

Stage 5

Group Members: M Tausif Tajwar Bhuiyan
Luke Rodriguez
Dinh Nhat Nguyen

November 10, 2025

Contents

1 Market and Sales Discoveries	3
1.1 Query 1.1: High-Value Customer States	3
1.2 Query 1.2: Top Selling Categories by Revenue	3
1.3 Query 1.3: States with Customers, No Orders	3
1.4 Query 1.4: Highest Sales Geolocation	4
2 Seller Performance and Network	5
2.1 Query 2.1: Seller Success Rate by Review Score	5
2.2 Query 2.2: Unused Product Catalog	5
2.3 Query 2.3: States with Customers, No Sellers	6
2.4 Query 2.4: Single-Product Sellers	6
3 Order and Review Quality	7
3.1 Query 3.1: Order Review Rate	7
3.2 Query 3.2: Worst Rated Product Category	7
3.3 Query 3.3: Orders Paid in Full (Single Payment)	7
3.4 Query 3.4: Review Score Extremes (1-Star vs. 5-Star)	8
4 Payments and Transactions	9
4.1 Query 4.1: Most Common Payment Type by State	9
4.2 Query 4.2: Average Number of Installments per Payment Type	9
5 Customer Behavior and Loyalty	10
5.1 Query 5.1: Repeat Purchase Customers	10
5.2 Query 5.2: Average Time Between Orders (Returning Customers)	10
6 Additional Queries	12
6.1 Query 6.1: Delivery Performance by State	12
6.2 Query 6.2: Category Performance by Quarter	12
6.3 Query 6.3: Revenue by State and Year	13
6.4 Query 6.4: Seller Inventory Check	13

1 Market and Sales Discoveries

1.1 Query 1.1: High-Value Customer States

Explanation: Find the states that have the highest average order value (AOV), showing where customers spend the most per transaction.

```

1 SELECT
2     g.geolocation_state AS state,
3     s.state_name,
4     COUNT(DISTINCT o.order_id) AS total_orders,
5     AVG(oi.price + oi.freight_value) AS avg_order_value
6 FROM CUSTOMERS c
7 JOIN GEOLOCATION g ON c.customer_zip_code_prefix = g.zip_code_prefix
8 JOIN STATES s ON g.geolocation_state = s.state_code
9 JOIN ORDERS o ON c.customer_id = o.customer_id
10 JOIN ORDER_ITEMS oi ON o.order_id = oi.order_id
11 GROUP BY g.geolocation_state, s.state_name
12 ORDER BY avg_order_value DESC
13 LIMIT 10;

```

Interestingness to Analysts: This query shows geographic markets where customers have higher purchasing power or willingness to spend, that indicates targeted marketing campaigns, premium product launches, and strategic resource allocation to high-value regions. States with high AOV may require increased seller recruitment, premium inventory, or specialized customer service investments.

(returns top 10 states by default)

1.2 Query 1.2: Top Selling Categories by Revenue

Explanation: Determine which product categories account for the largest total revenue overall.

```

1 SELECT
2     c.category_name_english,
3     COUNT(DISTINCT oi.order_id) AS orders_count,
4     COUNT(DISTINCT oi.product_id) AS products_count,
5     SUM(oi.price) AS total_revenue,
6     AVG(oi.price) AS avg_product_price
7 FROM CATEGORIES c
8 JOIN PRODUCTS p ON c.category_id = p.category_id
9 JOIN ORDER_ITEMS oi ON p.product_id = oi.product_id
10 GROUP BY c.category_id, c.category_name_english
11 ORDER BY total_revenue DESC
12 LIMIT 15;

```

Interestingness to Analysts: Understanding revenue distribution across categories helps prioritize inventory investments, identify growth opportunities, and allocate marketing budgets effectively.

(returns top 15 categories)

1.3 Query 1.3: States with Customers, No Orders

Explanation: Identify states that have registered customers but have zero recorded orders from those residents, highlighting untapped potential.

```

1 SELECT
2     s.state_code,
3     s.state_name,
4     COUNT(DISTINCT c.customer_id) AS registered_customers
5 FROM STATES s

```

```

6 JOIN GEOLOCATION g ON s.state_code = g.geolocation_state
7 JOIN CUSTOMERS c ON g.zip_code_prefix = c.customer_zip_code_prefix
8 WHERE NOT EXISTS (
9     SELECT 1
10    FROM ORDERS o
11   WHERE o.customer_id = c.customer_id
12 )
13 GROUP BY s.state_code, s.state_name
14 HAVING COUNT(DISTINCT c.customer_id) > 0
15 ORDER BY registered_customers DESC;

```

Why Interesting to Analysts: This identifies market product failures where customer acquisition succeeded but conversion did not, suggesting issues with product availability, delivery logistics, payment methods, or customer experience.

1.4 Query 1.4: Highest Sales Geolocation

Explanation: Pinpoint the single most active geolocation (specific zip code) associated with the highest number of customer orders.

```

1 SELECT
2     g.zip_code_prefix,
3     g.geolocation_city,
4     g.geolocation_state,
5     g.geolocation_lat,
6     g.geolocation_lng,
7     COUNT(DISTINCT o.order_id) AS total_orders,
8     COUNT(DISTINCT c.customer_id) AS unique_customers,
9     SUM(oi.price + oi.freight_value) AS total_revenue
10    FROM GEOLOCATION g
11   JOIN CUSTOMERS c ON g.zip_code_prefix = c.customer_zip_code_prefix
12   JOIN ORDERS o ON c.customer_id = o.customer_id
13   JOIN ORDER_ITEMS oi ON o.order_id = oi.order_id
14 GROUP BY g.zip_code_prefix, g.geolocation_city, g.geolocation_state,
15           g.geolocation_lat, g.geolocation_lng
16 ORDER BY total_orders DESC
17 LIMIT 1;

```

Why Interesting to Analysts: Identifying the single hottest location showing local strategies like establishing pickup points, positioning warehouse inventory, partnering with local influencers, or studying customer demographics in that specific area to replicate success elsewhere.

2 Seller Performance and Network

2.1 Query 2.1: Seller Success Rate by Review Score

Explanation: Calculate the average review score for each seller, but only using orders that actually received a review to assess quality.

```

1 SELECT
2     s.seller_id,
3     g.geolocation_city AS seller_city,
4     g.geolocation_state AS seller_state,
5     COUNT(DISTINCT oi.order_id) AS orders_fulfilled,
6     COUNT(DISTINCT r.review_id) AS reviews_received,
7     AVG(r.review_score) AS avg_review_score,
8     ROUND(COUNT(DISTINCT r.review_id) * 100.0 /
9             COUNT(DISTINCT oi.order_id), 2) AS review_rate_percentage
10    FROM SELLERS s
11    JOIN GEOLOCATION g ON s.seller_zip_code_prefix = g.zip_code_prefix
12    JOIN ORDER_ITEMS oi ON s.seller_id = oi.seller_id
13    LEFT JOIN ORDER_REVIEWS r ON oi.order_id = r.order_id
14    GROUP BY s.seller_id, g.geolocation_city, g.geolocation_state
15    HAVING COUNT(DISTINCT r.review_id) > 0
16    ORDER BY avg_review_score DESC, orders_fulfilled DESC
17    LIMIT 20;

```

Why Interesting to Analysts: This shows top-performing sellers who could be featured prominently, rewarded with lower commission rates, or studied as best-practice examples. Low-scoring sellers may need quality improvement interventions, additional training, or potential removal from the platform.

(returns top 20 sellers)

2.2 Query 2.2: Unused Product Catalog

Explanation: List all existing products that have never been included in any order item, indicating dead inventory.

```

1 SELECT
2     p.product_id,
3     c.category_name_english AS category,
4     p.product_weight_g,
5     p.product_length_cm,
6     p.product_height_cm,
7     p.product_width_cm,
8     p.product_photos_qty
9     FROM PRODUCTS p
10    JOIN CATEGORIES c ON p.category_id = c.category_id
11 WHERE NOT EXISTS (
12     SELECT 1
13     FROM ORDER_ITEMS oi
14     WHERE oi.product_id = p.product_id
15 )
16 ORDER BY c.category_name_english, p.product_id;

```

Why Interesting to Analysts: Dead inventory represents wasted catalog space, missed opportunities, and potential issues with product listing quality, pricing, or categorization. This list helps to target actions like removing inactive products, repricing underperformers, improving product descriptions/photos, or investigating why certain product types fail to attract buyers.

2.3 Query 2.3: States with Customers, No Sellers

Explanation: Identify states that have customers but do not have any registered sellers, highlighting a potential fulfillment gap.

```

1 SELECT
2     s.state_code,
3     s.state_name,
4     COUNT(DISTINCT c.customer_id) AS customer_count
5 FROM STATES s
6 JOIN GEOLOCATION g ON s.state_code = g.geolocation_state
7 JOIN CUSTOMERS c ON g.zip_code_prefix = c.customer_zip_code_prefix
8 WHERE NOT EXISTS (
9     SELECT 1
10    FROM SELLERS sel
11   JOIN GEOLOCATION g2 ON sel.seller_zip_code_prefix = g2.zip_code_prefix
12   WHERE g2.geolocation_state = s.state_code
13 )
14 GROUP BY s.state_code, s.state_name
15 ORDER BY customer_count DESC;

```

Why Interesting to Analysts: This identifies supply-side gaps where customer demand exists but local seller infrastructure is absent, leading to longer delivery times, higher shipping costs, and reduced customer satisfaction.

2.4 Query 2.4: Single-Product Sellers

English Explanation: Count how many sellers offer only one unique product in their entire portfolio.

```

1 SELECT
2     COUNT(*) AS single_product_seller_count
3 FROM (
4     SELECT
5         oi.seller_id,
6         COUNT(DISTINCT oi.product_id) AS unique_products
7     FROM ORDER_ITEMS oi
8     GROUP BY oi.seller_id
9     HAVING COUNT(DISTINCT oi.product_id) = 1
10 ) AS single_product_sellers;

```

Why Interesting to Analysts: Single-product sellers may represent specialists, or new sellers testing the platform. This metric helps understand seller diversity and platform maturity.

3 Order and Review Quality

3.1 Query 3.1: Order Review Rate

Explanation: Determine what percentage of total orders end up receiving at least one customer review.

```

1 SELECT
2     COUNT(DISTINCT o.order_id) AS total_orders,
3     COUNT(DISTINCT r.order_id) AS reviewed_orders,
4     ROUND(COUNT(DISTINCT r.order_id) * 100.0 /
5             COUNT(DISTINCT o.order_id), 2) AS review_rate_percentage
6 FROM ORDERS o
7 LEFT JOIN ORDER_REVIEWS r ON o.order_id = r.order_id;

```

Why Interesting to Analysts: Review rate is a engagement metric indicating customer satisfaction, platform trust, and feedback loop health.

3.2 Query 3.2: Worst Rated Product Category

Explanation: Find the product category that has the lowest average customer review score across all time.

```

1 SELECT
2     c.category_name_english AS category,
3     COUNT(DISTINCT p.product_id) AS products_in_category,
4     COUNT(DISTINCT r.review_id) AS total_reviews,
5     AVG(r.review_score) AS avg_review_score,
6     SUM(CASE WHEN r.review_score <= 2 THEN 1 ELSE 0 END) AS negative_reviews,
7     SUM(CASE WHEN r.review_score >= 4 THEN 1 ELSE 0 END) AS positive_reviews
8 FROM CATEGORIES c
9 JOIN PRODUCTS p ON c.category_id = p.category_id
10 JOIN ORDER_ITEMS oi ON p.product_id = oi.product_id
11 JOIN ORDER_REVIEWS r ON oi.order_id = r.order_id
12 GROUP BY c.category_id, c.category_name_english
13 HAVING COUNT(DISTINCT r.review_id) >= 10
14 ORDER BY avg_review_score ASC
15 LIMIT 1;

```

Why Interesting to Analysts: Identifying the worst-performing category reveals systemic quality issues, seller problems, or customer expectation mismatches requiring immediate attention.

3.3 Query 3.3: Orders Paid in Full (Single Payment)

Explanation: Count the number of orders that were covered by exactly one payment record, indicating straightforward transactions.

```

1 SELECT
2     COUNT(*) AS single_payment_orders
3 FROM (
4     SELECT
5         order_id,
6         COUNT(*) AS payment_count
7     FROM ORDER_PAYMENTS
8     GROUP BY order_id
9     HAVING COUNT(*) = 1
10 ) AS single_payments;

```

Why Interesting to Analysts: Single-payment orders represent simpler, lower-risk transactions with less administrative overhead, fewer failure points, and reduced customer confusion.

3.4 Query 3.4: Review Score Extremes (1-Star vs. 5-Star)

Explanation: Report the total count of 1-star reviews and 5-star reviews to quickly gauge the extremes of customer sentiment.

```
1 SELECT
2   SUM(CASE WHEN review_score = 1 THEN 1 ELSE 0 END) AS one_star_reviews,
3   SUM(CASE WHEN review_score = 5 THEN 1 ELSE 0 END) AS five_star_reviews,
4   ROUND(SUM(CASE WHEN review_score = 1 THEN 1 ELSE 0 END) * 100.0 /
5         COUNT(*), 2) AS one_star_percentage,
6   ROUND(SUM(CASE WHEN review_score = 5 THEN 1 ELSE 0 END) * 100.0 /
7         COUNT(*), 2) AS five_star_percentage,
8   COUNT(*) AS total_reviews
9 FROM ORDER_REVIEWS;
```

Why Interesting to Analysts: Extreme ratings reveal polarized customer experiences that is high 5-star counts indicate strong satisfaction and product-market fit, while high 1-star counts signal serious problems requiring urgent intervention.

4 Payments and Transactions

4.1 Query 4.1: Most Common Payment Type by State

Explanation: Determine which payment method (credit card, boleto, etc.) is most frequently used by customers in each state.

```

1 WITH state_payment_counts AS (
2     SELECT
3         g.geolocation_state,
4         s.state_name,
5         op.payment_type,
6         COUNT(*) AS payment_count
7     FROM CUSTOMERS c
8     JOIN GEOLOCATION g ON c.customer_zip_code_prefix = g.zip_code_prefix
9     JOIN STATES s ON g.geolocation_state = s.state_code
10    JOIN ORDERS o ON c.customer_id = o.customer_id
11    JOIN ORDER_PAYMENTS op ON o.order_id = op.order_id
12    GROUP BY g.geolocation_state, s.state_name, op.payment_type
13),
14 ranked_payments AS (
15     SELECT
16         geolocation_state,
17         state_name,
18         payment_type,
19         payment_count,
20         ROW_NUMBER() OVER (PARTITION BY geolocation_state
21                             ORDER BY payment_count DESC) AS rank
22     FROM state_payment_counts
23)
24 SELECT
25     geolocation_state,
26     state_name,
27     payment_type AS most_common_payment_type,
28     payment_count
29 FROM ranked_payments
30 WHERE rank = 1
31 ORDER BY geolocation_state;

```

Why Interesting to Analysts: Payment method preferences vary significantly by region due to banking infrastructure, credit systems, cultural preferences, and socioeconomic factors. Understanding state-level payment preferences enables optimized payment gateway integration, targeted financial partnership strategies.

4.2 Query 4.2: Average Number of Installments per Payment Type

English Explanation: Compute the average installment count grouped by payment type to understand financing behaviors.

```

1 SELECT
2     payment_type,
3     COUNT(*) AS total_payments,
4     AVG(payment_installments) AS avg_installments,
5     MIN(payment_installments) AS min_installments,
6     MAX(payment_installments) AS max_installments,
7     ROUND(AVG(payment_value), 2) AS avg_payment_value
8 FROM ORDER_PAYMENTS
9 GROUP BY payment_type
10 ORDER BY avg_installments DESC;

```

Why Interesting to Analysts: Installment behavior reveals customer financial constraints, purchasing power, and payment preferences.

5 Customer Behavior and Loyalty

5.1 Query 5.1: Repeat Purchase Customers

Explanation: Identify customers who placed more than one order, to measure customer loyalty and repeat business.

```

1  SELECT
2      c.customer_id,
3      c.customer_unique_id,
4      g.geolocation_city,
5      g.geolocation_state,
6      COUNT(DISTINCT o.order_id) AS total_orders,
7      MIN(o.order_purchase_timestamp) AS first_order_date,
8      MAX(o.order_purchase_timestamp) AS last_order_date,
9      SUM(oi.price + oi.freight_value) AS total_lifetime_value
10     FROM CUSTOMERS c
11     JOIN GEOLOCATION g ON c.customer_zip_code_prefix = g.zip_code_prefix
12     JOIN ORDERS o ON c.customer_id = o.customer_id
13     JOIN ORDER_ITEMS oi ON o.order_id = oi.order_id
14     GROUP BY c.customer_id, c.customer_unique_id,
15             g.geolocation_city, g.geolocation_state
16     HAVING COUNT(DISTINCT o.order_id) > 1
17     ORDER BY total_orders DESC, total_lifetime_value DESC
18     LIMIT 50;

```

Why Interesting to Analysts: This query identifies the most valuable customer segment for VIP programs, exclusive offers, loyalty rewards, and retention campaigns.
 (returns top 50 repeat customers)

5.2 Query 5.2: Average Time Between Orders (Returning Customers)

Explanation: Calculate the average number of days between consecutive orders for customers who made multiple purchases.

```

1  WITH customer_orders AS (
2      SELECT
3          customer_id,
4          order_id,
5          order_purchase_timestamp,
6          LAG(order_purchase_timestamp) OVER
7              (PARTITION BY customer_id ORDER BY order_purchase_timestamp)
8              AS previous_order_date
9      FROM ORDERS
10 ),
11 order_gaps AS (
12     SELECT
13         customer_id,
14         DATEDIFF(day, previous_order_date, order_purchase_timestamp)
15             AS days_between_orders
16     FROM customer_orders
17     WHERE previous_order_date IS NOT NULL
18 )
19 SELECT
20     COUNT(DISTINCT customer_id) AS returning_customers,
21     AVG(days_between_orders) AS avg_days_between_orders,
22     MIN(days_between_orders) AS min_days_between_orders,
23     MAX(days_between_orders) AS max_days_between_orders,
24     PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY days_between_orders)
25         AS median_days_between_orders
26     FROM order_gaps;

```

Why Interesting to Analysts: Purchase frequency is a retention metric where shorter intervals indicate higher engagement and loyalty, while longer gaps suggest weakening customer relationships.

6 Additional Queries

6.1 Query 6.1: Delivery Performance by State

Explanation: Calculate average delivery time (actual vs. estimated) for each state to identify logistics performance issues.

```

1  SELECT
2      g.geolocation_state,
3      s.state_name,
4      COUNT(DISTINCT o.order_id) AS total_delivered_orders,
5      AVG(DATEDIFF(day, o.order_purchase_timestamp,
6          o.order_delivered_customer_date)) AS avg_actual_delivery_days,
7      AVG(DATEDIFF(day, o.order_purchase_timestamp,
8          o.order_estimated_delivery_date)) AS avg_estimated_delivery_days,
9      AVG(DATEDIFF(day, o.order_estimated_delivery_date,
10         o.order_delivered_customer_date)) AS avg_delay_days,
11     SUM(CASE WHEN o.order_delivered_customer_date >
12             o.order_estimated_delivery_date THEN 1 ELSE 0 END) AS late_deliveries,
13     ROUND(SUM(CASE WHEN o.order_delivered_customer_date >
14             o.order_estimated_delivery_date THEN 1 ELSE 0 END) * 100.0 /
15             COUNT(*), 2) AS late_delivery_percentage
16   FROM ORDERS o
17   JOIN CUSTOMERS c ON o.customer_id = c.customer_id
18   JOIN GEOLOCATION g ON c.customer_zip_code_prefix = g.zip_code_prefix
19   JOIN STATES s ON g.geolocation_state = s.state_code
20   WHERE o.order_delivered_customer_date IS NOT NULL
21   GROUP BY g.geolocation_state, s.state_name
22   ORDER BY avg_delay_days DESC;
```

Why Interesting to Analysts: Delivery performance directly impacts customer satisfaction and review scores. States with high delay percentages ($\geq 20\%$) or long average delays (≥ 3 days) indicate logistics problems requiring infrastructure investment, carrier partnerships, or warehouse positioning.

6.2 Query 6.2: Category Performance by Quarter

Explanation: Track quarterly revenue trends for each product category to identify seasonal patterns and growth trajectories.

```

1  SELECT
2      c.category_name_english,
3      YEAR(o.order_purchase_timestamp) AS order_year,
4      QUARTER(o.order_purchase_timestamp) AS order_quarter,
5      COUNT(DISTINCT oi.order_id) AS orders,
6      SUM(oi.price) AS revenue,
7      AVG(oi.price) AS avg_item_price
8   FROM CATEGORIES c
9   JOIN PRODUCTS p ON c.category_id = p.category_id
10  JOIN ORDER_ITEMS oi ON p.product_id = oi.product_id
11  JOIN ORDERS o ON oi.order_id = o.order_id
12  WHERE o.order_status = 'delivered'
13  GROUP BY c.category_name_english, YEAR(o.order_purchase_timestamp),
14                  QUARTER(o.order_purchase_timestamp)
15  ORDER BY c.category_name_english, order_year, order_quarter;
```

Why Interesting to Analysts: Quarterly trends reveal seasonal demand patterns essential for inventory planning, marketing budget allocation, and seller recruitment timing.

6.3 Query 6.3: Revenue by State and Year

Explanation: Calculates the total revenue and order count for a specific state for a user-specified year. Parameters * ? (State Code) * ? (Year)

```
1 SELECT
2     COUNT(DISTINCT T1.order_id) AS Total_Orders ,
3     SUM(T2.payment_value) AS Total_Revenue
4 FROM ORDERS T1
5 JOIN ORDER_PAYMENTS T2 ON T1.order_id = T2.order_id
6 JOIN CUSTOMERS T3 ON T1.customer_id = T3.customer_id
7 WHERE
8     T3.customer_state = ?
9     AND T1.order_purchase_timestamp LIKE ? || '%'
10 GROUP BY
11     T3.customer_state;
```

Why Interesting to Analysts: This query helps analysts uncover geographic and temporal revenue trends. By allowing users to specify both a state and a year

6.4 Query 6.4: Seller Inventory Check

Explanation: Lists all sellers who offer products whose descriptions contain specific user-defined keywords. Parameters * ? (Product Keyword)

```
1 SELECT DISTINCT
2     T1.seller_id
3 FROM ORDER_ITEMS T1
4 JOIN PRODUCTS T2 ON T1.product_id = T2.product_id
5 WHERE
6     T2.product_description_length IS NOT NULL
7     AND T2.product_description_length LIKE '%' || ? || '%';
```

Why Interesting to Analysts: This query gives analysts visibility into the market supply side, specifically, which sellers are offering products that match certain keywords or product themes