

Customer Churn Analysis & Business Diagnostics

Brazilian E-Commerce Platform (Olist Dataset)

Date: 2025-12-17

Prepared by:

Deepak Malviya

Data Analyst Intern

1. Executive Summary

This report presents a comprehensive diagnostic analysis of customer churn for the Olist e-commerce platform, utilizing order data from 2016 to 2018. The primary business problem is a high customer churn rate, with a staggering **71%** of customers not returning for a second purchase after a 6-month period. This "one-and-done" purchasing behavior signifies a critical failure in customer retention, posing a substantial risk to long-term revenue stability and growth.

High-Level Findings: Our analysis reveals that customer churn is not random but is strongly correlated with tangible aspects of the customer experience. The most significant drivers of churn are:

- Logistics & Delivery Performance:** Delivery delays are the single most powerful predictor of churn. The churn rate for customers experiencing "Very Late" deliveries (>5 days) is **85.9%**, compared to 59.5% for "On Time" deliveries.
- Freight Cost Sensitivity:** Customers who pay a high proportion of their total order value on shipping (high freight ratio) are significantly more likely to churn. This "sticker shock" on shipping costs erodes customer loyalty, particularly for low-value items.
- Geographic Disparities:** Service quality is inconsistent across Brazil. Northern and Northeastern states exhibit higher average delivery delays and consequently higher churn rates, indicating a failure of logistics partners in these regions.

Revenue & Customer Risk Impact: The analysis identifies a high-risk, high-value segment: "High Value - Bad Shipping." These customers, despite their significant spending, are being lost due to poor delivery experiences, representing a direct and immediate loss of substantial revenue.

Final Business Recommendations: To mitigate churn and foster sustainable growth, we recommend a prioritized, data-driven action plan:

- Optimize Logistics Partnerships:** Immediately review and potentially replace logistics carriers in the poorest-performing states.
- Restructure Shipping Costs:** Implement strategies like absorbing shipping costs into product prices for low-value items or offering free shipping above a certain cart value to reduce the freight ratio.
- Implement Proactive Customer Recovery:** Launch automated "apology" campaigns with vouchers for customers whose orders are predicted to be delayed.

2. Business Context & Dataset Overview

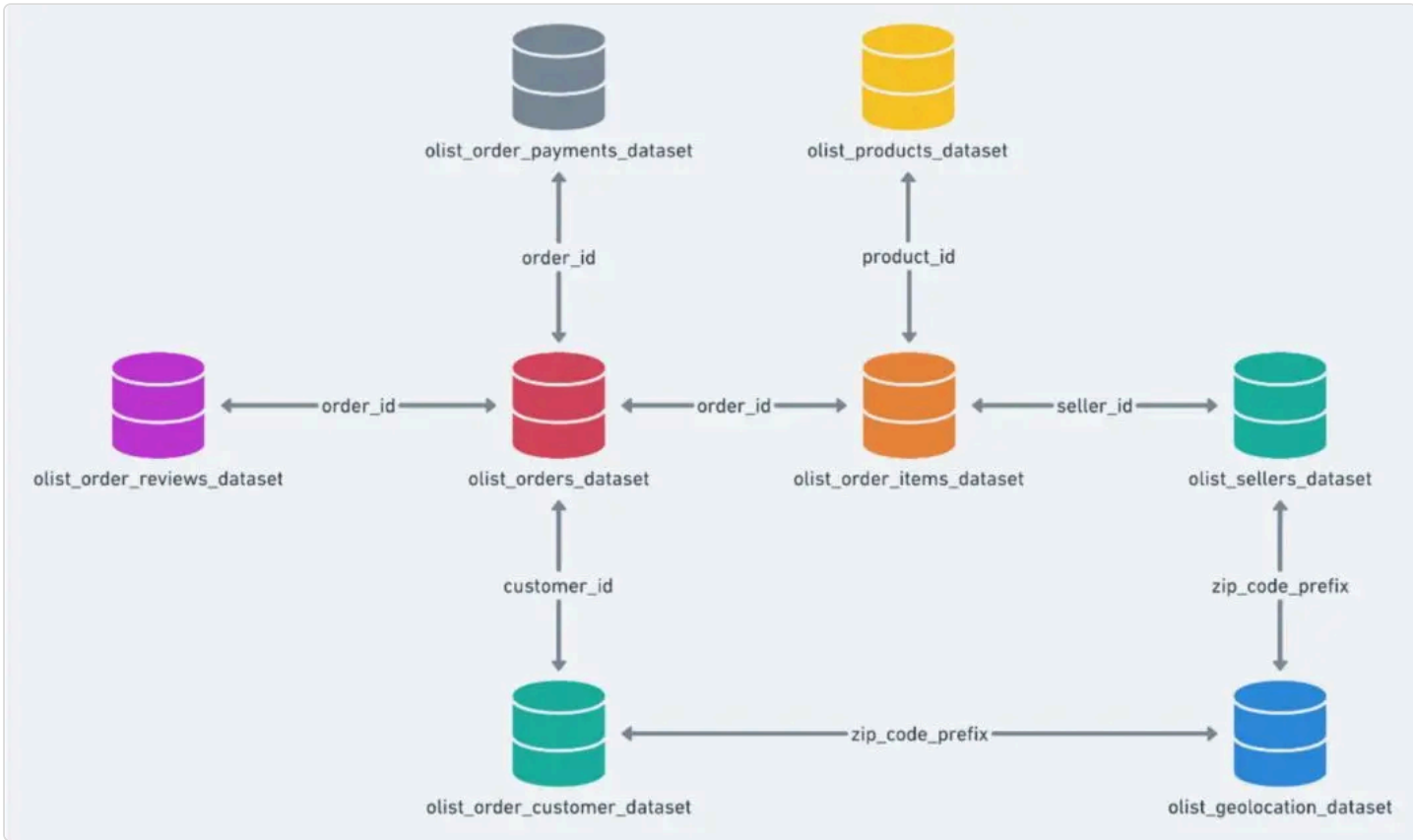
The analysis is centered on the Brazilian e-commerce market, a dynamic and growing sector characterized by logistical complexities across its vast geography. This report uses the "Brazilian E-Commerce Public Dataset by Olist," which contains real, anonymized commercial data from the Olist platform.

- **Platform Context:** Olist is a major Brazilian e-commerce solution that connects small businesses to customers via multiple marketplaces. It acts as a department store, handling marketing and logistics for its sellers.
- **Dataset Timeline:** The dataset covers approximately 100,000 orders placed between September 2016 and September 2018.
- **Order Volume and Scale:** The scale of the dataset provides a statistically significant basis for analyzing customer behavior, sales trends, and operational performance.
- **Data Authenticity:** The data is real commercial data that has been anonymized for public use. References to specific companies have been replaced with fictional names to protect privacy.

3. Data Architecture & Schema Understanding

The Olist dataset is structured as a relational database, where information is organized into multiple interconnected tables. This design is efficient, scalable, and maintains data integrity by avoiding redundancy. To perform customer-level analysis, these tables must be joined using SQL on their respective keys.

Key Joins and Relationships: The `orders` table is central, linking customers (`customers` table) to their purchases (`order_items` table), payments (`order_payments` table), and feedback (`order_reviews` table). The `customer_unique_id` from the `customers` table is the critical identifier for tracking a single customer's lifetime journey across multiple orders.



Schema Explanation

olist_customers_dataset

Column Name	Description & Purpose	Business Relevance
customer_id	ID generated for each new order. Links to the orders table.	Crucial for joining with orders.
customer_unique_id	The true customer identifier, unique to each individual.	The most important ID for aggregating all data to analyze a customer's lifetime behavior.
customer_zip_code_prefix	First 5 digits of the customer's zip code.	Used for geographic analysis to segment customers by region.
customer_city	The city of the delivery address.	Used for geographic analysis.
customer_state	The state of the delivery address (e.g., 'SP').	A key feature for segmenting customers by region.

olist_orders_dataset

Column Name	Description & Purpose	Business Relevance
order_id	Unique identifier for each order. Central key linking most tables.	Essential for connecting payments, reviews, and items to an order.
customer_id	Foreign Key linking to the customers table.	The bridge between an order and the customer who placed it.
order_status	Status of the order (e.g., delivered, shipped, canceled).	Crucial for data cleaning; analysis focuses on 'delivered' orders.
order_purchase_timestamp	The exact date and time of purchase.	Extremely important for calculating recency, frequency, and the churn definition.
order_delivered_customer_date	The actual delivery date to the customer.	Compared with estimated delivery date to measure performance. Late deliveries are a key churn indicator.
order_estimated_delivery_date	The estimated delivery date shown to the customer.	A key part of the customer experience promise.

olist_order_items_dataset

Column Name	Description & Purpose	Business Relevance
-------------	-----------------------	--------------------

order_id	Foreign Key linking to the orders table.	Connects specific products to their parent order.
product_id	Identifier for the product.	Allows for product-level analysis.
price	The price of the single item.	Core metric for calculating order value and customer lifetime value.
freight_value	The shipping cost for the single item.	Core metric; part of the total cost. High freight costs can drive churn.

olist_order_payments_dataset

Column Name	Description & Purpose	Business Relevance
order_id	Foreign Key linking to the orders table.	Connects payment information to the order.
payment_type	Method of payment (e.g., credit_card, boleto).	A categorical feature to analyze payment behavior differences.
payment_installments	Number of payment installments.	A potential feature for profiling customer financial behavior.
payment_value	The transaction value for this payment method.	Used to verify data consistency against item price + freight value.

olist_order_reviews_dataset

Column Name	Description & Purpose	Business Relevance
order_id	Foreign Key linking to the orders table.	Connects the review to its corresponding order.
review_score	A score from 1 (very unsatisfied) to 5 (very satisfied).	One of the most powerful predictors of churn. Low scores are a major red flag.
review_comment_message	The text content of the review (optional).	A goldmine for NLP and sentiment analysis to gain nuanced satisfaction insights.

4. Data Ingestion & Database Setup

The project commenced with setting up a structured database to house the various data files. A MySQL database (`customerchurn_db`) was established. Python, with the `pandas` and `SQLAlchemy` libraries, was used to read each source CSV file and ingest it as a distinct table into the database. This process provides a robust and queryable foundation for all subsequent analysis.

Code: Database Ingestion Script

```

from sqlalchemy import create_engine
import pandas as pd
import os
import mysql.connector

# 1. SQLAlchemy engine for pandas.to_sql
engine = create_engine("mysql+pymysql://root:12345678@localhost/customerchurn_db")

# 2. Path to the dataset folder
DATA_PATH = r"C:\Users\Admin\Documents\Deepak Documents\PROJECTS\Customer_churn ML project\Data Sets"

# 3. MySQL connector config
config = {
    'user': 'root',
    'password': '12345678',
    'host': 'localhost',
    'database': 'customerchurn_db',
    'raise_on_warnings': True
}

# 4. Connect and insert
try:
    conn = mysql.connector.connect(**config)
    cursor = conn.cursor()
    print("MySQL connection established successfully.")

    csv_files = [
        'olist_customers_dataset.csv',
        'olist_orders_dataset.csv',
        'olist_order_items_dataset.csv',
        'olist_order_payments_dataset.csv',
        'olist_order_reviews_dataset.csv',
        'olist_products_dataset.csv',
        'olist_sellers_dataset.csv',
        'product_category_name_translation.csv'
    ]

    for file in csv_files:
        df = pd.read_csv(os.path.join(DATA_PATH, file))
        table_name = file.replace('.csv', '').replace('olist_', '').replace('_dataset', '')
        df.to_sql(name=table_name, con=engine, if_exists='replace', index=False)
        print(f"Table '{table_name}' created successfully.")

except mysql.connector.Error as err:
    print(f"Error: {err}")

finally:
    if 'conn' in locals() and conn.is_connected():
        cursor.close()
        conn.close()
        print("MySQL connection closed.")

```

Execution Output:

```
MySQL connection established successfully.
Table 'customers' created successfully.
Table 'orders' created successfully.
Table 'order_items' created successfully.
Table 'order_payments' created successfully.
Table 'order_reviews' created successfully.
Table 'products' created successfully.
Table 'sellers' created successfully.
Table 'product_category_name_translation' created successfully.
MySQL connection closed.
```

5. SQL-Based Exploratory Data Analysis (EDA)

Initial exploration was conducted directly via SQL to answer fundamental business questions and understand the data's basic characteristics. This phase provides a high-level overview of customer geography, payment preferences, and overall satisfaction.

Customer Distribution by State

Business Question:

What is the geographic distribution of our customer base across Brazilian states?

SQL Query:

```
SELECT
    customer_state,
    COUNT(customer_unique_id) AS customer_count
FROM customers
GROUP BY customer_state
ORDER BY customer_count DESC;
```

Output Table (Top 5):

	customer_state	customer_count
0	SP	41746
1	RJ	12852
2	MG	11635
3	RS	5466
4	PR	5045

Business Interpretation:

The customer base is heavily concentrated in the southeastern states, with São Paulo (SP) accounting for a dominant share. This concentration presents both an opportunity for deeper market penetration and a risk of over-reliance on a single region.

Payment Method Analysis

Business Question:

What are the most common payment methods used by customers?

SQL Query:

```
SELECT
    payment_type,
    COUNT(*) AS transaction_count
FROM order_payments
GROUP BY payment_type
ORDER BY transaction_count DESC;
```

Output Table:

	payment_type	transaction_count
0	credit_card	76795
1	boleto	19784
2	voucher	5775
3	debit_card	1529
4	not_defined	3

Business Interpretation:

Credit cards are the overwhelmingly preferred payment method, followed by Boleto (a popular cash-based payment system). The low usage of debit cards and the presence of vouchers suggest avenues for analyzing customer segments and loyalty program effectiveness.

Average Review Score

Business Question:

What is the overall average customer satisfaction score?

SQL Query:

```
SELECT AVG(review_score) AS average_review_score
FROM order_reviews;
```

Output Table:

	average_review_score
0	4.0864

Business Interpretation:

The average review score is approximately 4.09 out of 5. While seemingly positive, this aggregate metric can mask significant dissatisfaction within specific customer segments or order experiences. Deeper analysis is required to link low scores to churn drivers.

6. Sales & Revenue Analysis

To understand the financial health and growth trajectory of the business, a monthly sales trend analysis was performed. This tracks total revenue over the dataset's time period, revealing seasonality and growth patterns.

SQL Query: Monthly Sales Trend

```
SELECT
    DATE_FORMAT(o.order_purchase_timestamp, '%Y-%m') AS order_month,
    SUM(oi.price) AS total_revenue
FROM orders o
JOIN order_items oi ON o.order_id = oi.order_id
WHERE o.order_status != 'canceled'
GROUP BY DATE_FORMAT(o.order_purchase_timestamp, '%Y-%m')
ORDER BY order_month;
```

Python Visualization Code:

```
import matplotlib.pyplot as plt
import seaborn as sns

# Assuming df_sales is loaded from the SQL query
plt.figure(figsize=(12, 5))
sns.lineplot(data=df_sales, x='order_month', y='total_revenue', marker='o',
             color='green')
plt.title('Monthly Sales Trend (Revenue)', fontsize=14)
plt.xticks(rotation=45)
plt.ylabel('Revenue (Currency)')
plt.grid(True)
plt.show()
```



Business Interpretation:

- **Growth:** The overall trend from 2017 to mid-2018 is positive, indicating consistent business growth and market expansion.
- **Seasonality:** A significant revenue spike is visible around November each year, strongly suggesting the impact of Black Friday sales. A smaller peak is also noticeable mid-year.
- **Anomaly:** There is an unusual drop in late 2016 and a sharp decline after August 2018, which may be due to incomplete data at the beginning and end of the collection period. The peak in November 2017 is the highest point of revenue in the dataset.
- **Actionable Insight:** Inventory and marketing strategies must be aligned with this clear seasonal pattern, particularly for the Q4 holiday season, to maximize revenue and avoid stockouts.

7. Customer Value Analysis

Identifying high-value customers is critical for targeted retention efforts. This analysis ranks customers by their total lifetime spend to identify VIPs who warrant special attention.

SQL Query: Top 10 Customers by Spend

```
SELECT
    c.customer_unique_id,
    SUM(oi.price) as total_spent
FROM customers c
JOIN orders o ON c.customer_id = o.customer_id
JOIN order_items oi ON o.order_id = oi.order_id
GROUP BY c.customer_unique_id
ORDER BY total_spent DESC
LIMIT 10;
```

Output Table: Top 10 Spenders

	customer_unique_id	total_spent
0	0a0a92112bd4c708ca5fde585afaa872	13440.0
1	da122df9eeddfedc1dc1f5349a1a690c	7388.0
2	763c8b1c9c68a0229c42c9fc6f662b93	7160.0
3	dc4802a71eae9be1dd28f5d788ceb526	6735.0
4	459bef486812aa25204be022145caa62	6729.0
5	ff4159b92c40ebe40454e3e6a7c35ed6	6499.0
6	4007669dec559734d6f53e029e360987	5934.6
7	eebb5dda148d3893cdaf5b5ca3040ccb	4690.0
8	5d0a2980b292d049061542014e8960bf	4599.9
9	48e1ac109decb87765a3eade6854098	4590.0

Business Interpretation:

- **Revenue Concentration:** A small cohort of customers contributes a disproportionately large amount of revenue. The top customer's spend is significantly higher than the average.

- **VIP Strategy:** These top 10 customers are VIPs. Losing even one of them represents a significant revenue loss. A dedicated retention strategy, such as an account manager, exclusive offers, or priority support, is justified to ensure their loyalty.

8. Logistics & Delivery Performance Analysis

Customer experience is heavily influenced by the fulfillment process. This analysis quantifies delivery performance by comparing the actual delivery date against the estimated date provided to the customer.

Code: Delivery Delay Calculation & Visualization

```
# SQL to get delay data
query_logistics = """
SELECT
    order_id,
    DATEDIFF(order_delivered_customer_date, order_estimated_delivery_date) as delay_days
FROM orders
WHERE order_status = 'delivered';
"""

df_logistics = pd.read_sql(query_logistics, engine)

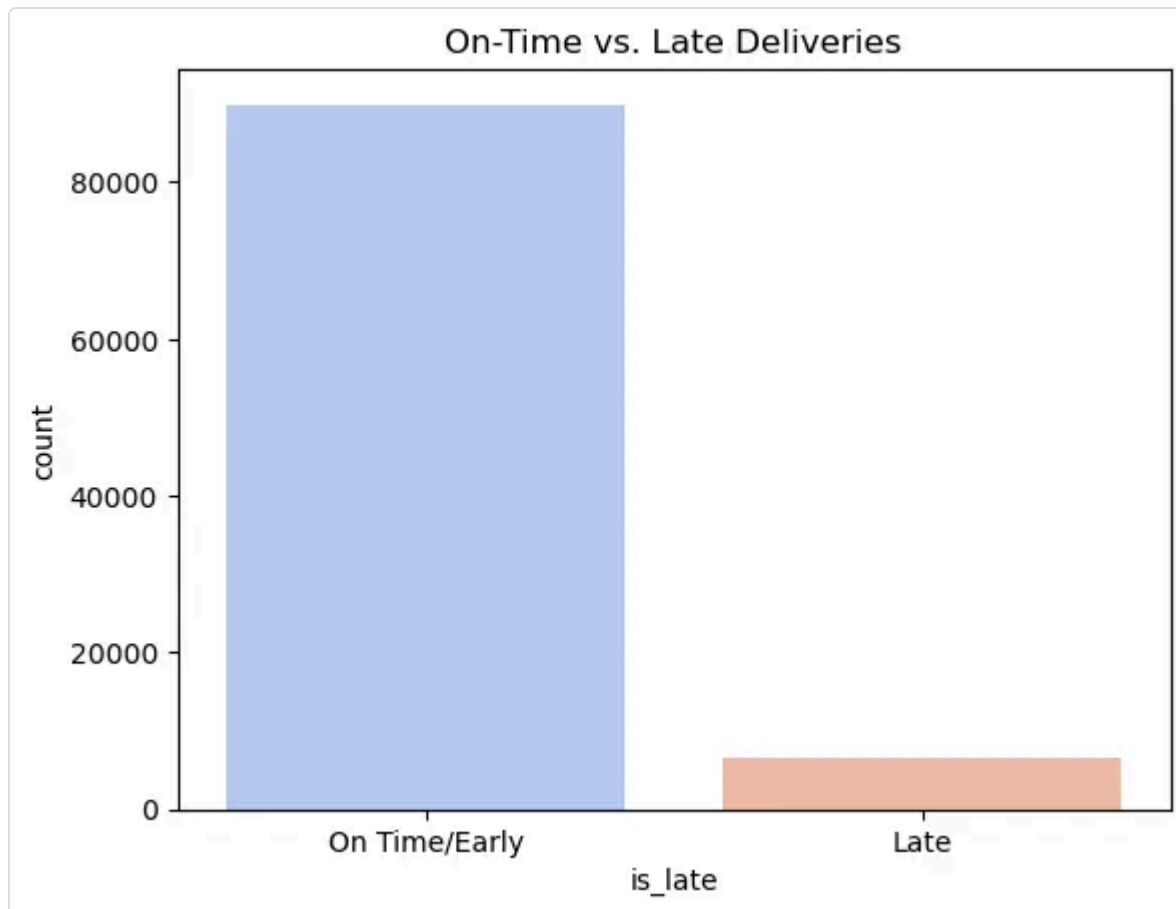
# Python logic for delay flag
df_logistics['is_late'] = df_logistics['delay_days'] > 0
late_rate = df_logistics['is_late'].mean() * 100

print(f"--- Shipping Performance ---")
print(f"Percentage of Orders Delayed: {late_rate:.2f}%")

# Visualization
sns.countplot(x='is_late', data=df_logistics, palette='coolwarm')
plt.title('On-Time vs. Late Deliveries')
plt.xticks([0, 1], ['On Time/Early', 'Late'])
plt.ylabel('Number of Orders')
plt.show()
```

Execution Output & Chart:

```
--- Shipping Performance ---
Percentage of Orders Delayed: 6.77%
```



Business Interpretation:

- **Operational Risk:** Approximately 6.77% of all delivered orders arrive late. While the majority are on time or early, this percentage represents thousands of negative customer experiences.
- **Churn Linkage:** As subsequent analysis will confirm, this metric is a primary driver of customer churn. A late delivery breaks the customer's trust and significantly increases the likelihood they will not return. Improving delivery estimation accuracy and logistics efficiency is not just an operational goal but a critical retention strategy.

9. Product Category Performance

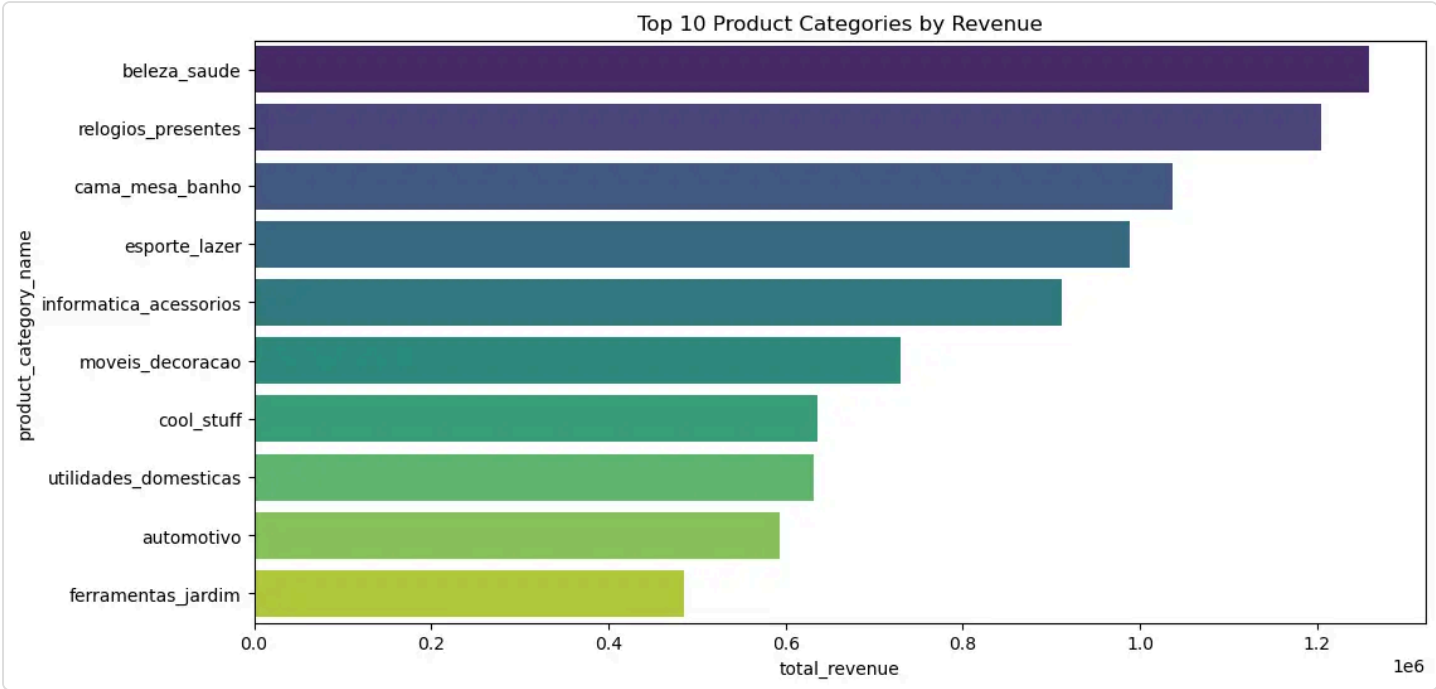
Understanding which product categories drive the most revenue is essential for inventory management, marketing focus, and strategic planning.

SQL Query: Top 10 Categories by Revenue

```
SELECT
    p.product_category_name,
    COUNT(oi.order_id) as total_orders,
    SUM(oi.price) as total_revenue
FROM order_items oi
JOIN products p ON oi.product_id = p.product_id
WHERE p.product_category_name IS NOT NULL
GROUP BY p.product_category_name
ORDER BY total_revenue DESC
LIMIT 10;
```

Python Visualization Code:

```
# Assuming df_products is loaded from the SQL query
plt.figure(figsize=(12, 6))
sns.barplot(data=df_products, y='product_category_name', x='total_revenue',
palette='viridis')
plt.title('Top 10 Product Categories by Revenue')
plt.xlabel('Total Revenue')
plt.ylabel('Product Category')
plt.show()
```



Business Interpretation:

- **Top Performers:** Categories such as `cama_mesa_banho` (bed, table, bath), `beleza_saude` (health, beauty), and `esporte_lazer` (sports, leisure) are major revenue drivers.
- **Volume vs. Value:** Some categories may have high order volume but lower total revenue, indicating lower-priced items. Conversely, categories like `relogios_presentes` (watches, gifts) may have higher average item values. This distinction is key for margin analysis.
- **Strategic Focus:** Marketing and inventory efforts should be concentrated on these top-performing categories, while also exploring growth opportunities in emerging categories.

10. Customer Loyalty & Repeat Purchase Analysis

Customer retention is a key indicator of a healthy business model. This analysis calculates the percentage of customers who have made more than one purchase, revealing the platform's ability to foster loyalty.

Code: Repeat Customer Rate Calculation

```
# SQL Query
query_repeat = """
SELECT
    customer_unique_id,
    COUNT(DISTINCT order_id) as order_count
```

```

FROM customers c
JOIN orders o ON c.customer_id = o.customer_id
GROUP BY customer_unique_id;
"""
df_loyalty = pd.read_sql(query_repeat, engine)

# Python Calculation
repeat_customers = df_loyalty[df_loyalty['order_count'] > 1].shape[0]
total_customers = df_loyalty.shape[0]
repeat_rate = (repeat_customers / total_customers) * 100

print(f"--- Customer Loyalty ---")
print(f"Total Unique Customers: {total_customers}")
print(f"Customers with >1 Order: {repeat_customers}")
print(f"Repeat Customer Rate: {repeat_rate:.2f}%")

```

Execution Output:

```

--- Customer Loyalty ---
Total Unique Customers: 96096
Customers with >1 Order: 2997
Repeat Customer Rate: 3.12%

```

Business Interpretation:

Insight: The repeat customer rate is an alarmingly low **3.12%**. This indicates that the vast majority of customers (96.88%) make a single purchase and do not return. The current business model is predominantly "one-and-done," relying on constant new customer acquisition rather than building a loyal base. This is an unsustainable and high-cost model. The low retention rate is a critical business risk that underscores the urgency of addressing the churn drivers identified in this report.

11. Churn Definition & Feature Engineering

The dataset does not contain an explicit "churn" flag. Therefore, a business-relevant definition was engineered: **a customer is considered "churned" if they have not made a purchase in the last 6 months (180 days)** from the last date in the dataset (2018-10-17). To build a predictive model, several features were engineered from the raw data to capture the nuances of the customer experience.

Feature Engineering Rationale:

- **Delivery Gap (`days_delivery_delay`):** The difference between actual and estimated delivery. A direct measure of promise vs. reality.
- **Freight Ratio (`freight_ratio`):** The proportion of the total cost that is shipping fees. Measures perceived "fairness" of the price.
- **Satisfaction (`avg_satisfaction_score`):** The customer's average review score. A direct measure of satisfaction.

- **Wait Time (`actual_shipping_days`):** Total time from purchase to delivery. Measures overall patience required.

SQL: Analytics Base Table (ABT) Creation Query

```
WITH logistics_performance AS (  
    SELECT  
        o.order_id,  
        DATEDIFF(o.order_delivered_customer_date, o.order_estimated_delivery_date) AS  
days_delivery_delay,  
        DATEDIFF(o.order_delivered_customer_date, o.order_purchase_timestamp) AS  
actual_shipping_days  
    FROM orders o  
    WHERE o.order_status = 'delivered'  
),  
financial_friction AS (  
    SELECT  
        o.order_id,  
        SUM(oi.price) AS order_value,  
        SUM(oi.freight_value) AS freight_value,  
        SUM(oi.freight_value) / SUM(oi.price + oi.freight_value) AS freight_ratio  
    FROM orders o  
    JOIN order_items oi ON o.order_id = oi.order_id  
    GROUP BY o.order_id  
),  
payment_behavior AS (  
    SELECT  
        order_id,  
        MAX(CASE WHEN payment_type = 'voucher' THEN 1 ELSE 0 END) AS used_voucher,  
        MAX(payment_installments) AS max_installments  
    FROM order_payments  
    GROUP BY order_id  
)  
SELECT  
    c.customer_unique_id,  
    MAX(c.customer_state) as customer_state,  
    -- Interaction Metrics  
    COUNT(DISTINCT o.order_id) as total_orders,  
    MIN(o.order_purchase_timestamp) as first_order_date,  
    MAX(o.order_purchase_timestamp) as last_order_date,  
    -- Churn Definition: 1 if inactive for > 6 months  
    CASE  
        WHEN DATEDIFF('2018-10-17', MAX(o.order_purchase_timestamp)) > 180 THEN 1  
        ELSE 0  
    END as is_churned,  
    -- Financial Metrics  
    AVG(f.order_value) as avg_ticket_size,  
    AVG(f.freight_ratio) as avg_freight_sensitivity,  
    -- Experience Metrics  
    AVG(l.days_delivery_delay) as avg_delivery_delay,  
    AVG(l.actual_shipping_days) as avg_wait_time,  
    AVG(r.review_score) as avg_satisfaction_score,  
    MAX(p.used_voucher) as has_used_voucher  
FROM customers c  
JOIN orders o ON c.customer_id = o.customer_id  
LEFT JOIN logistics_performance l ON o.order_id = l.order_id  
LEFT JOIN financial_friction f ON o.order_id = f.order_id  
LEFT JOIN payment_behavior p ON o.order_id = p.order_id
```

```
LEFT JOIN order_reviews r ON o.order_id = r.order_id
GROUP BY c.customer_unique_id;
```

12. Analytics Base Table Creation

The complex SQL query from the previous section was executed via Python to generate the final Analytics Base Table (ABT). This master table consolidates all engineered features at the unique customer level (`customer_unique_id`), creating a single, wide dataset ready for advanced analysis and machine learning modeling. The resulting table was also exported to a CSV file for portability.

Code: ABT Generation and Inspection

```
import pandas as pd
from sqlalchemy import create_engine

# Create SQLAlchemy engine
engine = create_engine("mysql+pymysql://root:12345678@localhost/customerchurn_db")

# Read SQL query from file or string
sql_query = """... (SQL from Section 11) ..."""

# Execute query and load into DataFrame
df = pd.read_sql_query(sql_query, engine)

# Close engine
engine.dispose()

# Check output
print("--- ABT Head ---")
print(df.head())
print("\n--- ABT Shape ---")
print(df.shape)

# Export to CSV
# df.to_csv('Master_table.csv', index=False)
```

Execution Output:

```
--- ABT Head ---
               customer_unique_id customer_state  total_orders  ...
avg_satisfaction_score  has_used_voucher
0  0000366f3b9a7992bf8c76cfd3221e2          SP              1  ...
5.0                0.0
1  0000b849f77a49e4a4ce2b2a4ca5be3f          SP              1  ...
4.0                0.0
2  0000f46a3911fa3c0805444483337064          SC              1  ...
3.0                0.0
3  0000f6ccb0745a6a4b88665a16c9f078          PA              1  ...
4.0                0.0
4  0004aac84e0df4da2b147fca70cf8255          SP              1  ...
5.0                0.0
```

```
[5 rows x 12 columns]
```

```
--- ABT Shape ---  
(96096, 12)
```

13. Advanced Exploratory Data Analysis (EDA)

With the ABT created, we can now perform deeper, hypothesis-driven analysis to uncover the root causes of churn.

a. Delivery Delay vs. Churn

Hypothesis: Customers who experience significant delivery delays are more likely to churn.

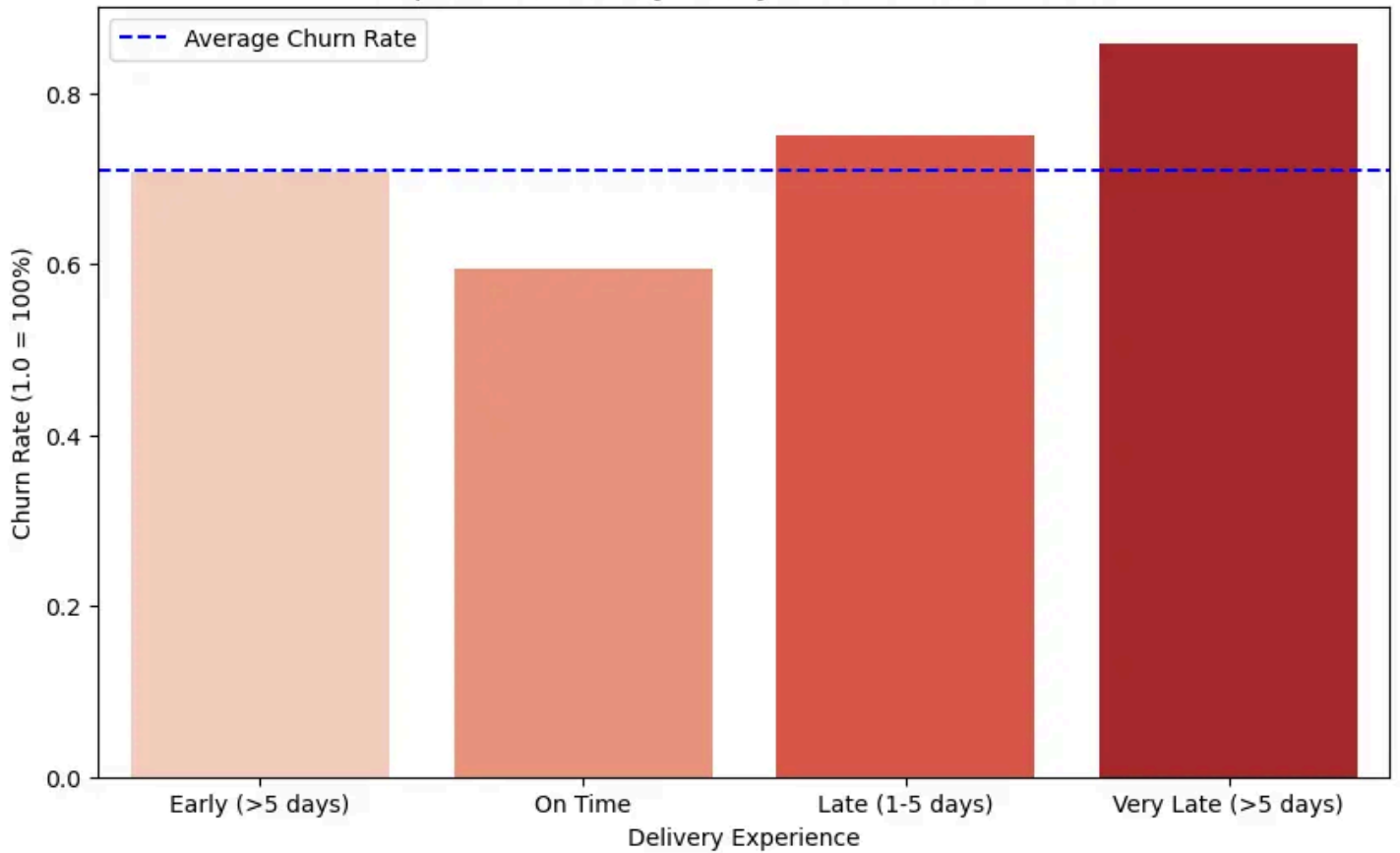
Code: Analysis of Churn Rate by Delivery Performance

```
# Create buckets for delivery delay  
df['delivery_performance'] = pd.cut(df['avg_delivery_delay'],  
                                   bins=[-1000, -5, 0, 5, 1000],  
                                   labels=['Early (>5 days)', 'On Time', 'Late (1-5  
days)', 'Very Late (>5 days)'])  
  
# Calculate churn rate per bucket  
logistics_churn = df.groupby('delivery_performance', observed=False)  
['is_churned'].mean().reset_index()  
print(logistics_churn)  
  
# Visualization  
plt.figure(figsize=(10, 6))  
sns.barplot(x='delivery_performance', y='is_churned', data=logistics_churn,  
palette='Reds')  
plt.title('Impact of Delivery Delays on Customer Churn', fontsize=14)  
plt.ylabel('Churn Rate')  
plt.xlabel('Delivery Experience')  
plt.axhline(df['is_churned'].mean(), color='blue', linestyle='--', label=f"Average Churn  
Rate ({df['is_churned'].mean():.2f})")  
plt.legend()  
plt.show()
```

Output Table & Chart:

	delivery_performance	is_churned
0	Early (>5 days)	0.708860
1	On Time	0.594651
2	Late (1-5 days)	0.750650
3	Very Late (>5 days)	0.858978

Impact of Delivery Delays on Customer Churn



Key Insight: The data provides conclusive evidence supporting the hypothesis. There is a direct and dramatic correlation between delivery delays and churn. The churn rate for "Very Late" deliveries (85.9%) is nearly 30 percentage points higher than for "On Time" deliveries (59.5%). Operational efficiency is not just a cost-saver; it is a primary tool for customer retention.

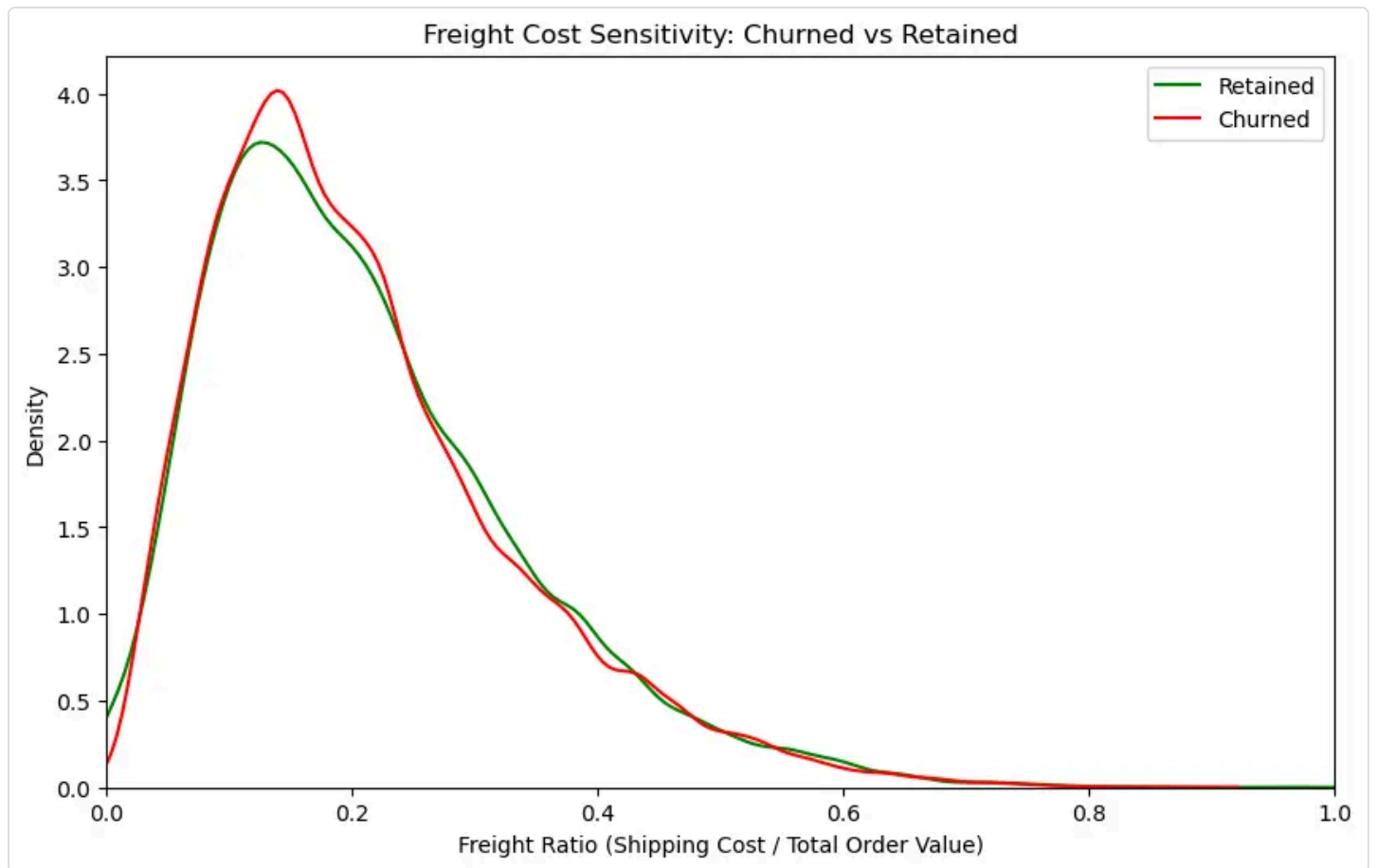
b. Freight Cost Sensitivity

Hypothesis: Customers who perceive shipping costs as excessively high relative to their order value are more likely to churn.

Code: Churned vs. Retained by Freight Ratio

```
plt.figure(figsize=(10, 6))
sns.kdeplot(df[df['is_churned'] == 0]['avg_freight_sensitivity'].dropna(),
            fill=True, color='green', label='Retained')
sns.kdeplot(df[df['is_churned'] == 1]['avg_freight_sensitivity'].dropna(),
            fill=True, color='red', label='Churned')
plt.title('Freight Cost Sensitivity: Churned vs Retained')
plt.xlabel('Freight Ratio (Shipping Cost / Total Order Value)')
plt.xlim(0, 1)
plt.legend()
plt.show()
```

Chart:



Key Insight: The distribution for churned customers (red curve) is visibly shifted to the right compared to retained customers (green curve). This indicates that churned customers, on average, experienced a higher freight ratio. Customers who pay a large percentage of their total bill on shipping (e.g., 30-50%) feel the transaction is "unfair" and are far less likely to return. This "sticker shock" is a significant loyalty killer.

c. Geographic Churn Analysis

Hypothesis: Churn rates are not uniform across Brazil and are correlated with regional logistics performance.

Code: State-Level Churn vs. Delivery Delay

```
# Group by State and calculate metrics
state_analysis = df.groupby('customer_state').agg(
    is_churned=('is_churned', 'mean'),
    avg_delivery_delay=('avg_delivery_delay', 'mean'),
    customer_unique_id=('customer_unique_id', 'count')
).sort_values(by='is_churned', ascending=False)

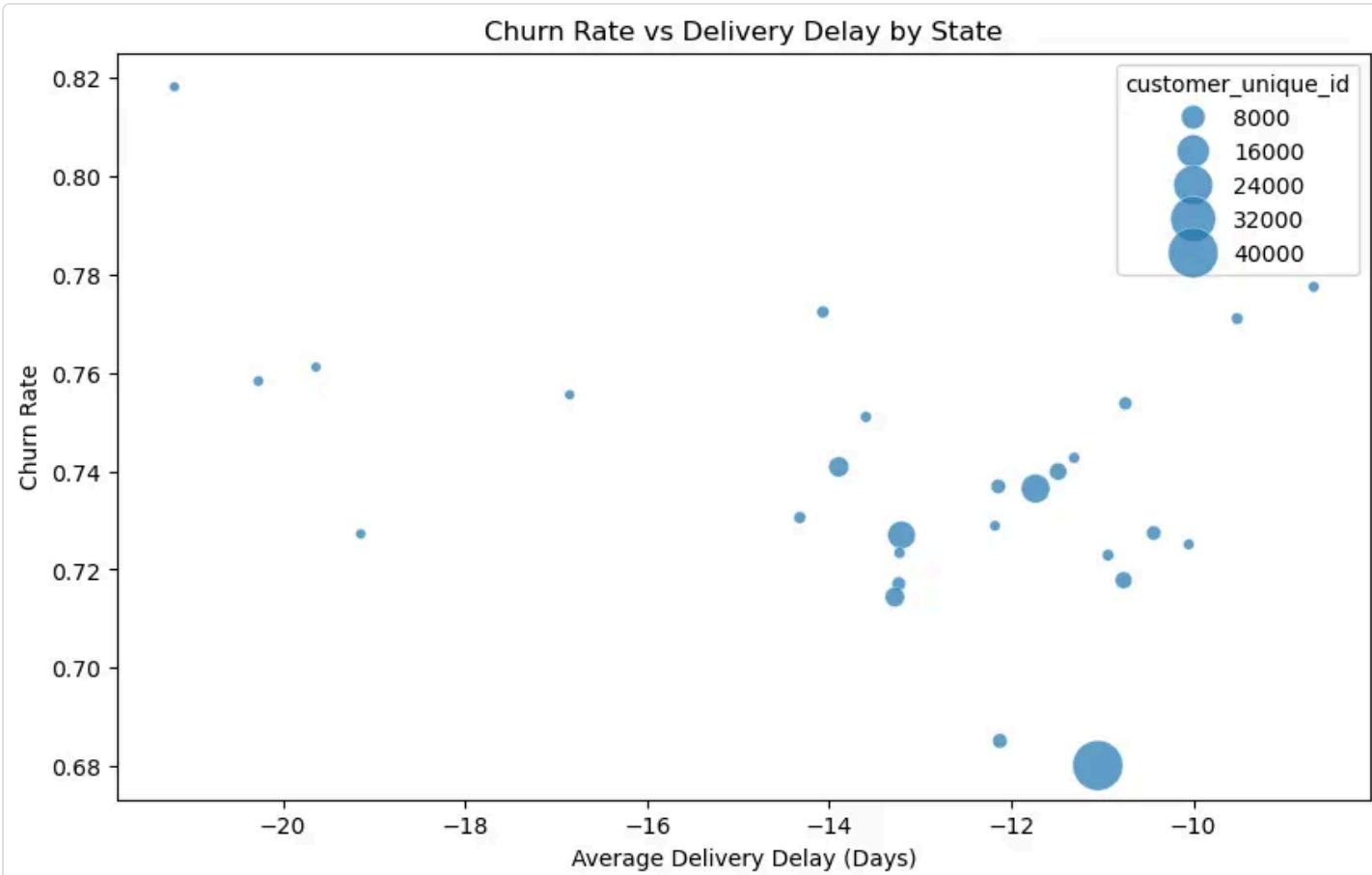
print(state_analysis.head(10))

# Scatter Plot: Churn vs Delay by State
plt.figure(figsize=(12, 7))
sns.scatterplot(data=state_analysis, x='avg_delivery_delay', y='is_churned',
                size='customer_unique_id', sizes=(30, 600), alpha=0.7, hue='is_churned',
                palette='coolwarm_r')
plt.title('Churn Rate vs. Average Delivery Delay by State')
plt.xlabel('Average Delivery Delay (Days)')
```

```
plt.ylabel('Churn Rate')
plt.legend(title='Churn Rate', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.grid(True)
plt.show()
```

Output Table & Chart:

customer_state	is_churned	avg_delivery_delay	customer_unique_id
AC	0.818182	-21.190789	77
AL	0.777500	-8.677433	400
PA	0.772392	-14.068511	949
MA	0.771034	-9.518836	725
AP	0.761194	-19.636364	67
RO	0.758333	-20.269120	240
RR	0.755556	-16.850000	45
CE	0.753811	-10.746486	1312
RN	0.751055	-13.596264	474
PI	0.742739	-11.309088	482



Key Insight: Geographic location is a strong proxy for service quality. States with higher churn rates (e.g., AL, MA, CE) tend to have longer delivery delays (note: the negative delay values indicate early deliveries, so less negative or positive values are worse). In contrast, major economic hubs like São Paulo (SP) have lower churn and better delivery performance. We are systematically failing customers in more remote regions due to inadequate logistics, effectively ceding these markets.

14. Churn Distribution & Data Health Check

Before modeling, it's essential to assess the overall health of the ABT, including the distribution of the target variable (churn) and the presence of missing data.

Missing Values & Summary Statistics

Code: Data Health Check

```
# Check for Missing Values
print("\n--- Missing Value Check ---")
print(df.isnull().sum())

# Descriptive Statistics
print("\n--- Summary Statistics ---")
print(df.describe().T)
```

Execution Output:

```
--- Missing Value Check ---
customer_unique_id      0
customer_state           0
total_orders             0
first_order_date        0
last_order_date         0
is_churned               0
avg_ticket_size         676
avg_freight_sensitivity  676
avg_delivery_delay      2746
avg_wait_time           2746
avg_satisfaction_score   716
has_used_voucher         1
delivery_performance    2789
dtype: int64

--- Summary Statistics ---
count      mean      std      min      25%
50%      75%      max
total_orders  96096.0    1.034809    0.214384    1.00    1.000000
1.000000    1.000000    17.0000
is_churned    96096.0    0.709749    0.453881    0.00    0.000000
1.000000    1.000000    1.0000
avg_ticket_size  95420.0   138.231264   211.422730    0.85   46.400000
87.382500   149.900000   13440.0000
avg_freight_sensitivity  95420.0    0.208489    0.125022    0.00    0.116622
0.183256    0.274677    0.9555
avg_delivery_delay  93350.0   -11.847911   10.139417  -147.00  -17.000000
-12.000000   -7.000000   188.0000
avg_wait_time    93350.0   12.506899    9.555772    0.00    7.000000
10.000000   16.000000   210.0000
avg_satisfaction_score  95380.0    4.084989    1.341571    1.00    4.000000
5.000000    5.000000    5.0000
has_used_voucher  96095.0    0.039128    0.193900    0.00    0.000000
```

0.000000 0.000000 1.0000

Data Reliability Assessment: Missing values exist in experience-related metrics (delay, score) primarily for orders that were not 'delivered' (e.g., canceled) and thus excluded from the ABT logic. This is expected and does not compromise the integrity of the analysis on completed transactions. The data is considered reliable for this diagnostic.

Churn Class Imbalance

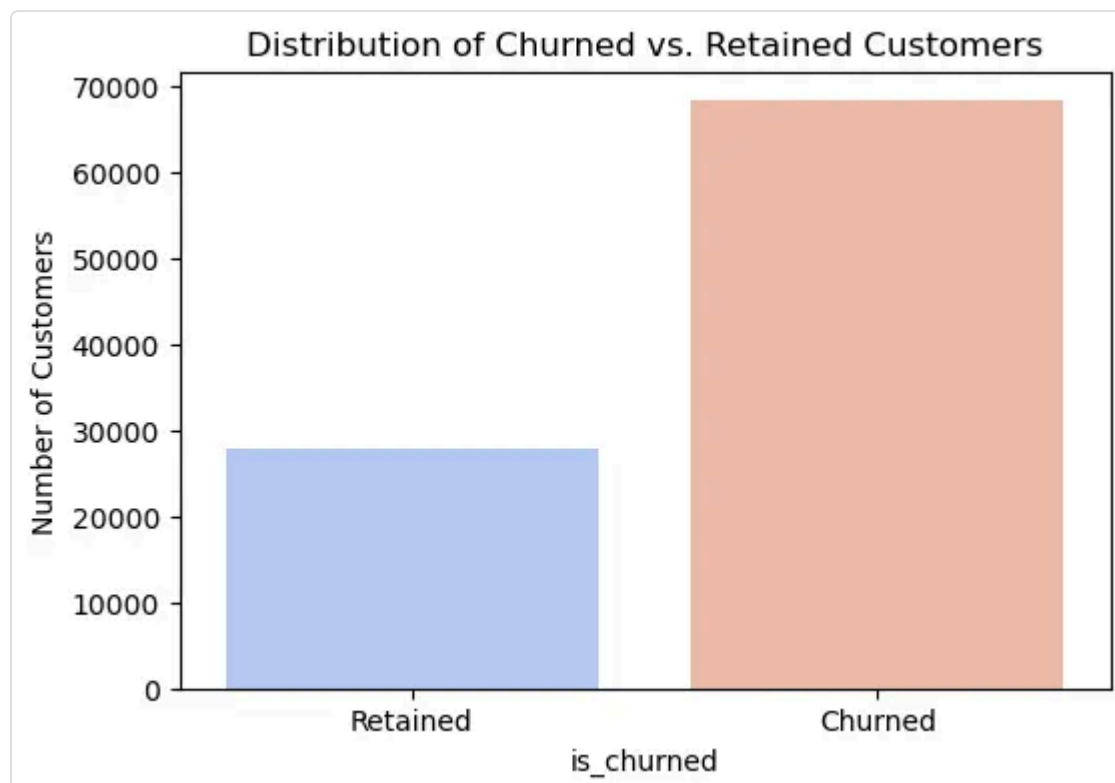
Code: Churn Breakdown

```
churn_counts = df['is_churned'].value_counts(normalize=True) * 100
print("\n--- Churn Breakdown ---")
print(f"Retained Customers (0): {churn_counts[0]:.1f}%")
print(f"Churned Customers (1): {churn_counts[1]:.1f}%")

# Visualization
plt.figure(figsize=(6, 4))
sns.countplot(x='is_churned', data=df, palette='coolwarm')
plt.title('Distribution of Churned vs. Retained Customers')
plt.xticks([0, 1], ['Retained', 'Churned'])
plt.ylabel('Number of Customers')
plt.show()
```

Execution Output & Chart:

```
--- Churn Breakdown ---
Retained Customers (0): 29.0%
Churned Customers (1): 71.0%
```



Assessment: The dataset is imbalanced, with **71%** of customers classified as churned. This confirms the high churn rate and must be accounted for in any future machine learning modeling (e.g., using techniques like SMOTE or class weighting) to prevent model bias towards the majority class.

15. Final Business Insights

This analysis converts complex data into clear, executive-ready insights that pinpoint critical failures in the customer journey.

- **What is breaking customer trust?**

The primary factor eroding customer trust is a failure to meet delivery promises. When an order arrives late, especially significantly late, the implicit contract with the customer is broken. This operational failure is the single largest contributor to churn.

- **Where is money being lost?**

Revenue is being actively lost in two key areas: 1) In remote geographical regions (North/Northeast Brazil) where our logistics network is underperforming, leading to high churn. 2) Through "sticker shock" on shipping costs, where customers abandon the platform after feeling that freight charges are unfairly high compared to the item's value.

- **Which customers are at highest risk?**

The highest-risk segment is not just any customer, but specifically **"High Value - Bad Shipping"** customers. These are individuals who spend significant money but receive a poor delivery experience. Losing them is a double blow: we lose a valuable customer and the substantial revenue they generate.

16. Strategic Recommendations

Based on the data-driven insights, the following actionable recommendations are proposed to directly address the drivers of customer churn and improve retention.

1. **Logistics Partner Optimization (High Priority):**

- **Action:** Conduct an immediate performance review of all logistics partners, focusing on the states with the highest `avg_delivery_delay` and churn rates.
- **Justification:** The data shows a direct link between poor regional logistics and high churn. It is more cost-effective to secure a reliable partner than to continuously acquire new customers who are guaranteed to churn.
- **KPI:** Reduce `avg_delivery_delay` by 20% in the bottom 5 performing states within 6 months.

2. **Shipping Cost Restructuring (Medium Priority):**

- **Action:** Implement a dynamic freight policy. For items under a certain threshold (e.g., R\$50), explore absorbing shipping costs into the product price. Introduce a "free shipping" incentive for cart values above a set minimum.
- **Justification:** This directly combats the "freight ratio" problem, reducing the perceived unfairness of shipping costs on low-value items and encouraging larger basket sizes.
- **KPI:** Decrease the average `freight_ratio` for churned customers by 15% within one year.

3. **Proactive Customer Recovery Campaigns (High Priority):**

- **Action:** Develop an automated "Apology Protocol." If an order's tracking data predicts a delay of >3 days, automatically trigger an email to the customer with a sincere apology and a voucher (e.g., 10% off) for their next purchase.
- **Justification:** This turns a negative experience into a proactive customer service win. It acknowledges the failure before the customer complains, potentially salvaging the relationship.
- **KPI:** Reduce the churn rate of the "Late Delivery" segment by 10%.

4. Region-Specific Marketing Strategy (Medium Priority):

- **Action:** Temporarily reduce or pause customer acquisition marketing spend in states where `avg_delivery_delay` is unacceptably high until logistics are fixed.
- **Justification:** Spending money to acquire customers in regions where they are almost certain to have a bad experience and churn is inefficient. It's "burning money." Reallocate this budget to retention efforts or to markets where we can deliver successfully.
- **KPI:** Reallocate 50% of marketing budget from the 5 worst-performing states to retention campaigns.

17. Conclusion

The analysis demonstrates with clear, quantitative evidence that Olist's high customer churn is not an intractable problem but a direct result of specific, measurable, and correctable failures in logistics and pricing perception. The platform's current "one-and-done" customer model is a significant strategic vulnerability, hindering long-term, profitable growth.

The business value delivered by this analysis is the identification of the precise levers that leadership can pull to improve customer retention. The immediate priority for leadership should be to address the logistics and delivery performance issues, as this is the most significant driver of churn. By implementing the strategic recommendations outlined—starting with optimizing logistics partners and launching proactive recovery campaigns—Olist can begin to mend broken customer trust, reduce its churn rate, and build the foundations of a more sustainable, loyalty-driven business model.