# PROJECT TITLE: CUSTOMER SEGMENTATION AND BEHAVIOURAL ANALYSIS FOR INSTACART BASKET CASE STUDY
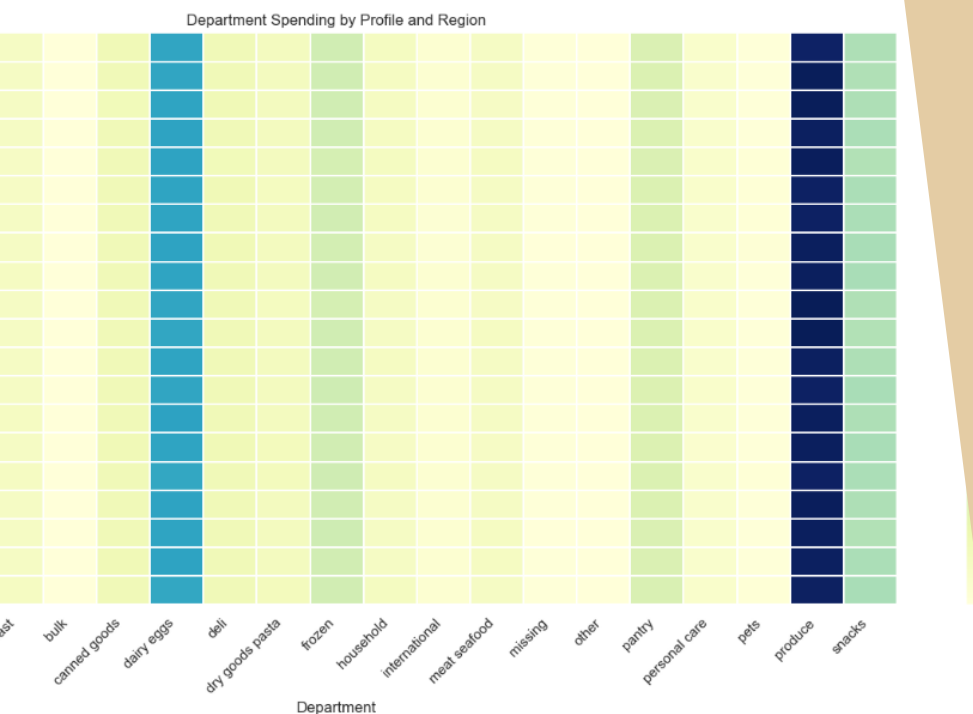
HENRY ARKOH

CLAYMAN

## Overview

Segmenting Instacart customers using demographic and behavioral data to identify patterns in spending habits, time-of-day activity, and regional preferences.

## Purpose and Context

By analyzing millions of real purchase records, we uncovered:

- **When** different types of customers shop

- **What** departments they prefer

- **How** income and region affect price sensitivity and product choice

## Objective

The goal was to identify, categorize, and analyze customer segments based on key demographic indicators (e.g., age, dependents, income) and behavioral trends (e.g., time of purchase, department preferences).

## Duration

The project spanned **4 weeks**. It was completed on time, structured in weekly milestones including data preparation, analysis, visualization, and presentation.

## Credit

Role: Data Analyst

## Tools, Skills, and Methodologies

**Technologies**: Python (pandas, matplotlib, seaborn), Tableau, Jupyter Notebook

Distribution of Customer Profiles

**Data Cleaning and Preparation**

Ensured Consistency, Removed Invalid Entries, And dealt with all data privacy issues

**Customer Profile Segmentation**

Created profiles using age and number of dependents.

Below is a visual summary:

```
[24]: profile_counts
```

| | Profile | Count |
|---|---|---|
| 0 | Other | 10966597 |
| 1 | Young Parent | 6326031 |
| 2 | Middle-aged Family | 3980187 |
| 3 | Senior Solo | 2000467 |
| 4 | Young Single | 1137884 |

Most customers are either young parents or don't fit a spe

```
[47]:  # Transpose the DataFrame so hours (0-23) are on
       hour_crosstab_transposed = hour_crosstab.T

       # Plot each profile as a line
       plt.figure(figsize=(12, 6))
       for profile in hour_crosstab_transposed.columns:
           plt.plot(hour_crosstab_transposed.index, hour_c

       plt.title('Ordering Patterns by Hour of Day and Cust
       plt.xlabel('Hour of Day')
       plt.ylabel('Percentage of Orders')
       plt.xticks(range(0, 24))
       plt.grid(True)
       plt.legend(title='Profile')
       plt.tight_layout()
       plt.savefig(os.path.join(path, '04 Analysis', 'Visuali
       plt.show()
```

**Behavioral Analysis**

Explored order behavior by time of day across profiles. Key trends include evening activity for young parents and morning preference for seniors.

Young Singles and Parents tend to shop later in the morning and early afternoon, while Senior Solos and Middle-aged Families show steady activity earlier in the day

Ordering Patterns by Hour of Day and Customer Profile

Profile
— Middle-age
— Other
— Senior Solo
— Young Pare
— Young Singl

Hour of Day
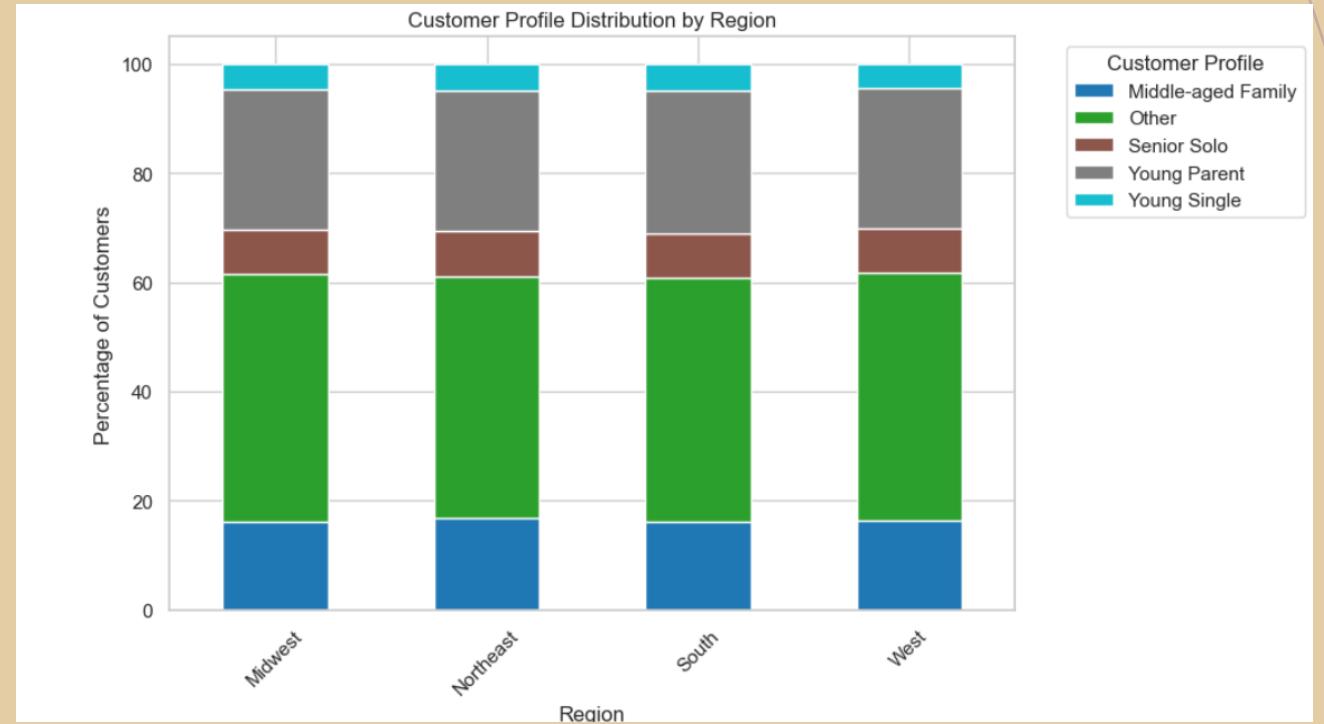
# Income vs. Price sensitivity

Explored how income affects spending patterns. Higher income groups showed more stable purchasing behavior.
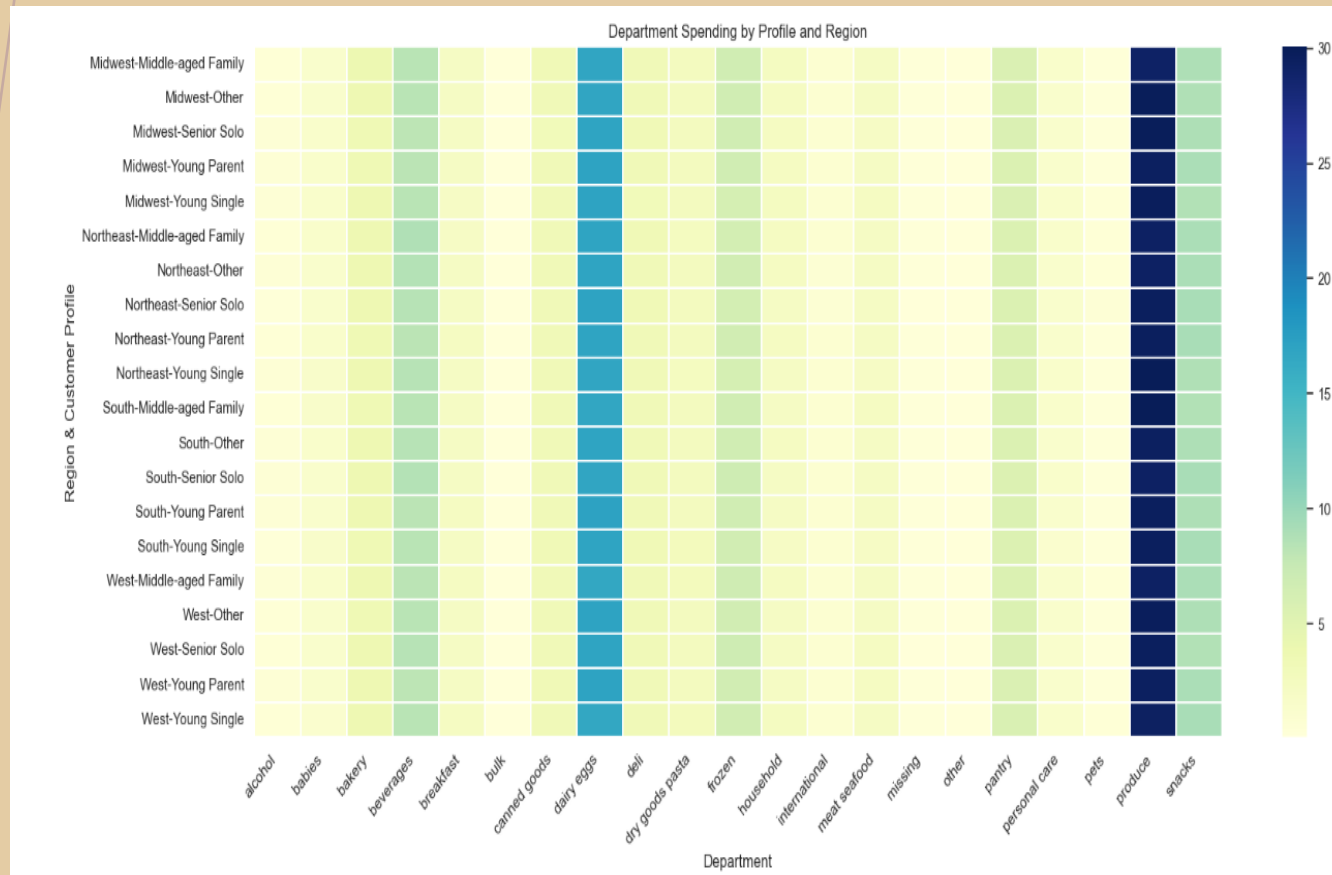
## Customer Profile vs Regional Distribution

The West might have more Young Singles, while the South has more Middle-aged Families. That helps us target the right products and promotions for each region.



Customer Profile Distribution by Region

**Heatmap showing average purchase frequency across regions and profiles.**



Department Spending by Profile and Region

Middle-aged Families and Young Parents are more prominent in the Midwest and South, where they frequently purchase from traditional, family-oriented departments like dairy, pantry, and produce. In contrast, Young Singles and Senior Solos are more common in the West and Northeast, showing stronger preferences for convenience-focused and self-care items such as snacks, beverages, and personal care. While regional differences exist, profile-based shopping behavior remains relatively consistent, making customer profiles a reliable basis for targeted marketing and regional merchandising strategies.

## Challenges Faced:

1. **Data Quality and Completeness** – Inaccurate, missing, or inconsistent demographic or behavioral data can hinder reliable profiling and analysis.

2. **Privacy and Ethical Considerations** – Handling sensitive personal information necessitates strict adherence to data privacy laws and ethical standards.

3. **Overgeneralization** – There's a risk of creating profiles that are too broad, leading to ineffective targeting or excluding nuanced customer behaviors.

4. **Dynamic Customer Behavior** – Preferences and purchasing patterns can change over time, making it difficult to maintain accurate, up-to-date profiles.

5. **Scalability** – Applying profiling methods effectively across large datasets while maintaining performance can be computationally challenging.

## Final Reflections & Next Steps

1. This project reinforced my ability to extract insights from messy real-world data. I developed strengths in behavioral segmentation, data visualization, and strategic communication.

Next Steps:

- Automate real-time dashboards for ongoing monitoring.

- Apply clustering (e.g., K-means) to refine profile creation.

- Expand analysis to include loyalty and promotional responses.

This case study showcases my capability to turn raw data into actionable insight with real business value.