

profitLens

ProfitLens: Visualizing E-Commerce Trends

Understanding which products, regions, and discount strategies drive the highest profitability through data-driven insights.



The Business Challenge

Our Mission

Analyze e-commerce data from Kaggle to uncover how products, regions, and discounts impact sales and profit. Discover which business areas generate the most revenue and enable smarter, data-driven decisions.

Why It Matters

Companies need clear insights into profitability drivers to optimize strategies, reduce waste, and maximize returns across categories and regions.

Project Objectives



Data Preparation

Clean and organize raw data



Pattern Discovery

Identify sales, profit, and discount relationships



Profitability Analysis

Spot most profitable categories/regions



Visual Insights

Communicate results with clear charts



Predictive Modeling

Forecast future profit and sales

_Team_Project_eComm_Sales

m_Sales

Current iteration

Roadmap

My items

View 6

In progress 2 / 5

Estimate: 0

This is actively being worked on

UofT_DSI_C7_Team_Project_eComm_Sales #23

complete the presentation

UofT_DSI_C7_Team_Project_eComm_Sales #25

Update the readme

Done

This has

#7
add cor

UofT #16

Create l

UofT #17

regressi

UofT #18

Revamp

UofT #19

make p

UofT #20

review a

Team & Responsibilities



Vikrama

Business proposal development, project naming, data cleaning, experiments, and code review



Iryna

Dataset selection, data exploration, KPI development, and setting business goals



Paul

Repository management, project tracking, data exploration and cleaning, slides, visualization, and presentation management

Key Data Parameters

1

Price

Pricing & revenue

2

Product Quantity

Volume & inventory

3

State

Regional performance

4

Repeat orders

Driving returning customers

5

Date

Temporal trends

```
# Train/test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Fit regression model
model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)

# Evaluate
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
r2 = r2_score(y_test, y_pred)

print(f"\n SKU Regression Results:")
print(f"RMSE: {rmse:.2f}")
print(f"R²: {r2:.4f}")

# Feature importance
coefficients = pd.Series(model.coef_, index=X.columns).sort_values(ascending=False)
print("\n Top 15 Positive Features (drive higher sales):\n", coefficients.head(15))
print("\n Top 15 Negative Features (lower sales):\n", coefficients.tail(15))

# Find top-selling SKUs
print("\n Top 10 Most Popular SKUs:\n")
print(agg_df[['SKU', 'Category', 'Size', 'Qty']].sort_values('Qty', ascending=False).head(10))

# Read data
data_path = "Users\\Paul\\AppData\\Local\\Temp\\ipykernel_9260\\99505703.py:22: Dt"
df = pd.read_csv(data_path)

# Regression Results:
RMSE: 19.35
R²: 0.0455

Top 15 Positive Features (drive higher sales):
Category_Western Dress      14.529545
Category_Set                 7.448648
Size_6XL                    7.118155
Category_kurta              6.180175
```

Data Cleaning Strategy

Essential Steps

- Remove unnecessary columns (index, promotion-ids)
- Handle missing values in Courier Status, currency, Amount, ship-city, ship-state
- Convert Date to proper format for time-based analysis
- Standardize text formats to lowercase

Data Enhancement

- Check and remove duplicate Order IDs
- Group similar order statuses for consistency
- Ensure currency consistency (convert to INR)
- Add derived columns: Profit Margin, Month/Year extractions

Analysis Focus Areas



Product Insights

Top selling and highest revenue categories



Regional Insights

Most profitable states and cities



Order Performance

Amazon vs Merchant fulfillment comparison

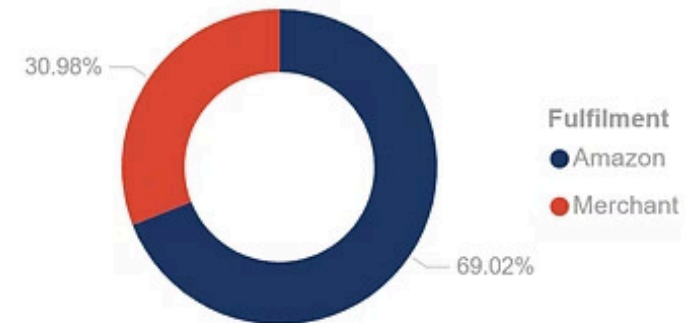


Sales Patterns

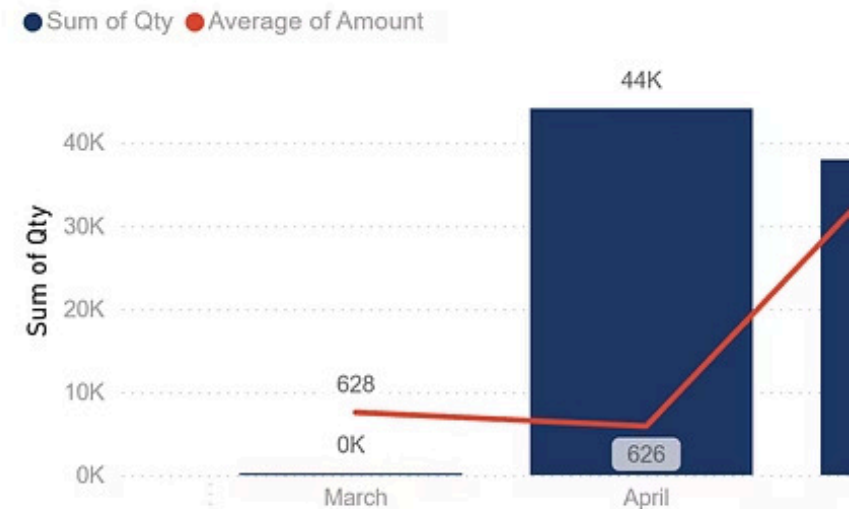
Impact of discounts, promotions and **repeat orders**

Amazon Sales Report

Orders by Fulfilment



Qty Sold Month Ov



unt) by Order Status

19,284

0M 20M 30M 40M 50M 60M

Total Sales (Amount)

Visualizations

- **Plotly, Seaborn, Matplotlib:** Compare product performance across types.
- **Power BI:** Track shipped, cancelled, and pending orders.
- **Monthly Sales Trends:** Identify seasonal spikes and patterns.
- **Correlation Heatmap:** Understand relationships between sales, quantity, and discounts.

Risks & Challenges

Profit Margin Calculation

No clear methodology to calculate profit margins given current dataset limitations

Cost Data Gaps

Lack of cost and Bill of Materials (BOM) data restricts deeper financial analysis

Time Constraints

Limited project timeline requires prioritization of key deliverables and analysis

PM methodology

GitHub Project used with agile methods.

DO NOT fall asleep at the

Open 0 / 1



pradziie opened 6 minutes ago

...



Sub-issues 0 of 1

More coffee! #24

Create sub-issue



pradziie self-assigned this 6 minutes ago



pradziie moved this to Todo in UofT_DSI_C7_Team_Project_eComm_Sales



pradziie added this to UofT_DSI_C7_Team_Project_eComm_Sales



pradziie added help wanted 6 minutes ago



Made with GAMMA

More coffee! #24

Expected Outcomes

01

Clean Dataset

Organized, standardized data

03

Visual Dashboard

Interactive data insights

05

Actionable Recommendations

Data-driven strategies for growth

02

Profitability Insights

Identify key profit drivers

04

Predictive Model

Future sales forecasting

06

Bonus

UofT **color schemes** for Excel and
Power BI (available in readme)

