# E-commerce Product Range Analysis

*By: Dorothy Kunth*

# Project Overview

## Project Goal

The goal of this project is to identify the top-performing products, product categories and customers.

## Project Scope

The project's objective is to carry out the following tasks and communicate data findings and recommendations:

- Data Preprocessing
- Exploratory Analysis
- Customer Segmentation using RFM
- Product Segmentation using RFM
- Product Categorization
- Product Category Analysis
- Statistical Hypotheses

## Dataset Overview

The dataset is an open online retail dataset from https://archive.ics.uci.edu/ml/datasets/online+retail#  which is a sales transactions history of an online store that sells household goods. It contains 541,909 transaction records from  2018-11-29 to 2019-12-07. The data has 7 attributes:

**InvoiceNo**: *reference number uniquely assigned to each transaction;*

> *If the number starts with C, it indicates a cancellation*

**StockCode**: *item code uniquely assigned to each distinct product*

**Description**: *product name*

**Quantity**: *quantities for each product per transaction*

**InvoiceDate**: *transaction date and time*

**UnitPrice**: *product price per unit*

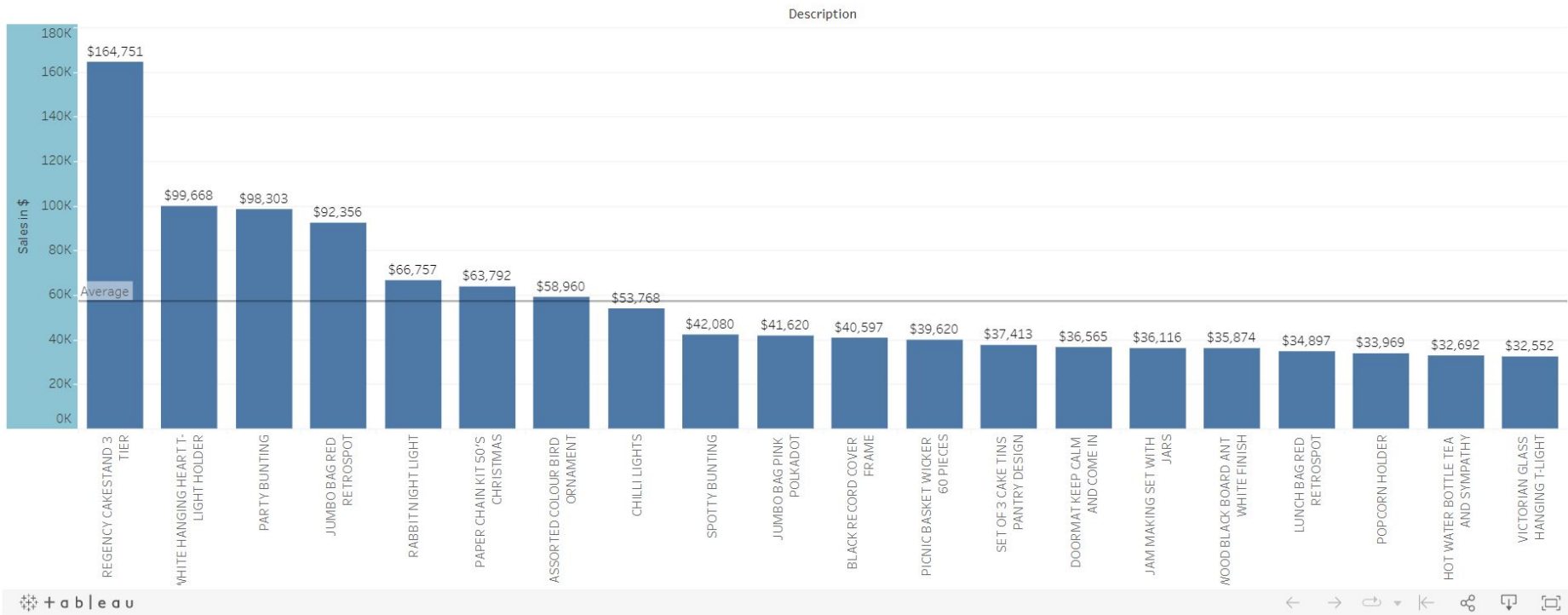**CustomerID**: *reference number uniquely assigned to each customer*

# Data Preprocessing

The first thing to manage before performing any analysis is to prepare the data. The data has been explored and preprocessed. And the following have been cleaned accordingly:

- Records where *InvoiceNo* starts with A were removed as these are related to bad debts adjustment.

- Records where *StockCode* is non-numeric were removed as these are not related to sales of products.

- Records where UnitPrice is 0 were removed as these are not valid sales transactions and mostly related to inventory adjustments, damages, etc.

- A *TotalSales* column was created. The value is the product of *Quantity* and *UnitPrice*.

- Invoices that starts with C are cancellations and most of the cancellations have a matching original invoice. For every customer where there is one cancellation and one matching invoice, the cancellation pairs (1:1) were removed.

- There are cases where for each customer, there is one cancellation and multiple matching invoices. The search was done for each customer for negative quantities below -100. All identified pairs were also removed.

- Individual columns were created for the extracted month, day, weekday and period. The date-only was also extracted from the InvoiceDate.

# What are the top-performing products and product categories?

Top 20 Most Purchased Products based on $ Sales
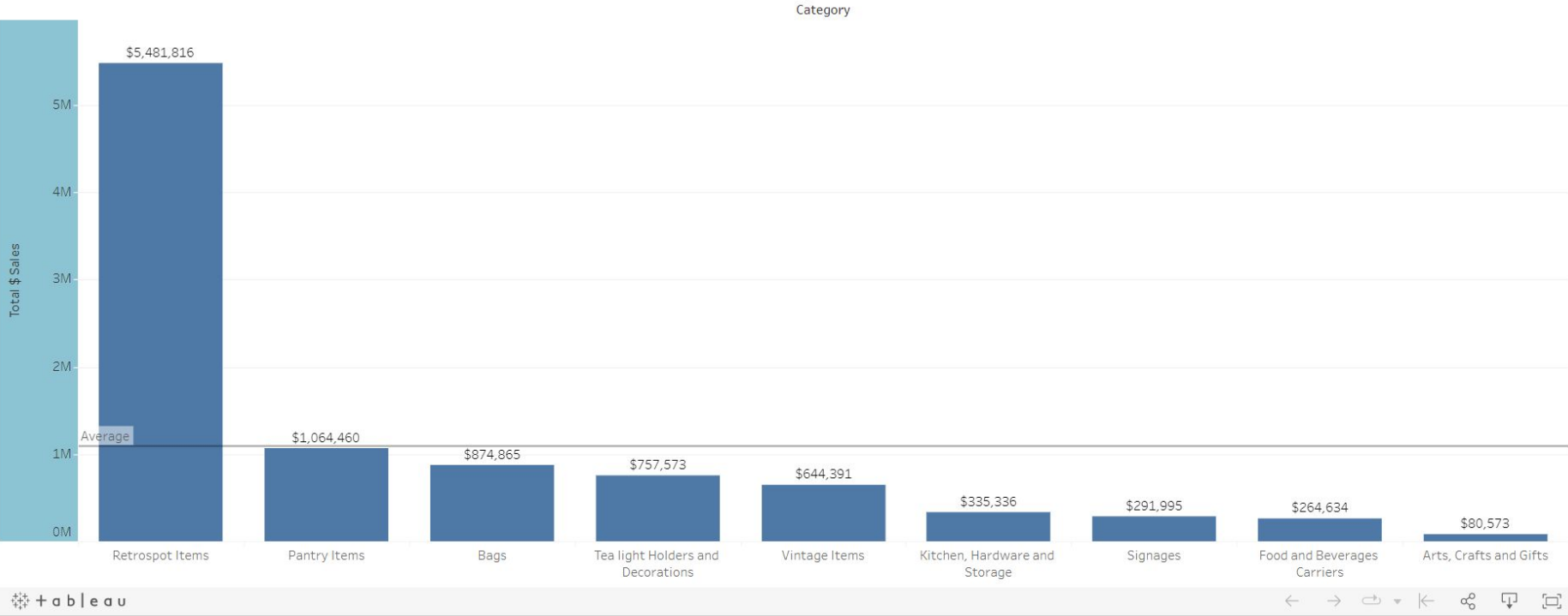
Description

# Top 20 Most Purchased Products based on Quantity



Link:
https://public.tableau.com/views/Top20MostPurchasedProductsbasedonQuantity/Top20MostPurchasedProductsbasedonQuantity?:language=en-US&:display_count=n&:origin=viz_share_link

# Sales Performance per Category based on $ Sales

Category



Total $ Sales

$5,481,816

5M

4M

3M

2M

Average

1M

$1,064,460

$874,865

$757,573

$644,391

$335,336

$291,995

$264,634

$80,573

0M

Retrospot Items | Pantry Items | Bags | Tea light Holders and Decorations | Vintage Items | Kitchen, Hardware and Storage | Signages | Food and Beverages Carriers | Arts, Crafts and Gifts

+ableau

Link: https://public.tableau.com/views/SalesPerformanceperCategorybasedonSales/SalesPerformanceperCategorybasedonSales?:language=en-US&:display_count=n&:origin=viz_share_link

# Sales Performance per Category based on Quantity Sold

Category



| | | |
|---|---|---|
| Quantity | | |

Retrospot Items — 2,975,736

Tea light Holders and Decorations — 522,259

Bags — 477,782

Pantry Items — 474,154

Vintage Items — 314,794

Kitchen, Hardware and Storage — 172,790

Arts, Crafts and Gifts — 168,589

Signages — 155,951

Food and Beverages Carriers — 53,711

Average

Link:
https://public.tableau.com/views/SalesPerformanceperCategorybasedonQuantitySold/SalesPerformanceperCategorybasedonQuantitySold?:language=en-US&:display_count=n&:origin=viz_share_link

# Top-performing Products and Product Categories

The top-performing **products** based on **Total $ Sales** are:

1. REGENCY CAKESTAND 3 TIER
2. WHITE HANGING HEART T-LIGHT HOLDER
3. PARTY BUNTING
4. JUMBO BAG RED RETROSPOT

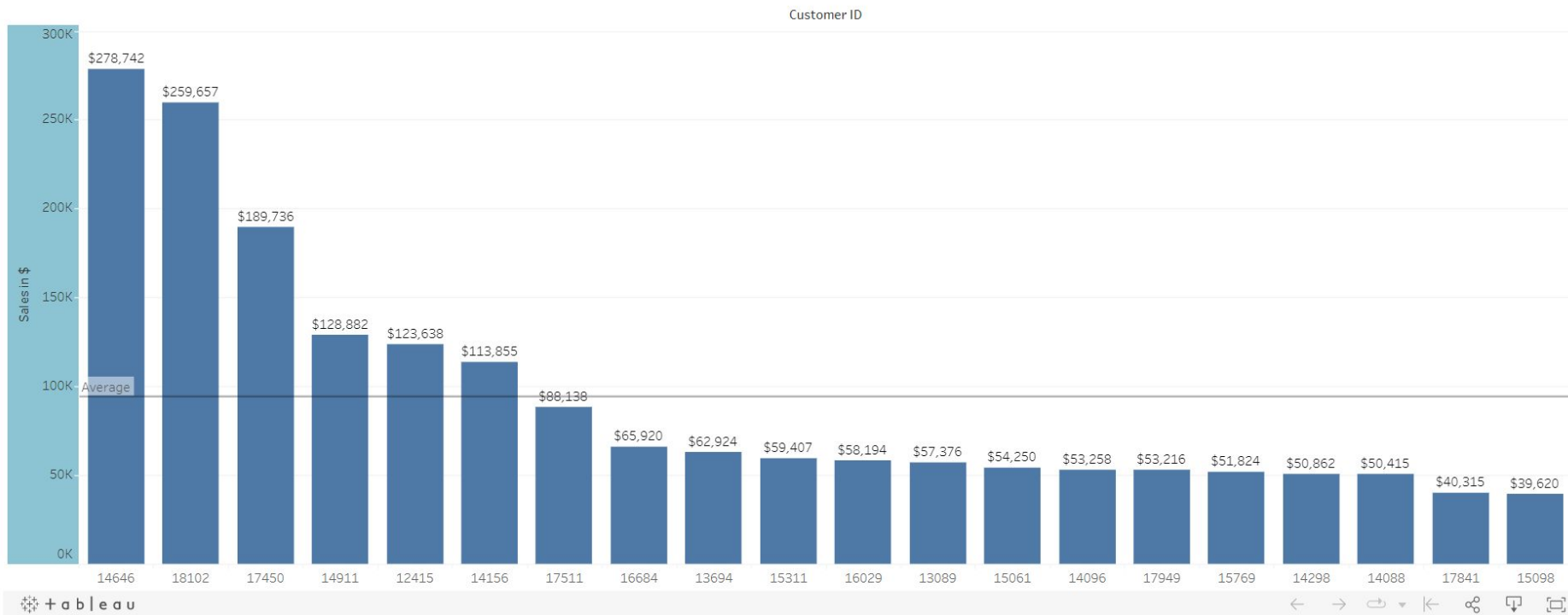The top-performing **products** based on **Quantity** are:

1. WORLD WAR 2 GLIDERS ASSTD DESIGNS
2. JUMBO BAG RED RETROSPOT.
3. ASSORTED COLOUR BIRD ORNAMENT
4. POPCORN HOLDER
5. PACK OF 72 RETROSPOT CAKE CASES
6. WHITE HANGING HEART T-LIGHT HOLDER

The top-performing **product category** based on both **Total $ Sales** and **Quantity** is:

- Retrospot Items

# Who are the top-performing customers?

## Top 20 Performing Customers based on $ Sales



Customer ID

| Customer ID | Sales in $ |
|---|---|
| 14646 | $278,742 |
| 18102 | $259,657 |
| 17450 | $189,736 |
| 14911 | $128,882 |
| 12415 | $123,638 |
| 14156 | $113,855 |
| 17511 | $88,138 |
| 16684 | $65,920 |
| 13694 | $62,924 |
| 15311 | $59,407 |
| 16029 | $58,194 |
| 13089 | $57,376 |
| 15061 | $54,250 |
| 14096 | $53,258 |
| 17949 | $53,216 |
| 15769 | $51,824 |
| 14298 | $50,862 |
| 14088 | $50,415 |
| 17841 | $40,315 |
| 15098 | $39,620 |

Average line at approximately 100K.

Link: https://public.tableau.com/views/Top20PerformingCustomersbasedonSales/Top20PerformingCustomersbasedonSales?:language=en-US&:display_count=n&:origin=viz_share_link

## Top 20 Performing Customers based on Quantity Purchased

Customer ID



| | |
|---|---|
| **196,556** | 14646 |
| 76,946 | 12415 |
| 76,848 | 14911 |
| 69,041 | 17450 |
| 64,124 | 18102 |
| 63,014 | 17511 |
| 61,808 | 13694 |
| 58,021 | 14298 |
| 57,013 | 14156 |
| 49,391 | 16684 |
| 37,720 | 15311 |
| 33,584 | 16422 |
| 32,320 | 17404 |
| 32,203 | 16029 |
| 32,184 | 16333 |
| 30,787 | 13089 |
| 28,638 | 15061 |
| 27,660 | 15769 |
| 27,572 | 17949 |
| 25,646 | 17381 |

Average

+ableau

Link:
https://public.tableau.com/views/Top20PerformingCustomersbasedonQuantityPurchased/Top20PerformingCustomersbasedonQuantityPurchased?:language=en-US&:display_count=n&:origin=viz_share_link

# Top-performing Customers

Based on **Total $ Sales**, the **top 5 customers** are:

1. 14646
2. 18102
3. 17450
4. 14911
5. 12415

Based on total **Quantity** purchased, the **top 5 customers** are:

1. 14646
2. 12415
3. 14911
4. 17450
5. 18102

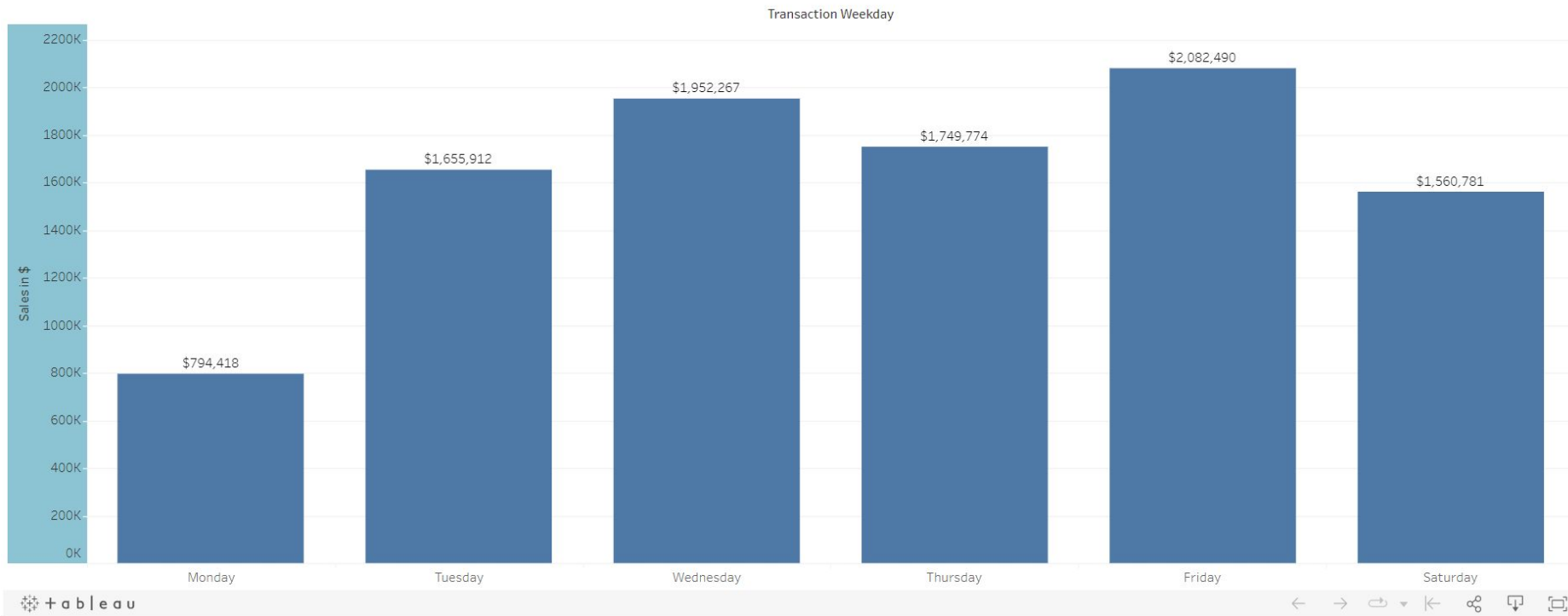Noticeably, they are the same set of customers.

# Monthly $ Sales



Transaction Period

| Transaction Period | Sales in $ |
|---|---|
| 2018-11 | $104,506 |
| 2018-12 | $655,192 |
| 2019-01 | $605,993 |
| 2019-02 | $495,073 |
| 2019-03 | $660,891 |
| 2019-04 | $487,770 |
| 2019-05 | $774,195 |
| 2019-06 | $684,711 |
| 2019-07 | $712,093 |
| 2019-08 | $736,125 |
| 2019-09 | $948,169 |
| 2019-10 | $1,126,241 |
| 2019-11 | $1,468,130 |
| 2019-12 | $336,553 |

+ a b l e a u

Link: https://public.tableau.com/views/MonthlySales_16598862693950/MonthlySales?:language=en-US&:display_count=n&:origin=viz_share_link

The dataset coverage is from 2018-11-29 to 2019-12-07. It is expected to have lower sales figures for 11-2018 and 12-2019 periods since these cover only a few days and not the full month. Surprisingly, the sales figures for the period of 12-2018 is relatively low as compared to sales figures in 09-2019, 10-2019 and 11-2019. 11-2019 is the best month.

# $ Sales per Day of the Week

**Transaction Weekday**



| | | | | | | |
|---|---|---|---|---|---|---|
| $794,418 | $1,655,912 | $1,952,267 | $1,749,774 | $2,082,490 | $1,560,781 |
| Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |

Link:

Most of the customers made their purchases during Fridays and Wednesdays but never on Sunday
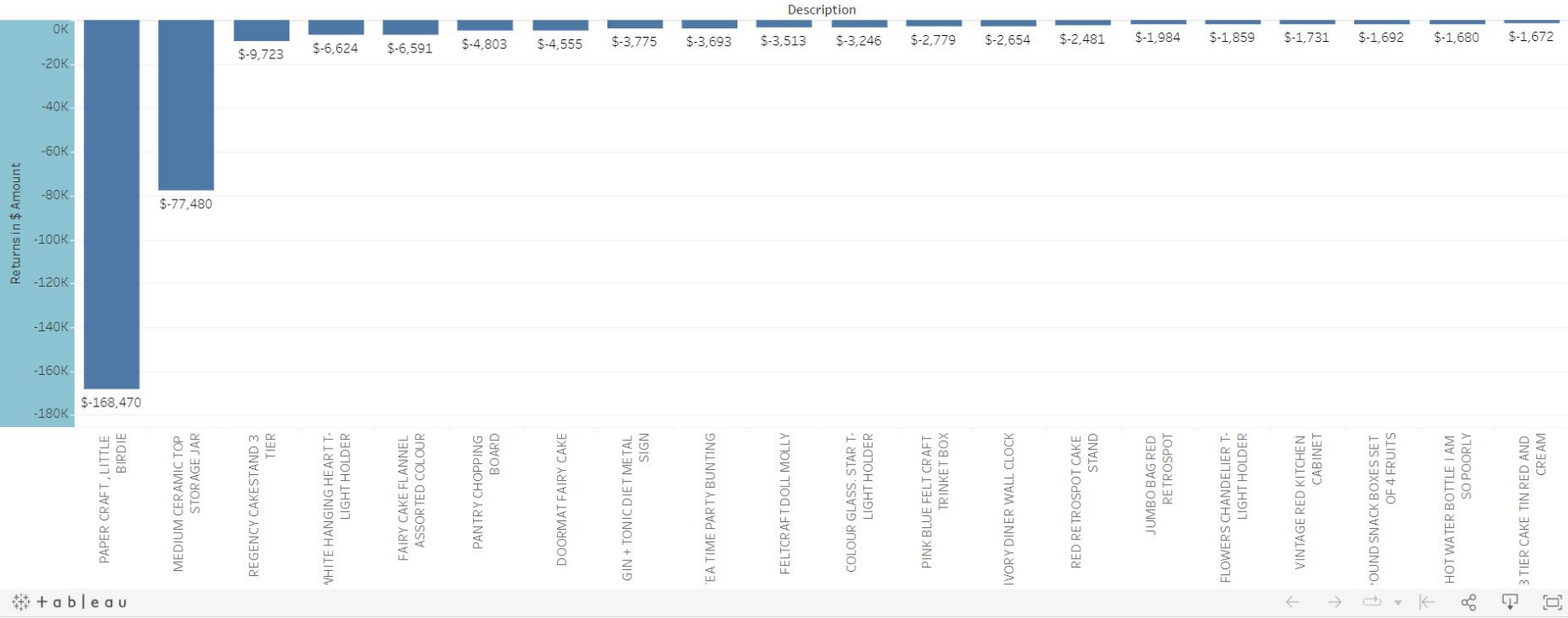
Daily Sales Metrics

The line charts display the changes for the number of customers, number of orders and total sales over the period of one year. The uptrend line is a good indication that the company's sales performance is increasing over time.
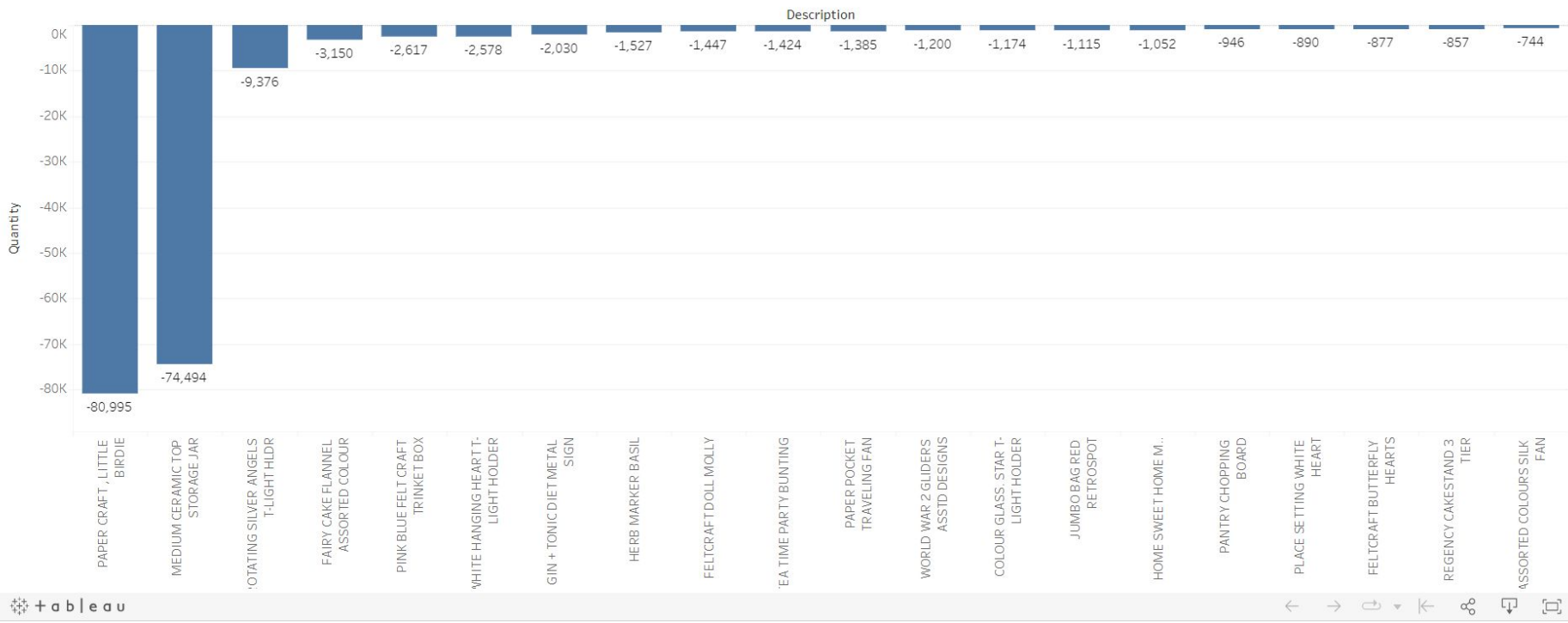
# Which products have the highest value and largest number of returns, and which items are most frequently returned?

Top 20 Product Returns based on $ Amount



Returns in $ Amount

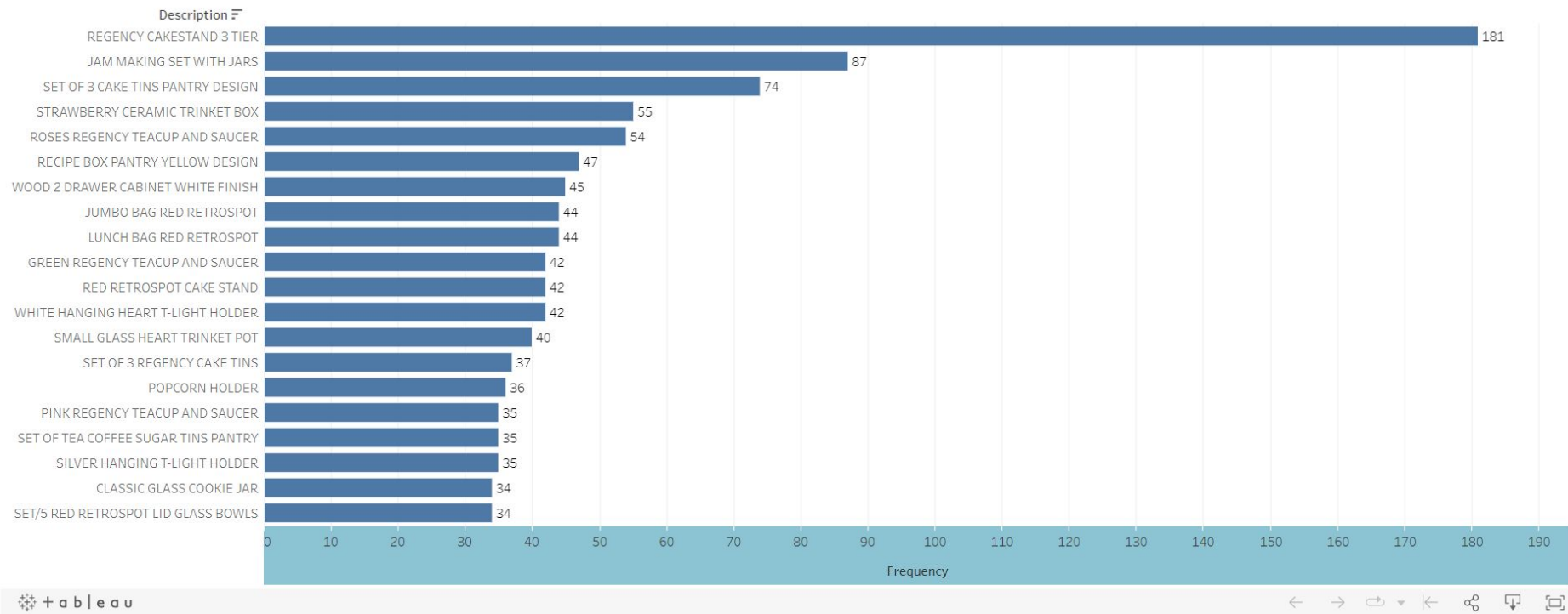| Description | Value |
|---|---|
| PAPER CRAFT , LITTLE BIRDIE | $-168,470 |
| MEDIUM CERAMIC TOP STORAGE JAR | $-77,480 |
| REGENCY CAKESTAND 3 TIER | $-9,723 |
| WHITE HANGING HEART T-LIGHT HOLDER | $-6,624 |
| FAIRY CAKE FLANNEL ASSORTED COLOUR | $-6,591 |
| PANTRY CHOPPING BOARD | $-4,803 |
| DOORMAT FAIRY CAKE | $-4,555 |
| GIN + TONIC DIET METAL SIGN | $-3,775 |
| TEA TIME PARTY BUNTING | $-3,693 |
| FELTCRAFT DOLL MOLLY | $-3,513 |
| COLOUR GLASS. STAR T-LIGHT HOLDER | $-3,246 |
| PINK BLUE FELT CRAFT TRINKET BOX | $-2,779 |
| IVORY DINER WALL CLOCK | $-2,654 |
| RED RETROSPOT CAKE STAND | $-2,481 |
| JUMBO BAG RED RETROSPOT | $-1,984 |
| FLOWERS CHANDELIER T-LIGHT HOLDER | $-1,859 |
| VINTAGE RED KITCHEN CABINET | $-1,731 |
| ROUND SNACK BOXES SET OF 4 FRUITS | $-1,692 |
| HOT WATER BOTTLE I AM SO POORLY | $-1,680 |
| 3 TIER CAKE TIN RED AND CREAM | $-1,672 |

✳ +ableau

Link: https://public.tableau.com/views/Top20ProductReturnsbasedonAmount/Top20ProductReturnsbasedonAmount?:language=en-US&:display_count=n&:origin=viz_share_link
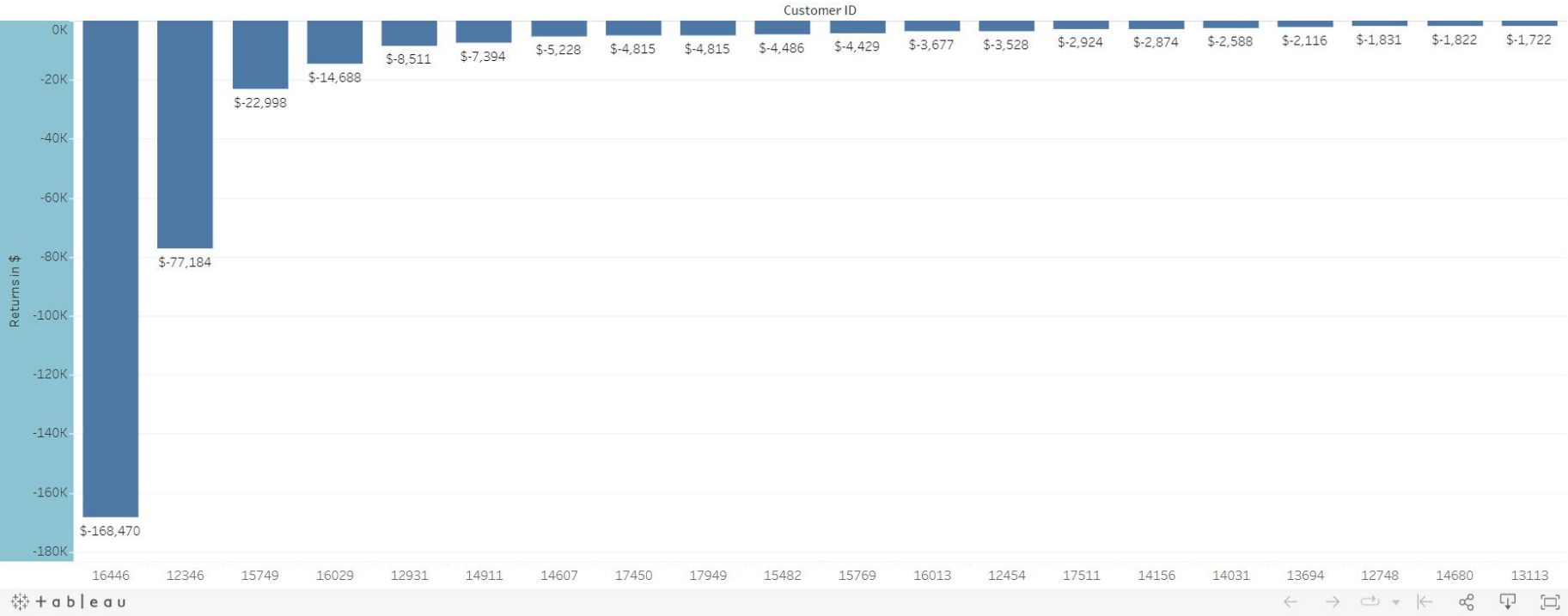
## Top 20 Product Returns based on Quantity



**Description**

| Value |
|---|
| -3,150 |
| -2,617 |
| -2,578 |
| -2,030 |
| -1,527 |
| -1,447 |
| -1,424 |
| -1,385 |
| -1,200 |
| -1,174 |
| -1,115 |
| -1,052 |
| -946 |
| -890 |
| -877 |
| -857 |
| -744 |
| -9,376 |
| -74,494 |
| -80,995 |

X-axis categories (left to right):
PAPER CRAFT, LITTLE BIRDIE | MEDIUM CERAMIC TOP STORAGE JAR | ROTATING SILVER ANGELS T-LIGHT HLDR | FAIRY CAKE FLANNEL ASSORTED COLOUR | PINK BLUE FELT CRAFT TRINKET BOX | WHITE HANGING HEART T-LIGHT HOLDER | GIN + TONIC DIET METAL SIGN | HERB MARKER BASIL | FELTCRAFT DOLL MOLLY | TEA TIME PARTY BUNTING | PAPER POCKET TRAVELING FAN | WORLD WAR 2 GLIDERS ASSTD DESIGNS | COLOUR GLASS. STAR T-LIGHT HOLDER | JUMBO BAG RED RETROSPOT | HOME SWEET HOME M.. | PANTRY CHOPPING BOARD | PLACE SETTING WHITE HEART | FELTCRAFT BUTTERFLY HEARTS | REGENCY CAKESTAND 3 TIER | ASSORTED COLOURS SILK FAN

+ableau

Link: https://public.tableau.com/views/Top20ProductReturnsbasedonQuantity/Top20ProductReturnsbasedonQuantity?:language=en-US&:display_count=n&:origin=viz_share_link

## Most Frequently Returned Products

| Description ⩫ | Frequency |
|---|---|
| REGENCY CAKESTAND 3 TIER | 181 |
| JAM MAKING SET WITH JARS | 87 |
| SET OF 3 CAKE TINS PANTRY DESIGN | 74 |
| STRAWBERRY CERAMIC TRINKET BOX | 55 |
| ROSES REGENCY TEACUP AND SAUCER | 54 |
| RECIPE BOX PANTRY YELLOW DESIGN | 47 |
| WOOD 2 DRAWER CABINET WHITE FINISH | 45 |
| JUMBO BAG RED RETROSPOT | 44 |
| LUNCH BAG RED RETROSPOT | 44 |
| GREEN REGENCY TEACUP AND SAUCER | 42 |
| RED RETROSPOT CAKE STAND | 42 |
| WHITE HANGING HEART T-LIGHT HOLDER | 42 |
| SMALL GLASS HEART TRINKET POT | 40 |
| SET OF 3 REGENCY CAKE TINS | 37 |
| POPCORN HOLDER | 36 |
| PINK REGENCY TEACUP AND SAUCER | 35 |
| SET OF TEA COFFEE SUGAR TINS PANTRY | 35 |
| SILVER HANGING T-LIGHT HOLDER | 35 |
| CLASSIC GLASS COOKIE JAR | 34 |
| SET/5 RED RETROSPOT LID GLASS BOWLS | 34 |

Frequency (0 — 190)

+ a b l e a u

Link: https://public.tableau.com/views/MostFrequentlyReturnedProducts/MostFrequentlyReturnedProducts?:language=en-US&:display_count=n&:origin=viz_share_link

# Who are the customers with the highest value and largest number of returns?

Top 20 Customers with Highest Returns based on $ Amount

Customer ID



Returns in $

| | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $-168,470 | $-77,184 | $-22,998 | $-14,688 | $-8,511 | $-7,394 | $-5,228 | $-4,815 | $-4,815 | $-4,486 | $-4,429 | $-3,677 | $-3,528 | $-2,924 | $-2,874 | $-2,588 | $-2,116 | $-1,831 | $-1,822 | $-1,722 |
| 16446 | 12346 | 15749 | 16029 | 12931 | 14911 | 14607 | 17450 | 17949 | 15482 | 15769 | 16013 | 12454 | 17511 | 14156 | 14031 | 13694 | 12748 | 14680 | 13113 |

+ableau

# Top 20 Customer with Highest Returns based on Quantity

Customer ID

| | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

- -80,995 (16446)
- -74,215 (12346)
- -9,360 (15838)
- -9,014 (15749)
- -8,004 (16029)
- -4,427 (12931)
- -3,768 (14607)
- -3,331 (14911)
- -2,878 (17949)
- -2,022 (15482)
- -2,012 (15769)
- -1,776 (12901)
- -1,594 (16013)
- -1,535 (12748)
- -1,535 (17511)
- -1,515 (16938)
- -1,504 (13694)
- -1,242 (14533)
- -1,228 (12607)
- -1,006 (12454)

Quantity axis: 0K, -10K, -20K, -30K, -40K, -50K, -60K, -70K, -80K

Customer ID axis: 16446, 12346, 15838, 15749, 16029, 12931, 14607, 14911, 17949, 15482, 15769, 12901, 16013, 12748, 17511, 16938, 13694, 14533, 12607, 12454

‡ + a b l e a u

Link:
https://public.tableau.com/views/Top20CustomerwithHighestReturnsbasedonQuantity/Top20CustomerwithHighestReturnsbasedonQuantity?:language=en-US&:display_count=n&:origin=viz_share_link

# What are the products and who are the customers who have significant impact on returns?

- **PAPER CRAFT; LITTLE BIRDIE** is the product with the **highest value** and **largest number** of returns. The total amount returned is **-$168,470** and the total quantity returned is **80,995**

- **REGENCY CAKESTAND 3 TIER** is the top product based on total sales. However, it is also the **most frequently returned product**. It was returned **181 times**.

- Customer **16446** is the customer with the **highest value** and **largest number** of returns. This is the customer who returned **80,995** pieces of the product **PAPER CRAFT; LITTLE BIRDIE with total returns amount of -$168,470**.

## Monthly Returns



Transaction Period

| Period | Returns |
|--------|---------|
| 2018-11 | $-1,839 |
| 2018-12 | $-15,769 |
| 2019-01 | $-92,292 |
| 2019-02 | $-8,048 |
| 2019-03 | $-11,514 |
| 2019-04 | $-32,048 |
| 2019-05 | $-9,622 |
| 2019-06 | $-13,580 |
| 2019-07 | $-11,271 |
| 2019-08 | $-22,880 |
| 2019-09 | $-16,743 |
| 2019-10 | $-45,039 |
| 2019-11 | $-25,208 |
| 2019-12 | $-172,872 |

+ableau

Surprisingly, the period 12–2019 has the highest amount of returns considering that this period has only 7 days. The total amount of returns is -$172,872.

# Customer Segmentation using RFM

Based on RFM scores, the customers are grouped into the following segments and here are the recommendations:

- **Top Customers** - Send birthday and anniversary (as being customer) cards with discount vouchers. Create a referral program where they can get a discount upon the first purchase of their referrals.

- **Active Customers** - Create loyalty rewards program for this group where they can earn point for every purchase and convert these points into discount or voucher.

- **Average Customers** - This group of customers are not consistently inactive or active. Send them emails with discounted items or promotional items for every occasion like their birthdays, mother's day, father's day, Christmas etc.

- **Customers at Risk** - This group of customers need attention as they made some purchases but it's been a long time since their last purchase. Send them personalized emails as well containing promotional items or free samples of products to try to encourage them to be more active.

- **Inactive Customers** - Send them personalized emails containing discounted items to encourage them to order and be active.

# Customers per Segments



The combined proportion of top, active and average customers is almost 60%. This is not bad as 60% of the customer base are making their purchases in a manner that can sustain the online store business.

However, the 40% requires attention. These customers need to be reactivated and encouraged to make more purchases and frequently.

# Product Segmentation using RFM

The products were grouped into 3 clusters based on RFM and K-Means clustering:

- **Cluster 0** - is the **worst-performing product group** with average **recency of 160 days**, average **frequency of 10 times** and average **sales of $167**.
- **Cluster 1** - is the **best-performing product group** with average **recency of 2.9 days**, average **frequency of 306 times**, and average **sales of $6,231.**
- **Cluster 2** - is quite good and not bad with average **recency of 15 days**, average **frequency of 58 times** and average **sales of $655**.

**Recommendations** for each cluster:

- **Cluster 0** - Review the products in this cluster, it probably needs to have more assortment or variety, current stocks are maybe in incomplete sizes or colors and maybe the products are out of season or outdated. Recency of 160 days is like more than 5 months since the last time the products were purchased. Frequency of only 10 times is like the products are being purchased not even once a month. Sales of 167 throughout the year is very low.

- **Cluster 1** - For Cluster 1 products, probably add more variety and increase the quantities. Monitor the inventories so these products won't go out of stock. Continue whatever the existing advertising or marketing campaigns for these products.

- **Cluster 2** - For Cluster 2, with recency of 15, it is good. It means the products are marketable or saleable currently. However, the frequency of 58 times is quite rare and sales of 655 is low, it's like the products are being purchased about 5 times a month. This can be increased by creating more advertising and promotional strategies to target specific customer demographics.

# Product Categorization | Product Category Analysis

Product categorization is important as the quality of product analysis depends heavily on the ability to accurately cluster similar products. To create categories, the following were performed on the texts in the product description.

- Text-preprocessing - lowercasing all letters and removing stop words.
- Text vectorization algorithm - NLP TF-IDF Vectorizer to transform text into vectors.
- Clustering - K-Means (Elbow method) to identify the number of clusters and to cluster the texts.
- Most important words - WordCloud to identify the most important/common words in each cluster.

Here are the category names created:

- Cluster 0 - Retrospot Items
- Cluster 1 - Tea light Holder and Decorations
- Cluster 2 - Arts, Crafts and Gifts
- Cluster 3 - Kitchen, Hardware and Storage
- Cluster 4 - Signages
- Cluster 5 - Pantry Items
- Cluster 6 - Vintage Items
- Cluster 7 - Bags
- Cluster 8 - Food and Beverage Carriers

## Sales Performance per Category based on $ Sales

Category



Bar chart showing Total $ Sales per category:

| Category | Total $ Sales |
|---|---|
| Retrospot Items | $5,481,816 |
| Pantry Items | $1,064,460 |
| Bags | $874,865 |
| Tea light Holders and Decorations | $757,573 |
| Vintage Items | $644,391 |
| Kitchen, Hardware and Storage | $335,336 |
| Signages | $291,995 |
| Food and Beverages Carriers | $264,634 |
| Arts, Crafts and Gifts | $80,573 |

Average line at approximately $1M

+‡‡+ +ableau

Link: https://public.tableau.com/shared/QH4RG84HT?:display_count=n&:origin=viz_share_link

Retrospot items is the category that contributes significantly to total sales. More than half of the total sales are from Retrospot items.

The worst-performing is Arts, Craft and Gifts category

## Sales Performance per Category based on Quantity Sold

Category

The top-performing category based on quantity sold is Retrospot items which is almost 56% of the total quantity sold.

The worst-performing is Food and Beverages category.

# Monthly $ Sales Performance per Category

Retrospot items is the leading category every month. Retrospot's highest sales is for the period of 11-2019

# Average Monthly $ Sales per Category

Category



Average Monthly Sales

- Retrospot Items: $456,818
- Pantry Items: $88,705
- Bags: $72,905
- Tea light Holders and Decorations: $63,131
- Vintage Items: $53,699
- Kitchen, Hardware and Storage: $27,945
- Signages: $24,333
- Food and Beverages Carriers: $22,053
- Arts, Crafts and Gifts: $6,714

Retrospot items has the highest average monthly sales and Arts, Crafts and Gifts category has the lowest.

# Unique Products per Category



**Category**
- Arts, Crafts and Gifts
- Bags
- Food and Beverages C...
- Kitchen, Hardware an...
- Pantry Items
- Retrospot Items
- Signages
- Tea light Holders and ...
- Vintage Items

Retrospot Items
2,857

Pantry Items
333

Signages
68

Arts, Crafts and Gifts
67

Bags
67

Food and Beverages Carriers
25

Tea light Holders and Decorations
248

Kitchen, Hardware and Storage
140

Vintage Items
205

+ableau

And the reason for having Retrospot items as the top-performing category, is that it has the largest variety of products with 2857 unique items.

# Retrospot Items



Retrospot Items - Monthly Sales

# Statistical Hypothesis 1:
## *High-priced products contribute to higher sales than the low-priced products*

**Null Hypothesis:** There is no difference between the mean sales of low-priced products and high-priced products

**Alternative Hypothesis:** Mean sales of low-priced products is different than the mean sales of high-priced products

**Criteria for Decision:** alpha = 0.05
- Accept null hypothesis if p-value > alpha
- Reject null hypothesis if p-value < alpha

**Based on the data, we reject the null hypothesis:**

Based on the result of statistical test and excluding the outliers, we found out that there is a difference between the average sales of low-priced products and average sales of high-priced products. Low price values are unit prices below 20 and high price values are unit prices 20 and above. This could also mean that the average sales of low-priced products can be higher than the average sales of high-priced products or vice-versa.

# Statistical Hypothesis 2:
*Sales are higher during Christmas season (December) as compared to other months*

**Null Hypothesis:** There is no difference between the average sales in December and the average sales in other months

**Alternative Hypothesis:** The average sales in December is different than the average sales in other months
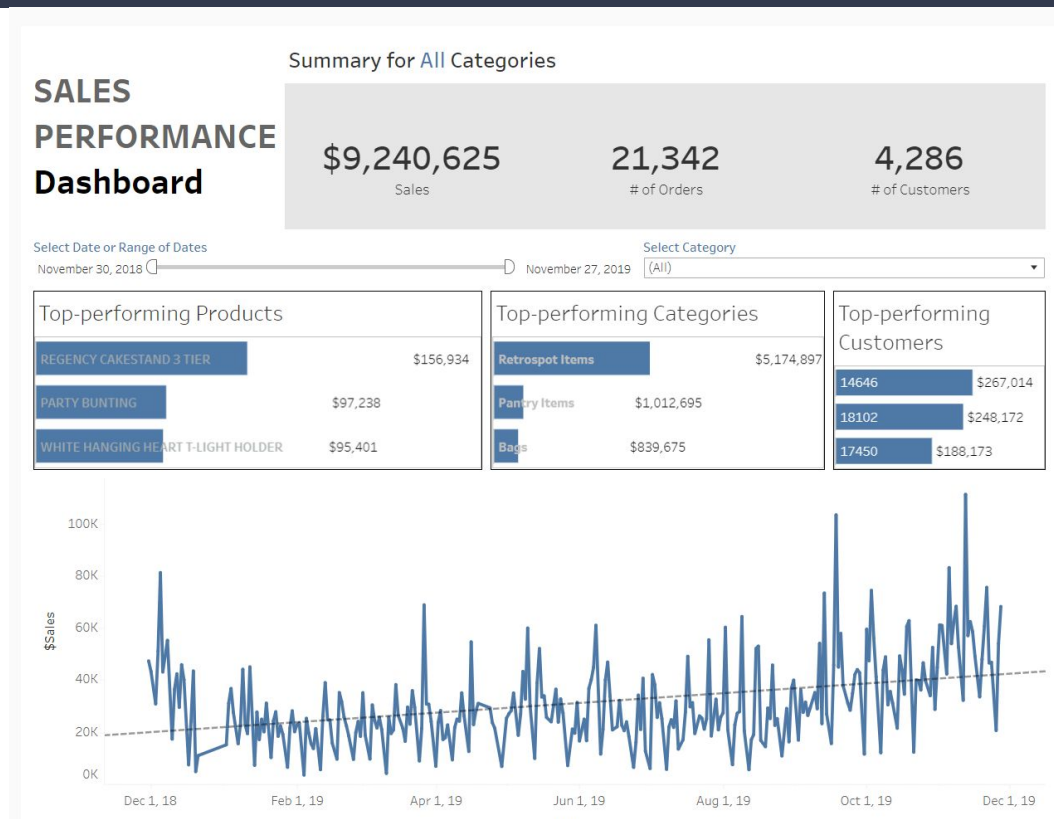
**Criteria for Decision:** alpha = 0.05
- Accept null hypothesis if p-value > alpha
- Reject null hypothesis if p-value < alpha

**Based on the data, we reject the null hypothesis:**
Based on the result of statistical test and excluding the outliers, we found out that there is a difference between the average sales generated in December and average sales in other months. This could also mean that the average sales in December can be higher than the average sales in other months or vice-versa.

# Dashboard

# Recommendations

**Returns**
For e-commerce business, customer returns are unavoidable. That's why it is important for any e-commerce companies to have a good and fair return policy. Returns are a cost factor. It negatively affects revenues and profit. For this reason, it is recommended for the company to perform monthly returns analysis. Analysis of returns is a collection of data about the return of a product. It helps to understand how big or small the impact of returns, find out the cause and manage it.

**Customer Segmentation**
It is recommended to perform customer segmentation using RFM scores. It helps the business to analyze their existing customers in terms of their spending habits. When customers are segmented, it is easier to target specific customers with communications or marketing strategies that are relevant based on data about a particular set of customer behaviors.

**Product Segmentation**
It is also recommended to segment the products into clusters based on the product's performance. It will help to analyze which products are currently marketable and which are generating high or low revenues.

**Product Categorization**
Product categorization plays an important role specifically for e-commerce. The quality of product analysis depends heavily on the ability to accurately cluster similar products. It also helps customers intuitively find what they want from the website. If products are not categorized properly, it may result to loss of sales opportunity or loss of customers as they couldn't find the items that they want in the category they're looking at.

# Limitations

- If I would be the future owner of the business or if I am the current owner of the business, I would be more interested on the profit analysis. The top performing products or top performing customers could be different based on profit. Profit is computed by subtracting the costs from the revenue. As the data doesn't have information about costs, the analysis is limited only on the unit price, quantity and sales amount.

- The dataset is only a year of sales transactions history. It would be better for comparison analysis if it is at least 3 years sales transaction history. Like we can compare the performance of the store or the products every year.

- There is no information about the cancellation or return policy. The dataset has a significant number of cancellations/returns and it would be easier to understand how to manage those if there were additional information.

- There is no information about the store's pricing policy. As I have observed, some stockcodes have multiple prices which I assume, it's either a price increase or a distinction between a wholesale price and retail price. Some prices are higher if there is no customer id.

- I wanted to use BERTopic or spacy but for some reason, I had issues with installing bertopic and spacy took a very long time to run and it caused my computer to freeze. I would probably explore more on other NLP or ML tools had I better computer resources and longer time.

# Resources

1. To get ideas on how to perform product categorization
   https://techblog.commercetools.com/boosting-product-categorization-with-machine-learning-ad4dbd30b0e8

2. How to deal with negative values in the logarithmic transformation to unskew the data
   https://campus.datacamp.com/courses/customer-segmentation-in-python/data-pre-processing-for-clustering?ex=4

3. How to interpret the p-value of 0.0
   https://www.statology.org/here-is-how-to-interpret-a-p-value-of-0-000/

4. How to find cancellation pairs
   https://stackoverflow.com/questions/38831088/remove-cancelling-rows-from-pandas-dataframe

5. To understand what is DOTCOM POSTAGE, DCGS (dotcomgiftshop) and Retrospot in the dataset
   https://www.rexlondon.com/blog/dotcomgiftshop-changing-to-rex-london

6. To decide whether to use or not to use lemmatization
   https://www.datacamp.com/tutorial/stemming-lemmatization-python

7. How to prepare data for hypothesis testing
   https://towardsdatascience.com/hypothesis-testing-the-discount-bump-4b8b6b4c4fec