# Optimizing Retail with Time-Based Insights and Machine Learning in Online Grocery

Team#027: Yaoyan Zhou, Wan Wang-Geissler, Hangxing Sha, Nai Ning Chi, Yan Huan

## 1. INTRODUCTION AND MOTIVATION

Online shopping has surged, especially during the COVID-19 pandemic, with Instacart at the forefront, streamlining grocery purchases with an easy-to-use interface connecting multiple stores. However, one vital element has often been neglected: the timing of customer purchases.

Our project is inspired by the vast potential for enhancing the online retail experience through time-based recommendation systems. Drawing on the wealth of data provided by Instacart and harnessing the power of advanced machine learning techniques, our aim is to inject the crucial 'when' factor into sales strategies. The focus lies in dynamically adjusting prices and effectively managing stock based on temporal insights derived from customer behavior.

## 2. PROBLEM DEFINITION

At the core of our initiative is the development of a tool designed to significantly elevate retailers' efficiency and competitive edge by incorporating personalized, time-sensitive recommendations. Simultaneously, this tool aims to empower consumers with timely and contextually relevant suggestions. Our primary objective is to craft a sophisticated machine-learning tool utilizing Instacart's comprehensive dataset. This tool will be dedicated to optimizing the timing of promotions, catering to the needs of both retailers and customers in the online grocery shopping domain.

## 3. LITERATURE SURVEY

Extant research on online grocery shopping has delved into understanding customer behavior, scrutinizing phenomena such as the disparity in customer page views across online and offline platforms [1], and the propensity for impulse purchases in the digital shopping environment [2]. The marketing domain has employed both qualitative and quantitative analyses to bolster customer retention [3]. Market Basket Analysis (MBA) stands out as a key data mining technique in retail, which helps in uncovering common item pairings. This has proven pivotal for strategies such as cross-selling and enhancing customer contentment [4], and is influential in optimizing store layouts [5].

MBA, in conjunction with predictive models like Markov Chains and clustering algorithms, aids in the creation of bespoke recommendation systems for online grocery patrons [6], enhancing both the consumer experience and business profitability. While these systems are grounded in historical purchase data [7], novel methodologies like the ARM-Predictor Algorithm and Transaction-Data Based Real-time Preference Inference Engine (TRPIE) have emerged to incorporate real-time data [8][9]. Despite the introduction of sophisticated techniques such as Learning-to-Rank (LRT), Determinantal Point Process (DPP) [10][11][12], Diversified Collaborative Filtering (DCF) [13], and linear combinations of rating functions with entropy regularizers [14], the domain lacks initiatives that factor in the timing of purchases. The recent works on diversified recommendations underscore the need for innovation in this area [15]. Bridging this gap by analyzing temporal consumer behavior is the cornerstone of our project's mission.

## 4. PROPOSED METHODS

### 4.1 Data Description and Cleaning

We sourced our data from a publicly available Kaggle competition, consisting of six interrelated CSV files that represent over 3 million grocery orders from 200,000+ Instacart users along with the order's day and time, department, and aisle information, all linked via unique identifiers.

Following the consolidation of the CSV files into a unified dataset, we embarked on an Exploratory Data Analysis (EDA) using Tableau to create an interactive dashboard. This visualization tool offered a user-friendly platform for insights that might be less accessible through traditional programming-based analysis. Due to the dashboard's resolution configuration (1920px x 1080px), direct embedding into this document may compromise visibility (refer to *Figure 1*). However, the Tableau workbook is available upon request

for detailed examination.

In alignment with our project's focus on temporal shopping patterns, we honed in on specific data dimensions pertinent to our research. The dashboard features a heat map that elucidates the correlation between order volumes, days of the week, and hours of the day. This visualization confirms peak order times predominantly span from 9 AM to 5 PM, with Mondays and weekends being the busiest periods. An area chart presents the ordering frequency, revealing prominent spikes at one-week and one-month intervals, suggesting weekly and monthly shopping cycles. A bubble chart details department-wise order volumes, where the size of each bubble is proportionate to the order count, identifying produce, dairy, eggs, snacks, beverages, and frozen goods as the top-selling categories.

Completing the dashboard are two bar charts, ranking aisles and products by order frequency in descending order. Interactivity is a cornerstone of this dashboard, allowing for dynamic filtering based on user ID and whether an item was reordered. This holistic interactive setup provides a multidimensional view of the dataset, setting the stage for the deeper temporal analysis that underpins our project's innovation.
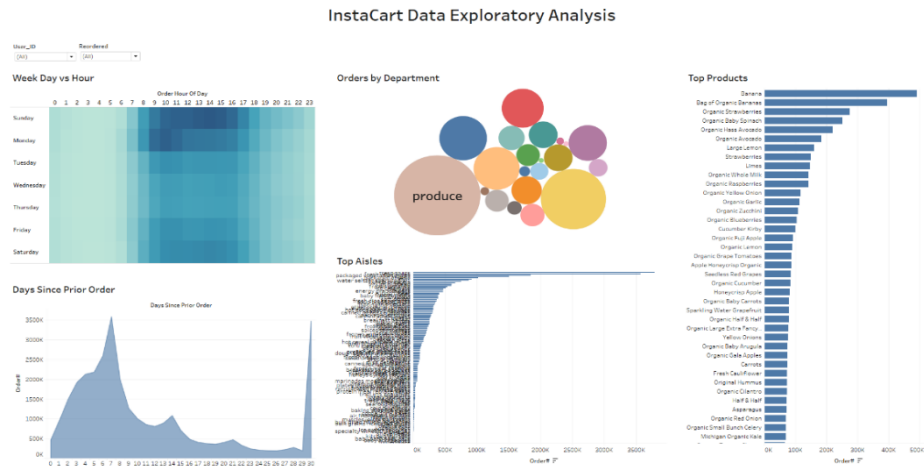


*Figure 1 Customized Tableau Dashboard for EDA*

## 4.2 Machine Learning Models/Algorithms

The EDA provided a macroscopic view of customer purchasing patterns, revealing prevalent shopping times throughout the week. To distill how these trends translate to individual stores and shoppers, we applied a range of machine learning techniques to refine our recommendations.

For new user acquisition by stores, we implemented K-means clustering, an unsupervised learning algorithm, to segment users into distinct groups. This method helped identify clusters based on re-order patterns, order timing, and purchase days, using the elbow method to determine the optimal number of clusters. For each cluster, we recommended the top three products, aiming to personalize the outreach for potential new customers.

In addressing existing customers, we turned to supervised learning algorithms, specifically Decision Trees, Logistic Regression, and AdaBoost, to predict customer behavior with greater precision. The Decision Tree model mapped decisions within a tree-like framework, Logistic Regression estimated the probability of purchase using a logistic function, and AdaBoost enhanced prediction accuracy by combining multiple learners. The collective output from these models was calibrated to minimize Mean Squared Error (MSE), ensuring the most accurate product recommendations for specific times and days.

Lastly, to advise on the optimal days and times for shopping in different departments, we employed Linear Regression, Random Forest, and XGBoost algorithms. Linear Regression sought relationships between time slots and purchasing behavior, while Random Forest and XGBoost—both ensemble tree-based methods—worked on subsets of data to prevent overfitting and used gradient boosting, respectively, for performance gains. Recommendations were finalized based on models with the lowest MSE, delivering

strategic insights into the best shopping times for various departments.

## 4.3 User Interfaces

As depicted in *Figure 2* and *Figure 3*, our interactive interface, titled "Time-Based Smart Product Recommendation," begins with an inviting welcome page. This page prominently features the title, along with the Instacart logo positioned at the top right and a decorative picture of grocery goods at the bottom. Central to this page is a blue button with white text stating "Go to Your Business Account", which serves as the portal to the main interface after the store owner logs in. This transition seamlessly redirects users to a sophisticated dashboard created with Tableau, accessible via a Python Dash application.
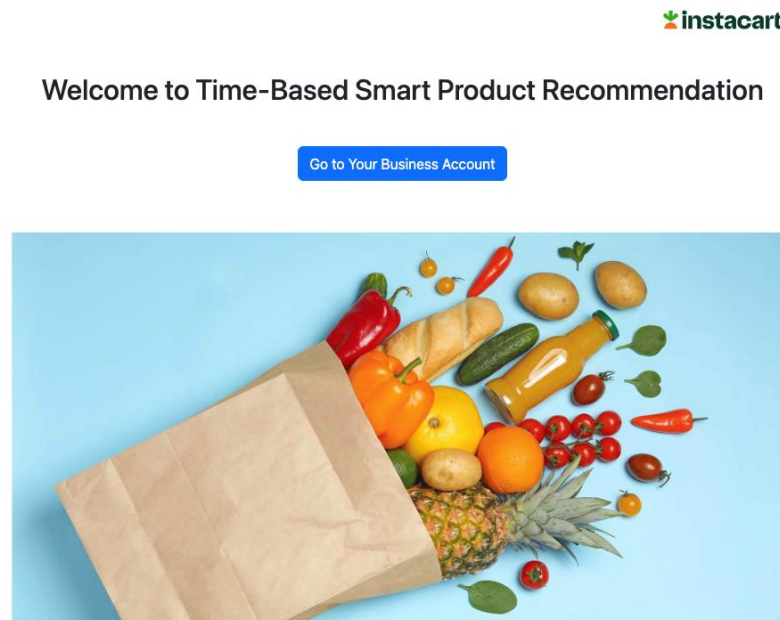


*Figure 2. Welcome Page in Python Dash application*

The dashboard's primary workspace, as illustrated in Figure 3, is meticulously organized with various interactive components. At the top, users find three filters—user_id, department, and product_id—allowing for a customized view of customer shopping behavior. A key feature, shown in the top left of the dashboard, is a heat map that reveals optimal shopping times by displaying data across a 24-hour x-axis and a 7-day y-axis. The user can toggle to the timeframe, and it will reveal the best day of week, best hour and count of items. Adjacent to this, in the lower left section, packed bubbles visually represent different customer groups. The right half of the dashboard is dedicated to a detailed form showcasing department, product_id, and product names. This layout ensures that any selection made via the filters or any interaction with the dashboard elements dynamically updates the displayed data, providing real-time insights and analytics tailored to the user's needs.
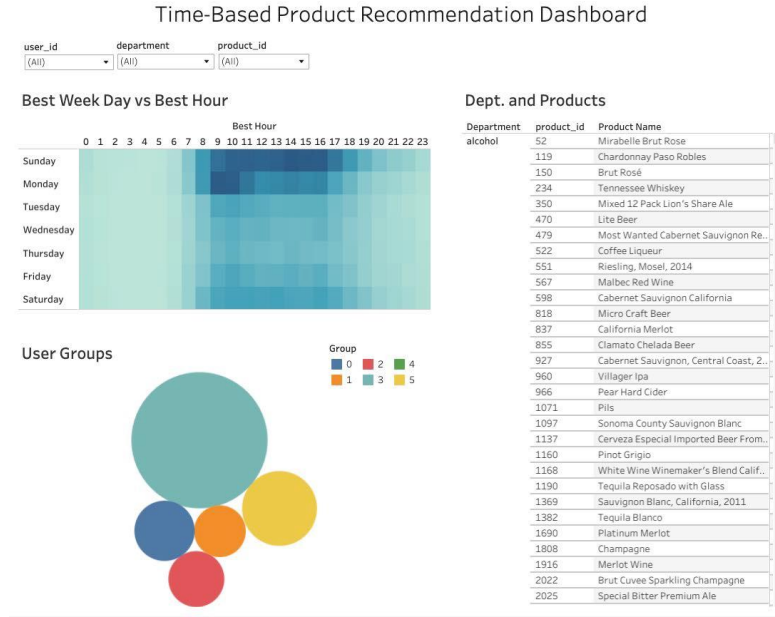
*Figure 3. Time-Based Product Recommendation Dashboard in Tableau.*

# 5. EXPERIMENTS AND EVALUATION

## 5.1 Evaluation of Prediction Models

Our customer profiling evaluation starts with K-means clustering, using the elbow method to determine the optimal number of clusters for precise segmentation by purchase timing. We allocate 80% of our dataset for training and 20% for testing for additional features. For regression models, we use Mean Squared Error (MSE) as the evaluation metric, which quantifies the accuracy of predictions in comparison to actual shopper behavior, thus offering a robust measure of our predictive methodology.

### 5.1.1 Customer Profiling and Cluster Visualization

Utilizing unsupervised learning through k-means clustering, we employ the Elbow method to determine the optimal number of clusters (k) for customer profiling. This approach allows for the creation of distinct customer clusters.

### 5.1.2 Smart Cart/Shopping List

For smart cart predictions, we utilize Decision Tree, Logistic Regression, and AdaBoost algorithms, achieving accuracy scores of 63.41%, 68.69%, and 70.18%, respectively. AdaBoost proves to be the most efficient, leading our smart cart generation process.

### 5.1.3 Dynamic Time-Sensitive Inventory Quantity

For predicting inventory levels, Decision Tree, Logistic Regression, and AdaBoost algorithms are employed. The accuracy scores for these algorithms are 64.07%, 68.73%, and 70.18%, respectively. AdaBoost stands out as the optimal algorithm for inventory prediction, influencing the generation of inventory recommendations in our model.

### 5.1.4 Best Time to Shop by Department

Linear Regression, Random Forest, and XGBoost models are employed to train and predict the best time to shop by department. Model performance is evaluated using Mean Squared Error (MSE). XGBoost emerges as the best model with the lowest MSE, providing valuable insights into the optimal timing for customer shopping preferences within different departments.

## 5.2 Evaluation of Interface

This section details the methods and observations from our comprehensive evaluation of the "Time-Based Smart Product Recommendation" interface. Our evaluation included three primary methods: a survey, individual interviews, and a case study walkthrough. Each method targeted participants with backgrounds in retail and business, using ten fictional store profiles (See *Appendices 7.1*) to simulate real-world applications.

### 5.2.1 Methodology

1.   **Survey**: Conducted with 30 participants, each assigned a fictional store profile. The survey consisted of structured questions using Likert scales and multiple-choice formats (See *Appendices 7.2*) to assess various aspects of the interface, including effectiveness, clarity, user-friendliness, and overall experience.
2.   **Interviews**: Five additional participants were selected and assigned fictional store profiles. They participated in detailed interviews with open-ended questions (See *Appendices 7.3*), providing qualitative insights into their experiences with the interface.
3.   **Case Study Walkthrough**: A focused case study involving a fictional liquor store owner was carried out to observe the interface's application in a specific retail context. The walkthrough highlighted how the dashboard tools, like filters and heat maps, aid in inventory and sales decision-making.

### 5.2.2 Observations

1.   **Survey Findings**:
     - The interface was rated highly for its effectiveness in inventory management, particularly by owners of diverse and specialized stores.
     - The heat map and data presentation tools were well received, with suggestions for more detailed segmentation.
     - Recommendations for improvements focused on enhancing user guidance and feature clarity.
2.   **Interview Insights**:
     - Participants provided unique perspectives based on their store profiles, praising the interface's user-friendliness and data visualization capabilities.
     - Suggestions for tailored features and enhanced algorithms for diverse customer bases were prominent.
     - The need for more customized functionalities for budget management and real-time updates was highlighted.
3.   **Case Study Walkthrough**:
     - Demonstrated the practical application of the interface for a liquor store owner, focusing on product-specific sales trends.
     - The interface facilitated strategic stocking and sales decisions, especially for high-demand products like Mirabelle Brut Rose.
     - The case study exemplified the interface's ability to provide granular insights for optimizing product management and enhancing business performance.

### 5.2.3                                                                                                     Conclusion

The comprehensive evaluation of the "Time-Based Smart Product Recommendation" interface through surveys, interviews, and case study walkthrough has demonstrated its effectiveness and utility in various retail contexts. The interface shows potential in assisting store owners in making data-driven decisions for inventory management and customer engagement. The feedback and suggestions collected will guide future enhancements, ensuring the interface remains adaptable and valuable for diverse retail needs.

## 6. CONCLUSIONS AND DISCUSSION

Our endeavor aimed to revolutionize the retail landscape by introducing time-based insights into online grocery shopping through a multifaceted approach. Leveraging Instacart's expansive dataset and employing advanced machine learning techniques, we developed a comprehensive recommendation system and an innovative visualization tool.

## 6.1 Conclusions

The integration of temporal insights into recommendation systems is pivotal in advancing the online retail sphere. By catering to the 'when' of customer behavior, retailers can enhance customer satisfaction, streamline inventory management, and bolster their competitive edge in the market.

- Temporal Insight Integration: Our study successfully unearthed intricate temporal patterns in customer purchasing behavior, highlighting peak shopping hours, weekly and monthly shopping cycles, and popular department-wise preferences.
- Machine Learning Recommendations: Through a range of machine learning models, we personalized recommendations for both new and existing customers, optimizing product suggestions based on timing, user segments, and departmental considerations.
- Innovative Visualization Tool: Our prototype showcased a user-friendly interface offering real-time insights into optimal shopping times, customer profiling, and inventory management.

## 6.2 Future Directions

The integration of temporal insights into recommendation systems is pivotal in advancing the online retail sphere. By catering to the 'when' of customer behavior, retailers can enhance customer satisfaction, streamline inventory management, and bolster their competitive edge in the market.

Our project marks a significant step toward temporal integration in online retail. The intersection of machine learning, temporal insights, and user-friendly interfaces holds immense promise in reshaping the future of online grocery shopping. However, there are limitations and challenges for future exploration and refinement:

- Real-time Adaptability: Focusing on strategies to incorporate real-time data updates into our models could significantly improve accuracy.
- User-Centric Refinement: Continual iterations of the visualization tool based on user feedback will ensure it remains intuitive and beneficial for both retailers and shoppers.
- Holistic Integration: Exploring ways to merge our recommendations with broader marketing strategies to create a seamless shopping experience for users.

Given that this project is part of a class assignment, we faced certain constraints in terms of time and access to real-world data. The evaluation of the interface relied on fictional stores provided to participants. In the future, it's crucial to:

- Conduct Real-World Testing: Engage actual business owners and real-world retail environments to test the interface. This will provide more authentic feedback and insights into how the tool performs under practical, everyday business conditions.
- Expand Dataset and Processing Capabilities: The project used a limited dataset due to local computing constraints. Future iterations should utilize comprehensive datasets and leverage cloud computing or advanced processing capabilities. This allows for more complex models and extensive parameter tuning to enhance our machine learning models' performance and accuracy.
- Iterative Model Improvement: With better computational resources, there's an opportunity to refine our models. This includes implementing advanced machine learning algorithms and exploring a range of parameters to optimize performance.
- Longitudinal Studies for Temporal Analysis: Conducting long-term studies to better understand how consumer behavior evolves over time. This will help in refining our temporal integration strategies and ensuring our models stay relevant and effective.

By addressing these areas, we aim to build upon the foundational work of this project, enhancing its applicability and effectiveness in the dynamic landscape of online retail.

Finally, all team members have contributed a similar amount of effort.

# 8. REFERENCES

1. Anesbury, Z., Nenycz-Thiel, M., Dawes, J. and Kennedy, R. (2016) How Do Shoppers Behave Online? An Observational Study of Online Grocery Shopping. Journal of Consumer Behavior, 15, 261-270. https://doi.org/10.1002/cb.1566.

2. Huyghe, E., Verstraeten, J., Geuens, M., & Van Kerckhove, A. (2017). Clicks as a Healthy Alternative to Bricks: How Online Grocery Shopping Reduces Vice Purchases. Journal of Marketing Research, 54(1), 61–74. http://www.jstor.org/stable/44878488.

3. Bapat, G. S., & Vishwanath Karad, M. I. T. Online Grocery Shopping: A STUDY OF CONSUMER BEHAVIOUR ON STAYING AND SWITCHING BETWEEN AMAZON, BIG BASKET AND GROFERS.

4. Zamil, A. M. A., Al Adwan, A., & Vasista, T. G. (2020). Enhancing customer loyalty with market basket analysis using innovative methods: a python implementation approach. International Journal of Innovation, Creativity and Change, 14(2), 1351-1368.

5. Halim, S., Octavia, T., & Alianto, C. (2019). Designing facility layout of an amusement arcade using market basket analysis. Procedia Computer Science, 161, 623-629.

6. Ilyas, Q. M., Mehmood, A., Ahmad, A., & Ahmad, M. (2022). A Systematic Study on a Customer's Next-Items Recommendation Techniques. Sustainability, 14(12), 7175.

7. Chabane, N., Bouaoune, A., Tighilt, R., Abdar, M., Boc, A., Lord, E., ... & Makarenkov, V. (2022). Intelligent personalized shopping recommendation using clustering and supervised machine learning algorithms. Plos one, 17(12), e0278364.

8. Kaur, M., & Kang, S. (2016). Market Basket Analysis: Identify the changing trends of market data using association rule mining. Procedia computer science, 85, 78-85.

9. Li, Y. R., Hwang, T. K., & Chang, S. C. (2018, October). Dynamic Inference of Personal Preference for Next-to-Purchase Items by Using Online Shopping Data. In CS & IT Conference Proceedings (Vol. 8, No. 14). CS & IT Conference Proceedings.

10. Chen, Y., Leong, Y. C., Yiing, L. S., & Xiao, Y. (2023). Study on the Influence of Knowledge-driven Technology on predicting consumer Repurchase Behaviour. International Journal of Communication Networks and Information Security, 15(1), 109-117.

11. Tatiana, K., & Mikhail, M. (2018). Market basket analysis of heterogeneous data sources for recommendation system improvement. Procedia Computer Science, 136, 246-254.

12. Ma, L., Sinha, N., Cho, J. H., Kumar, S., & Achan, K. (2023). Personalized diversification of complementary recommendations with user preference in online grocery. Frontiers in big Data, 6, 974072.

13. Lijing Qin and Xiaoyan Zhu. Promoting diversity in recommendation by entropy regularizer. In IJCAI, pages 2698–2704, 2013.

14. Peizhe Cheng, Shuaiqiang Wang, Jun Ma, Jiankai Sun, and Hui Xiong. Learning to recommend accurate and diverse items. In WWW, pages 183–192, 2017.

15. Wu Q, Liu Y, Miao C, Zhao Y, Guan L, Tang H. Recent advances in diversified recommendation. arXiv preprint arXiv:1905.06589. 2019

# 7. APPENDICES

## 7.1 Ten Fiction Store Profiles for Evaluation Purpose

1. Green Valley Market
Location: Suburban area
Specializes in: Organic and locally-sourced produce, artisanal breads, and eco-friendly products.
Target Customers: Health-conscious families and environmentally aware individuals.

2. Urban Oasis Grocers
Location: Downtown city
Specializes in: Gourmet foods, international cuisine ingredients, and ready-to-eat meals.
Target Customers: Busy professionals and culinary enthusiasts.

3. The Family Pantry
Location: Small town
Specializes in: Affordable staples, bulk goods, and family-sized packs.
Target Customers: Large families and budget shoppers.

4. Harvest Health Foods
Location: University town
Specializes in: Vegan, gluten-free, and dietary specific foods.
Target Customers: Health-conscious students and faculty.

5. Spice Route Emporium
Location: Culturally diverse neighborhood
Specializes in: Ethnic spices, halal meats, and exotic fruits.
Target Customers: Immigrant communities and culinary adventurers.

6. Blue Coast Seafood Market
Location: Coastal area
Specializes in: Fresh seafood, seasonal catches, and maritime delicacies.
Target Customers: Seafood lovers and high-end clientele.

7. Sweet Tooth Confectionery
Location: Tourist hotspot
Specializes in: Gourmet chocolates, artisanal sweets, and dessert baking supplies.
Target Customers: Tourists and dessert enthusiasts.

8. The Rustic Mill Bakery & Deli
Location: Rural area
Specializes in: Freshly baked bread, local cheeses, and handcrafted deli meats.
Target Customers: Locals seeking quality staples and artisanal products.

9. Peak Performance Sports Nutrition
Location: Near fitness centers and sports arenas
Specializes in: Sports nutrition products, energy snacks, and health supplements.
Target Customers: Athletes, fitness enthusiasts, and health-conscious individuals.

10. Bountiful Harvest Co-op
Location: Eco-friendly community
Specializes in: Community-sourced produce, bulk grains, and sustainable household items.

Target Customers: Eco-conscious consumers and supporters of local agriculture.

## 7.2 Survey Questions

1. How effective do you find the interface in managing your store's inventory?
   1 - Not Effective At All
   2 - Slightly Effective
   3 - Moderately Effective
   4 - Very Effective
   5 - Extremely Effective
2. On a scale of 1-5, rate the usefulness of the heat map in planning your inventory.
   1 - Not Useful At All
   2 - Slightly Useful
   3 - Moderately Useful
   4 - Very Useful
   5 - Extremely Useful
3. Was the purpose of each feature clear and understandable?
   1 - Not Clear At All
   2 - Slightly Clear
   3 - Moderately Clear
   4 - Very Clear
   5 - Extremely
4. Clear Rate the overall user-friendliness of the interface.
   1 - Not User-Friendly At All
   2 - Slightly User-Friendly
   3 - Moderately User-Friendly
   4 - Very User-Friendly
   5 - Extremely User-Friendly
5. How effectively does the heat map convey information about peak shopping times?
   1 - Not Effective At All
   2 - Slightly Effective
   3 - Moderately Effective
   4 - Very Effective
   5 - Extremely Effective
6. Are the packed bubbles helpful in understanding customer groups?
   1 - Not Helpful At All
   2 - Slightly Helpful
   3 - Moderately Helpful
   4 - Very Helpful
   5 - Extremely Helpful
7. What was your overall experience using the interface?
   1 - Very Poor
   2 - Poor
   3 - Average
   4 - Good
   5 - Excellent
8. Would you recommend this interface to other store owners?
   Yes/No/Maybe

## 7.3 Interview Question List
**Initial Impressions:**

- What were your first impressions upon seeing the welcome page?
- How did you find the overall look and feel of the interface?

**Functionality and Features:**
- Can you describe how you used the filters (user_id, department, product_id)?
- What did you find most useful about the heat map and packed bubbles?
- Were there any features you found unnecessary or confusing?

**Data Interpretation:**
- How did you interpret the data presented in the heat map and the long form?
- Were there any visual elements that enhanced or hindered your understanding of the data?

**User Experience and Improvement:**
- What was the most satisfying aspect of using the interface?
- Did you encounter any difficulties while navigating through the interface?
- How would you improve the interface based on your experience?

**Final Thoughts:**
- How do you think this interface would impact your store's operations if it were real?
- Do you have any additional comments or suggestions for the interface?