

Question 4

Perform a full join will retain all rows from both datasets, matching them based on the order_id, which ensures that every customer and every order will be included in the resulting dataset. By filtering the merged dataset to include only rows where the order_id column is missing, finding customers without orders. Similarly, by filtering the merged dataset to include only rows where the customer_id column is missing, is used to find orders without corresponding customer details. But both the result shown 0 rows, which means there is no customers without orders or orders without customers detail in the dataset given.

```
> customers_without_orders
[1] customer_id          customer_unique_id    customer_zip_code_prefix
[4] customer_city        customer_state        order_id
[7] order_status         order_purchase_timestamp order_approved_at
[10] order_delivered_carrier_date order_delivered_customer_date order_estimated_delivery_date
<0 rows> (or 0-length row.names)
> orders_without_customer_details
[1] customer_id          customer_unique_id    customer_zip_code_prefix
[4] customer_city        customer_state        order_id
[7] order_status         order_purchase_timestamp order_approved_at
[10] order_delivered_carrier_date order_delivered_customer_date order_estimated_delivery_date
<0 rows> (or 0-length row.names)
```

Question 5

Perform a semi join in this case will returns only the rows from olist_sellers_dataset where there is a match with the olist_order_items_dataset based on seller_id column. I first determine the total number of unique sellers in each dataset by counting the number of unique seller_id values. Then find the number of common sellers between the datasets by merging them based on the seller_id column and counting the number of unique sellers. After comparing the common sellers and total olist sellers, I found that they are 100% similar, which means that there is no seller being inactive, all the sellers have at least made 1 sales.

```
> similarity_percentage
[1] 100
```

Question 6

Perform an anti-join will compares two datasets based on customer_id and returns rows from the olist_customers_dataset that do not have matching keys in olist_orders_dataset. The result shown that the matched row is 0, means every customer ID from the olist_customers_dataset has a corresponding entry in the olist_orders_dataset, so there is no customer haven't placed an order.

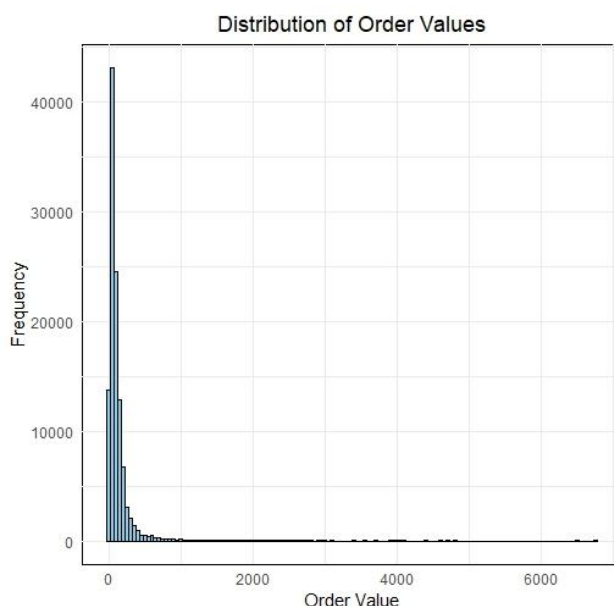
```
> customers_without_orders
[1] customer_id      customer_unique_id  customer_zip_code_prefix
[4] customer_city    customer_state
<0 rows> (or 0-length row.names)
```

Question 7

I first use an inner join to connect the olist_orders_dataset with the olist_order_items_dataset based on the common key order_id, which joins each order with its corresponding order items. Then use left join with olist_products_dataset to adds product details to each order item. Finally, I use right join with olist_sellers_dataset to add seller details to each order. By analysing the comprehensive dataset, I found that the flow between seller and product is there is 11.13 products sold per seller, which means that on average, one seller sells about 11 different products. The flow between products and order, is there is 1.04 products have been placed within an order, which indicates that the buyer usually only buy a single products per order, we can also observe this on the plot shown below, the order value at 1 having the highest frequency.

Average number of products sold per seller: 11.13021

Average number of products per order: 1.038098



Section B

Question 1-4

Chart generated in R studio.

Question 5

The bar chart depicts the global market trend. The United States leads the way with a staggering 1,000 sales, demonstrating its strong market presence and high level of consumption. In Europe, countries such as France, Germany, Spain, and the United Kingdom also have large consumer markets, indicating that these countries have stable economies and sizable consumer numbers. Japan and Singapore are the only two countries in the Asia-Pacific region that performed better, with sales lower than those of the countries, possibly due to their smaller populations. Conversely, countries such as Ireland and the Philippines have lower sales, possibly due to smaller populations or limited market coverage. However, countries such as Austria, Belgium, Denmark, and Norway are promising countries where companies can explore emerging markets and expand their sales channels.