounts optimize revenue.

- **Sales by Category:** *The percentage of sales for each product category—furniture, office supplies, and technology—is displayed in a pie chart. The graphic shows that technology products account for the largest portion of total sales. Since technology items account for the largest portion of sales, the company may benefit from concentrating on growing this market.*

- **Sales by Discount:** *The dashboard's "Sales by Category" shows technology leading at 48.96% of revenue, followed by furniture (37.95%) and office supplies (13.09%), guiding strategic focus on technology marketing, inventory, and growth opportunities.*

- **Team Collaboration**

*Team members worked closely together throughout the project, each contributing in a different way.*

<b><i>Sr. No</i></b>	<b><i>Team Member</i></b>	<b><i>Lead Role</i></b>	<b><i>Supporting Role</i></b>
<b><i>1.</i></b>	<b><i>Aswin Muthusamy</i></b>	<b><i>Data Cleaning, Scripting, Creating Dashboard</i></b>	<b><i>Proof-reading</i></b>
<b><i>2.</i></b>	<b><i>Ganga Hariharan</i></b>	<b><i>Creating the report, Drawing Insights from Refined Data</i></b>	<b><i>Proof-reading</i></b>
<b><i>3.</i></b>	<b><i>Omkar Krishnapurkar</i></b>	<b><i>Creating the Report, Scheduling Meetings</i></b>	<b><i>Creating Dashboards</i></b>
<b><i>4.</i></b>	<b><i>Adhish Pillai</i></b>	<b><i>Preparing Presentation</i></b>	<b><i>Drawing Insights from Refined Data</i></b>
<b><i>5.</i></b>	<b><i>Yashwanth Duddupuddi</i></b>	<b><i>Preparing Presentation</i></b>	

*Each member's committed efforts made the project successful, demonstrating an overall and cooperative effort throughout the process.*

#### **4. Conclusion**

*To facilitate strategic planning and data-driven decision-making, this study sought to extract useful insights from sales data. After carefully cleansing the data, we used Power BI to develop an interactive dashboard that allowed us to examine important trends like client profitability, geography sales, and product category success.*

*Important conclusions emphasized the relationship between discount tactics and profit margins as well as how client demographics and geographic location affect sales growth. These findings highlight the necessity of specialized advertising, focused marketing, and improved pricing.*

*All things considered, this study offers a basis for improving sales and profitability, allowing the company to hone its pricing, distribution, and segmentation initiatives to increase productivity and financial results.*