Optimizing Personalized Promotions at ACSE Supermarket: Focusing on Kellogg's Targeted

Recommender System Strategy

AI & ML at Scale Team 1

Lanston Chen, Songbo Hu, Kedi Lin, Jie Mei, Jayson Xu, Wenxi Xu

Emory Goizueta Business School

**Executive Summary**

ACSE Supermarket, leveraging its expansive presence across North America, is setting a new standard in marketing by transitioning from broad to highly personalized promotions. Our consulting team was tasked with developing a personalized recommender system targeting Kellogg's products to customers who typically buy General Mills items. Our approach began with a targeted analysis of customer transactions, concentrating on distinct purchasing patterns while also incorporating strategic data processing at the business level.

In the implementation, we utilized three types of models: Content-Based filtering, Collaborative Filtering, and a Naive Bayesian model. These models were designed to predict potential Kellogg's product purchases, streamlining the promotional strategy to be as relevant and effective as possible.

The initiative is expected to refine ACSE's promotional campaigns by making them more customer-specific, which in turn, enhances customer satisfaction and loyalty. This approach also presents a cost-effective marketing solution, optimizing promotional spending by targeting likely buyers rather than employing a one-size-fits-all strategy. The broader business impacts are significant, potentially increasing sales volume and improving inventory turnover through better-aligned product offerings and promotions. This project not only demonstrates ACSE's capability to adapt and innovate but also provides a scalable model that can be generalized to other products and campaigns in the future, ensuring sustained business growth and competitive advantage.

**Business Understanding**

**Problem Definition:**

The primary challenge addressed by our recommender system (RS) revolves around the competitive dynamics between two major consumer daily goods brands, General Mills (GM) and Kellogg's (K). Both brands cater to similar market segments, selling essential consumer food products across a broad demographic. The specific problem our RS aims to solve is capturing the

loyalty of consumers who predominantly purchase products from General Mills and persuading them to consider Kellogg's as an alternative. This strategic pivot is facilitated by leveraging detailed consumer purchasing data and behavioral insights to identify potential Kellogg's adopters among steadfast General Mills customers. The RS system is designed not only to suggest Kellogg's products that align with the consumers' preferences but also to enhance the visibility and appeal of Kellogg's offerings through personalized promotions.

**Expected Outcomes:**

Upon implementation of our RS, we anticipate several measurable outcomes that underscore the efficacy and impact of the system. Firstly, the RS is expected to accurately identify segments within the General Mills consumer base that show the highest propensity for crossover to Kellogg's products. By targeting these segments, ACSE can more effectively allocate marketing resources, thereby improving cost-efficiency. Additionally, by tracking the conversion rates and customer feedback, we will gain insights into the effectiveness of personalized promotions in real-time, allowing for ongoing optimization. The broader business impacts are substantial; we foresee improvements in customer retention rates, increased sales volumes for Kellogg's products, and a stronger market presence in competitive retail environments. The scalability and adaptability of this RS mean that its framework can be generalized to other product lines and promotional campaigns, enhancing overall business agility and strategic marketing capabilities.

## Data Understanding

**Data Overview:**

The refined dataset from ACSE Supermarket, now labeled "transactions_clean," has undergone significant cleaning and restructuring to optimize it for analytical purposes. This dataset, derived from the original transactions records spanning from 2017 to 2020, has been tailored to focus more sharply on relevant customer and product interactions. Key enhancements include the removal of duplicate trans_id and cust_id pairs, exclusion of products ranking in the

bottom 10% in sales, and elimination of non-relevant product categories such as 'Customer Service' and various types of coupons. Additionally, transactions from outlier store 8540 have been excluded, and only customer records showing consistent activity up until 2020 have been retained. These measures ensure that the data reflects genuine, meaningful consumer behavior and product performance.

**Transactions Table**

The Transactions Table, post-cleaning, provides a robust framework for analyzing customer purchase paths and store performance metrics across ACSE's network. It now contains only entries with positive sales figures and volumes, ensuring the accuracy of sales and demand analyses. The Products Table, streamlined to exclude less significant entries, offers detailed insights into over 100,000 products. This table is crucial for evaluating product assortment strategies, pricing decisions, and promotional activities, allowing for a nuanced understanding of market dynamics and consumer preferences.

**Products Table**

Containing information on over 100,000 products, the products table delves into the specifics of each item sold, including product ID, description, section, category, subcategory, type, manufacturer or brand code, unit quantity count, and the unit of measure in 150,000 entries. This rich product metadata allows for:

- Detailed performance analysis at the product level, facilitating strategic decisions regarding product assortment, pricing, and promotions.
- Insight into product categorization, which supports marketing strategies and the optimization of store layout according to product sections and categories.
- Brand performance assessment, guiding supplier negotiations and brand partnerships.

- Initial Observations: Share initial insights or interesting findings from a preliminary examination of the data. Discuss the quality of the data, including any potential issues like missing values or inconsistencies.
- Data Exploration Plan: Outline the steps for thorough data exploration, including statistical analyses, data visualization, and any specific techniques to uncover patterns or anomalies in the data.

**Data Prepossessing**

**Data Preprocessing Overview**

The data preprocessing stage for ACSE Supermarket's recommender system involved a meticulous approach to ensure the data's utility and accuracy for our analysis. Initially, we sampled transactions from 10,000 customers who had demonstrated significant engagement with General Mills products but had not purchased from ACSE's private label. This selective sampling was crucial for focusing our efforts on potential brand switchers. To account for seasonal buying patterns, which can significantly influence consumer purchasing behavior, we applied advanced feature engineering techniques during the modeling phase. Specifically, we incorporated 'day of the week' and 'month' as features to mitigate seasonal effects. The dataset was split temporally, with data from January to September serving as the training set, and data from October to December used for testing, aligning our evaluation period with critical sales quarters.

**Integration and Correction of Data Sources**

In the process of merging transactional and product data, we honed in on specific instances where our domain knowledge notably augmented the NLP findings, leading to precise corrections in product labels. For example, Product IDs 20640707002 and 20640707004, originally tagged under General Mills (GENM), were identified as Food Should Taste Good (FSTG) products through careful examination of product descriptions and corroborative domain

insights. Similarly, Product IDs 20313716 and 20318643, initially mislabeled as Nestle (NSTL) and Frito-Lay (FRTC) products, were correctly reclassified as General Mills items.

These particular ID-based corrections were just a fraction of the targeted adjustments made. Each adjustment was a result of a deliberate and methodical review process, involving both algorithmic NLP techniques and human expertise. This dual-faceted approach proved invaluable in standardizing product information and enhancing the overall integrity of our dataset, which, in turn, substantially benefited the modeling phase by aligning the data more closely with true customer and product interactions. Such meticulous data refinement practices underscore our commitment to data precision and the pursuit of a robust analytical foundation for ACSE Supermarket's recommender system.

**Enhancing Data Quality for Robust Analysis**

The refined data not only improved the accuracy of our analysis but also enhanced the reliability of the insights derived from it. By addressing data inconsistencies and standardizing product information, we ensured that our predictive models would be based on clean and relevant data, minimizing potential biases and errors in recommendation outputs. This stage of preprocessing was pivotal in preparing the dataset for complex analytical tasks, including predictive modeling and customer segmentation. The efforts to cleanse and integrate the data effectively set the stage for deploying a recommender system that could adeptly identify and target potential Kellogg's customers, thereby maximizing the impact of ACSE's marketing strategies.

## Modeling

**Modeling Overview**

In the modeling phase of the recommender system for ACSE Supermarket, we initiated our analysis by employing association rules to identify strong predictors of brand switching from General Mills (GM) to Kellogg (K). We first calculated five probability variations of GM and K products. Moreover, we divided customers into four groups: those who only bought K, both GM

and K, only GM, and neither K nor GM. Through this method, we analyzed the confidence levels of customers purchasing k subcategories after purchasing GM products. We identified six promising subcategories but selected the top two based on a statistical cutoff where the probability exceeded 0.3, with a noticeable gap between the top three and the remaining subcategories. The decision to exclude the third-ranked subcategory, 'breakfast', was due to its lack of alignment with GM's product range, thus focusing our efforts on the most relevant and potentially successful matches.

Further refining our approach, we performed advanced feature engineering on these two chosen subcategories, developing six key features to enhance the predictive accuracy of our models. These features included spending on subcategories, frequency of purchases within those subcategories, and a customer value classification based on purchasing behavior (cherry pickers, full value customers), which we categorized into high or low based on the distribution's long tail and the application of the 80/20 rule (Pareto principle). This categorization was crucial in understanding customer segments that might respond positively to targeted promotions involving Kellogg products while also considering the minor price variance between GM and Kellogg items, which supports the accuracy of our customer value assessment.

**Added Features**

**High_spending_longtail:**

- Definition: This feature aims to identify customers whose total spending in a specific subcategory is in the top 80%.
- Calculation Method: First, calculate each customer's total spending within the specific subcategory, then determine the 20th percentile of this spending distribution (i.e., customers below this value belong to the bottom 20%). Any spending exceeding this 20th percentile will result in the customer being marked with a 1 for this feature, indicating that these customers belong to the long tail of spending.

**High_frequency_longtail:**

- Definition: This feature is used to identify customers whose purchase frequency in a specific subcategory is in the top 80%.

- Calculation Method: First, calculate the number of purchases each customer makes within the specific subcategory, then determine the 20th percentile of this frequency distribution. Any customer whose purchase frequency exceeds this percentile will be marked with a 1, indicating that these customers belong to the long tail of purchasing frequency.

**No_Full_value_customer:**

- Definition: This feature is designed to identify low-risk and high-spending customers. It incorporates two dimensions: spending variance and average spending.

- Calculation Method: First, calculate the spending variance for each customer within the specific subcategory and find the 20th percentile of this distribution. Also, calculate the average spending amount and its standard deviation for each customer. If a customer's spending variance is in the bottom 20% and their average spending exceeds the mean by one standard deviation, they will be marked with a 0 for this feature. If these conditions are not met, they will be marked with a 1, indicating that these customers are not full-value customers, meaning they do not buy products solely based on price.

**Predictive Modeling and Evaluation**

With these refined features, we developed three predictive models: content-based filtering, collaborative filtering, and a naive Bayesian model, each designed to offer personalized product recommendations to individual customers. These models provided a list of the top five recommended Kellogg products for each customer, which were then tested against actual purchases within the test dataset to measure the effectiveness of our recommender system. The

accuracy of the models was evaluated using two metrics: the average hit rate, calculated as the number of hits divided by five (the number of recommendations per customer), and the overall success rate, which is the sum of successful recommendations divided by the total number of transactions. This dual-metric evaluation approach allowed us to comprehensively assess the performance of each model, ensuring that the recommender system could reliably predict and enhance customer transitions from General Mills to Kellogg, thereby aiming to boost Kellogg's market penetration and customer loyalty in a highly competitive segment.

| | Results | | |
|---|---|---|---|
| | Recommendation Accuracy | Transaction Hit Rate | Customer Purchase Rate |
| Naïve Bayes | 21.24% | **1.68%** | **72.64%** |
| Content-Based | 20.77% | 0.42% | 57.39% |
| CF | **25.36%** | 1.23% | 65.14% |

**Content-Based**

**Methodology and Design**:

The content-based filtering model deployed for ACSE Supermarket relies on the principle of item attributes to make recommendations. Unlike collaborative filtering, which utilizes user-user similarities, this model focuses on item-user relationships by examining the items that each user has engaged with. In the context of ACSE, the model was designed to analyze individual consumer purchase histories to identify the specific attributes of products purchased—such as brand, product subcategory, sales amount, and product specifications.

To design this model, we examined the cleaned transactional data to extract product attributes and customer interactions. Each product's metadata was leveraged to create a profile of interest for every customer, based on the products they had historically purchased. The model then used these profiles to generate predictions for each customer, recommending new items that shared similar attributes to those they had bought before.

The modeling process begins with preparing a unique text representation for each product by combining multiple features such as product subcategory, sales amount, unit value, and specific dummy variables indicating purchasing behavior patterns. These features are merged into a single string representation for each product, providing a comprehensive view of each product's characteristics.

The combined text strings are then transformed into a numerical format using the Term Frequency-Inverse Document Frequency (TFIDF) vectorization. This method emphasizes the importance of unique terms across products while penalizing common terms, helping to differentiate products based on distinctive features. The resulting sparse matrix encapsulates the textual information in a form suitable for calculating similarities.

Cosine similarity is used to measure the similarity between products. It computes the cosine of the angle between two vectors in the TFIDF vector space, effectively assessing how closely related two products are based on their textual features. For each product, similarities with all other products are calculated, forming a similarity matrix that serves as the basis for making recommendations.

**Processing Challenges**:

Managing Large Data Volumes: Given the extensive data from ACSE's millions of transactions and a vast number of products, computational efficiency is a prime concern. The code employs a TfidfVectorizer for feature transformation, which is inherently memory-efficient as it constructs a sparse matrix representation of the features. This approach minimizes memory consumption compared to dense matrix representations, crucial when dealing with large-scale data. Moreover, cosine similarity, calculated using the cosine_similarity function, is also efficient for sparse matrices, optimizing the computational load when calculating similarities between all pairs of products.

**Interpretation of Outcomes**:

The model's outcome, as indicated in the results provided, shows a Customer Recommendation Accuracy of 0.2077. This suggests that when recommending products, the content-based model correctly predicted 20.77% of the items that customers ended up purchasing. The Transaction Hit Rate of 0.0042 reflects the model's precision in predicting the exact transactional occurrence of a product purchase. The relatively low Transaction Hit Rate could be attributed to the individual nature of content-based filtering, which may not capture the broader trends within the customer base that collaborative efforts might.

The Customer Purchase Rate of 0.5739 indicates that over half of the customers made at least one purchase from the top five recommendations. This suggests that while the model has room for improvement, it successfully captured customer preferences to a reasonable degree, making it a valuable tool for targeted marketing campaigns. The results underscore the potential for content-based models to enhance personalization in marketing strategies, creating opportunities for ACSE to deepen customer engagement by recommending products more closely aligned with individual preferences.

## Collaborative Filtering

**Methodology and Design**:

Our collaborative filtering model is predicated on constructing a nuanced map of user-item interactions. This map, materializing as a sparse matrix, marks each customer's purchase history in terms of both product units and transactional value—characteristics that play a pivotal role in reflecting true consumer preference. Our methodology embraces the complexity of customer behavior, leveraging nuanced metrics like sales quantity, amount, as well as our self–constructed features (e.g. high_spending_longtail). This granularity allows for a multidimensional perspective of user interactions, vital for the refined operation of our model.

Prior to populating the interaction matrix, a rigorous normalization process is executed to address the issue of scale disparity across the transactional data. Employing *MinMaxScaler*, we

calibrate the data such that each value falls within a 0 to 1 range, preserving the original distribution while eliminating scale-induced bias. The normalization extends beyond continuous variables, encapsulating binary customer behavior features that inform on purchasing tendencies within specific product categories. This comprehensive normalization ensures a level playing field where each variable, irrespective of its original magnitude, contributes fairly to the outcome of the model.

Our collaborative filtering model is designed to be inherently scalable, with a focus on the stability of item-item relationships over the more volatile user-user interactions. By emphasizing similarities between products, we mitigate the need for frequent model recalibration, thereby reducing computational overhead. This item-based approach proves advantageous in handling ACSE's extensive and varied product catalog, allowing for an efficient scaling of the model in proportion to our growing customer base and inventory. The model's architecture is fine-tuned to ensure that the recommendations it generates are diverse and personalized, mirroring the unique purchasing patterns of each customer while also introducing them to new products that align with their established preferences.

**Processing Challenges**:

A significant challenge encountered during the deployment of our collaborative filtering model was the intricacy of ACSE Supermarket's transactional dataset. With a vast array of products and complex purchasing patterns, the construction of an accurate item-user interaction matrix was a sophisticated task. The matrix serves as the foundation of our model, and its precision is crucial for generating reliable recommendations. Given the specificity of our target—General Mills customers—our dataset was both rich and unwieldy, with high-dimensional data that required careful parsing and manipulation.

The challenge was further compounded by the need to account for the variability in sales quantities and amounts, where some transactions reflected bulk purchases while others were more modest. Normalizing these disparate transactional features to a uniform scale was essential to avoid skewed recommendations and to reflect genuine customer preferences accurately.

Moreover, handling the sparse nature of the interaction data, where numerous customers might have transactions with a limited subset of items, presented another layer of complexity. This sparsity necessitated advanced computational techniques to effectively fill in the gaps and draw inferences about customer preferences.

The variability and sparsity of the data posed significant computational demands. Optimizing the efficiency of the matrix generation process was imperative to ensure that the collaborative filtering model could operate at scale without compromising on performance. Achieving this balance between accuracy and computational tractability was a key focus of our methodological refinement. The process of transforming the raw transactional data into a form suitable for collaborative filtering—while maintaining its integrity and the nuances of customer behavior—required rigorous data engineering efforts and was central to our challenge in processing the data.

**Interpretation of Outcomes**:

The model exhibits a Customer Recommendation Accuracy of 0.2536, which conveys that approximately 25.36% of the items recommended by the model matched the items that customers actually purchased. This level of accuracy, while indicative of the model's predictive capabilities, also highlights areas where refinement could enhance performance.

The Transaction Hit Rate stands at 0.0123, signaling a modest precision in the model's ability to forecast specific product purchases within transactions. The modesty of this rate may stem from the collaborative nature of the filtering system, which, while robust in discerning shared customer behaviors, can sometimes miss the mark on individual purchase events due to the complex interplay of factors influencing a customer's decision.

Notably, the Customer Purchase Rate registers at 0.6514, underscoring that a majority of customers—over 65%—engaged with the model's top five recommendations to make at least one purchase. This rate is particularly telling of the model's efficacy, as it demonstrates a strong

alignment with customer preferences and suggests a significant potential for influencing purchasing decisions.

**Naive Bayes**

**Methodology and Design**:

The Naive Bayes model is a probabilistic approach that applies Bayes' theorem with the assumption of independence between predictors. For ACSE Supermarket, this model was designed to predict the likelihood that a customer would purchase a particular Kellogg's product based on their previous purchases from General Mills. This method is particularly well-suited for classification tasks where the dimensionality of the input data is high, as in the case of ACSE's extensive product catalog.

To build this model, we calculated the conditional probability of a customer buying a Kellogg's product, given their purchase history with General Mills products. The transactional data provided the necessary information to compute these probabilities. This approach allowed us to predict the top five products a customer was most likely to buy, given their profile.

**Processing Challenges**:

The independence assumption of Naive Bayes can be a challenge, especially in complex datasets where the buying patterns might be influenced by multiple interdependent factors. For example, a customer's decision to purchase a cereal may be dependent on the purchase of milk—a correlation that Naive Bayes may overlook due to its feature independence assumption.

Another challenge lies in handling continuous data, as Naive Bayes inherently assumes categorical data. To address this, we employed binning methods to convert continuous features, like sales amount and quantity, into categorical ranges, which could then be used effectively in the model.

**Interpretation of Outcomes**:

Looking at the results, the Naive Bayes model achieved a Customer Recommendation Accuracy of 0.2124, which indicates that approximately 21.24% of the recommended products matched the customers' actual purchases. This metric shows the model's strength in understanding customer preferences.

The Transaction Hit Rate at 0.0168 is indicative of the model's ability to capture the correct transactions within the recommended product sets. While not as high as might be desired, it demonstrates the model's utility in predicting purchase behavior at a granular level. Most notably, the Customer Purchase Rate of 0.7264 is the standout metric, as it implies that nearly 72.64% of customers made a purchase from the top recommendations provided by the model.

This high Customer Purchase Rate validates the Naive Bayes model as a strong contender for identifying potential cross-sell opportunities among customers. Its ability to interpret customer transactions and predict future behavior can be strategically utilized to enhance ACSE's personalized marketing efforts, driving both sales and customer satisfaction. The model's success in this area suggests that with further refinement, it could become an even more potent tool for influencing customer purchasing decisions and guiding promotional strategies.

## Revenue Expectation

The implementation of our recommender system (RS) within ACSE Supermarket's General Mills (GM) product transactions promises a significant uplift in revenue. The revenue expectation is calculated based on the sum of transactions within the targeted subcategories of GM products. By applying a hit rate of 1.68% to the total sales amount of these transactions, we have projected an additional annual revenue stream. The computation took into account customers characterized by all three features as '1'—indicating high frequency and spending, but not exclusively looking for full value—which represent the core target for the RS initiative. Our

model estimates that approximately 55% of the transactions within the subcategories are likely to be influenced by the RS, as per the average rate derived from the analysis of transaction data from 2017 to 2019. The estimated annual revenue from this cohort, using our RS, is calculated to be around $3,270,555, which is the aggregated sales amount adjusted by the hit rate and the proportion of transactions attributable to the target customer segment, averaged over a three-year period.

## Conclusion

As we reach the conclusion of our analytical journey with ACSE Supermarket, it's essential to reflect on the sophisticated methods we employed to create a recommender system capable of fostering a strategic marketing shift—enticing loyal General Mills customers to explore and purchase Kellogg's products. The system was articulated through the adept use of three distinct predictive models: content-based filtering, collaborative filtering, and a naive Bayesian model.

Upon a comparative analysis, each model demonstrated particular strengths. The collaborative filtering (CF) model stood out with the highest Customer Recommendation Accuracy (0.2536), indicating its superior ability to match recommendations with customers' eventual purchases. This model benefits from leveraging user behavior patterns, drawing on the wisdom of the crowd to predict individual preferences. The content-based model followed closely, with a slightly lower accuracy but was notably effective in considering item attributes, providing a nuanced understanding of individual customer preferences.

In contrast, the naive Bayesian model, despite its relatively lower recommendation accuracy, achieved the highest Customer Purchase Rate (0.7264). This suggests that the Naive Bayes approach was most effective in influencing actual customer purchasing decisions, perhaps due to its probabilistic foundations and the robust interpretation of customer transaction history.

However, when considering the Transaction Hit Rate, all models showed room for improvement, with the collaborative filtering model outperforming the others, albeit modestly. This metric indicates the potential precision with which these models can predict specific purchases, an area where refinement could yield substantial benefits.
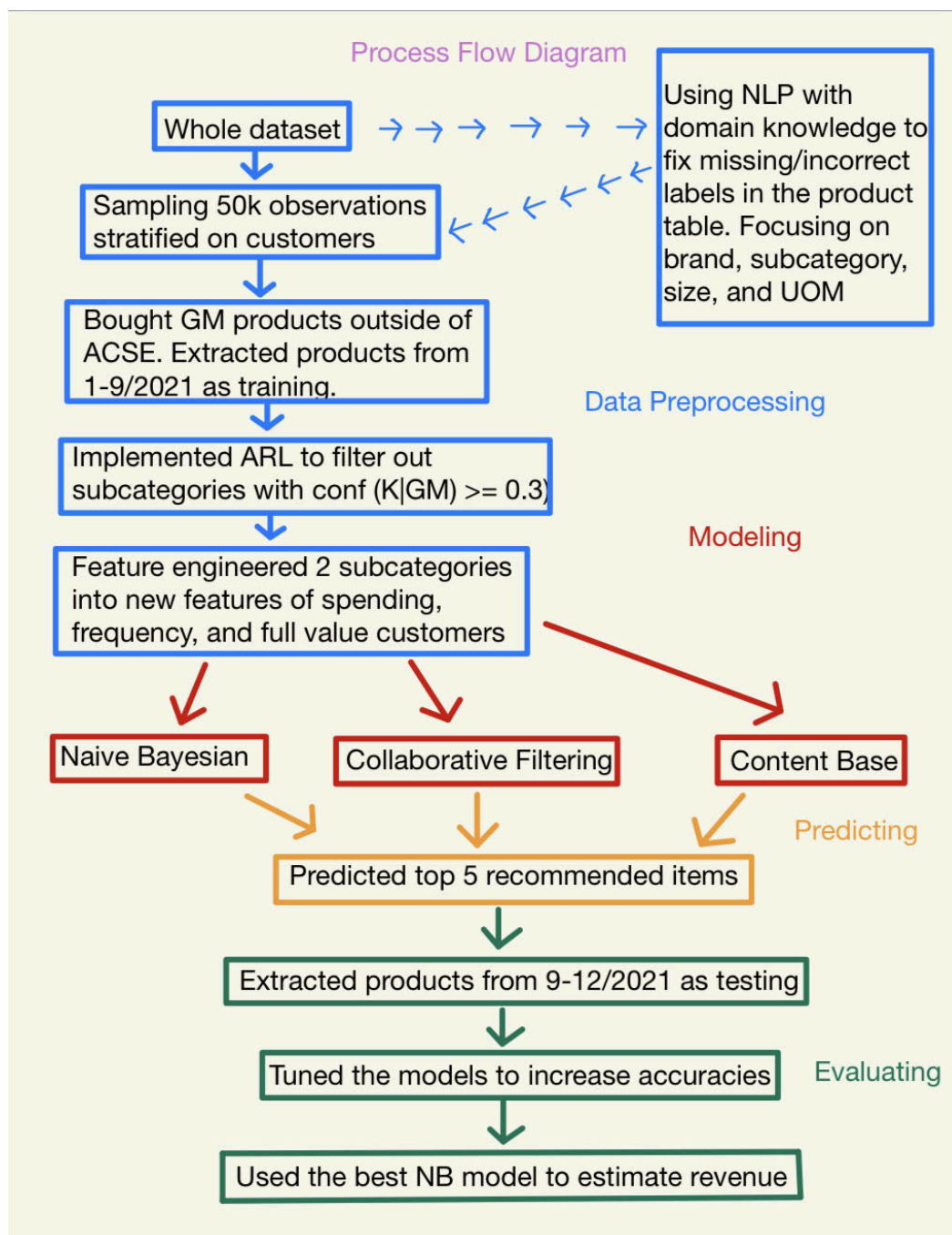
In synthesis, while the collaborative filtering model showed a balanced performance across metrics, the naive Bayesian model was particularly strong in driving purchases, and the content-based model offered valuable insights into individual preferences. Each model's unique attributes suggest that an ensemble approach, combining the strengths of each method, could further enhance the recommender system's efficacy.

This comprehensive analysis serves not only to conclude the current project but also to set a benchmark for future initiatives. By continually refining these models and exploring hybrid approaches, ACSE Supermarket can elevate its personalized marketing efforts, achieving a fine balance between customer satisfaction, increased sales, and strategic promotion of products, all while charting a course for sustained innovation in retail marketing strategies.
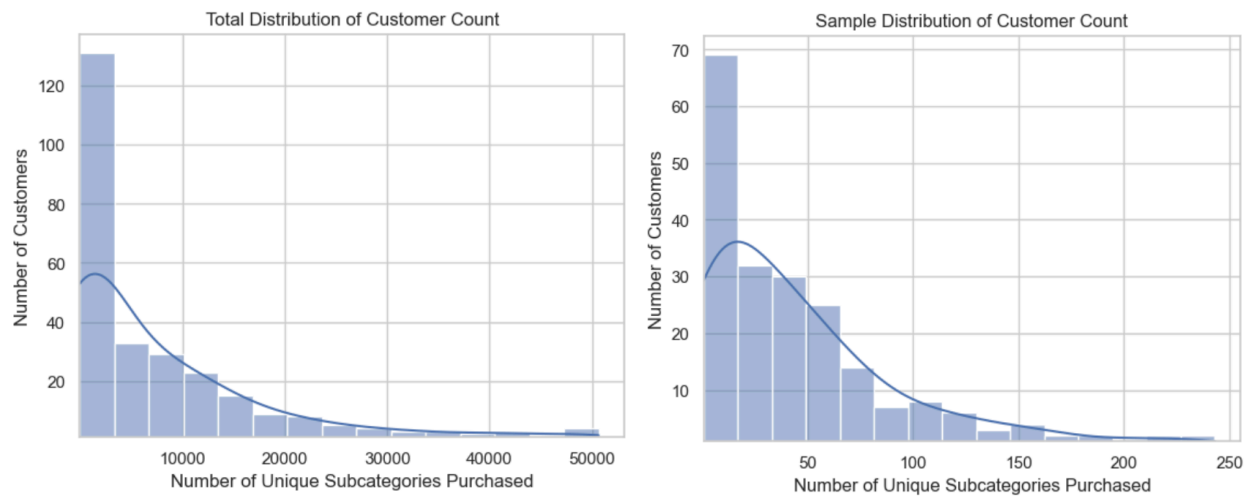
## Recommendations

Based on the comparative analysis of the recommender system models, we recommend ACSE Supermarket adopt an integrated approach that combines the strengths of collaborative filtering and naive Bayesian models. This hybrid model should leverage the high customer purchase rate of the Naive Bayes model and the superior accuracy and transaction hit rate of the collaborative filtering model. By integrating these models, ACSE can achieve a more precise and effective personalization strategy, enhancing customer engagement and increasing conversion rates for Kellogg's product promotions. This strategic implementation should be accompanied by continuous model tuning and evaluation to ensure sustained improvement in recommendation performance.

**Appendices**



Process Flow Diagram

Whole dataset → → → → → →

Using NLP with domain knowledge to fix missing/incorrect labels in the product table. Focusing on brand, subcategory, size, and UOM

Sampling 50k observations stratified on customers

← ← ← ← ← ← ←

Bought GM products outside of ACSE. Extracted products from 1-9/2021 as training.

Data Preprocessing

Implemented ARL to filter out subcategories with conf (K|GM) >= 0.3)

Modeling

Feature engineered 2 subcategories into new features of spending, frequency, and full value customers

Naive Bayesian          Collaborative Filtering          Content Base

Predicting

Predicted top 5 recommended items

Extracted products from 9-12/2021 as testing

Tuned the models to increase accuracies          Evaluating
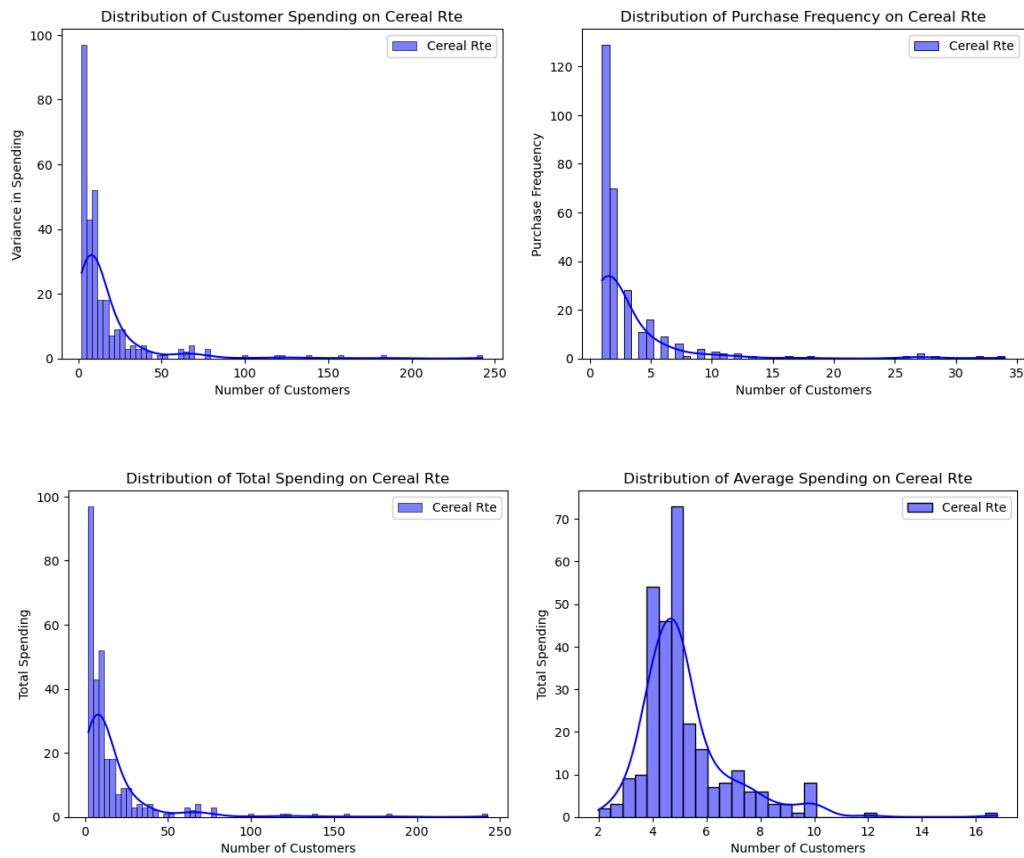
Used the best NB model to estimate revenue

Mind Map

The distributions show the number of unique subcategories each customer has purchased and then aggregate the number of customers for each subcategory count between the entire data set and sample data set. We observed the same distribution pattern in both the entire dataset and a sample, suggesting that the sample is representative of the whole. This consistency indicates reliable customer behavior across the dataset, ensuring that statistical analyses and predictions based on the sample are likely to be valid for the entire population.
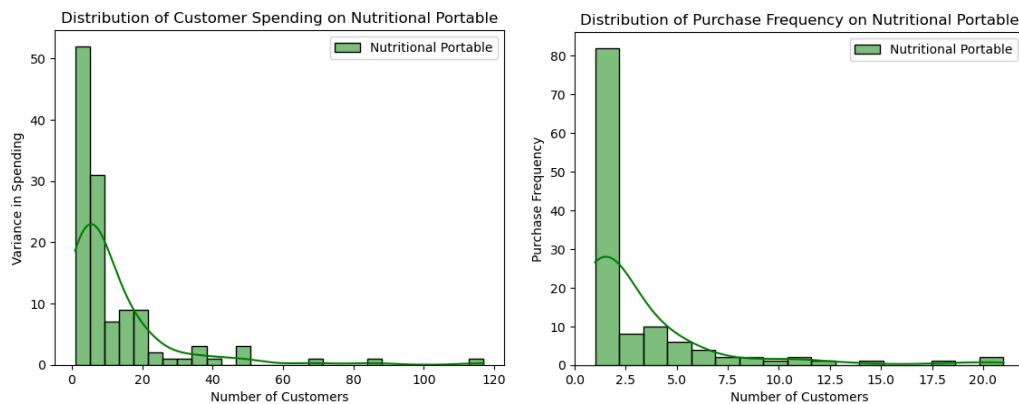


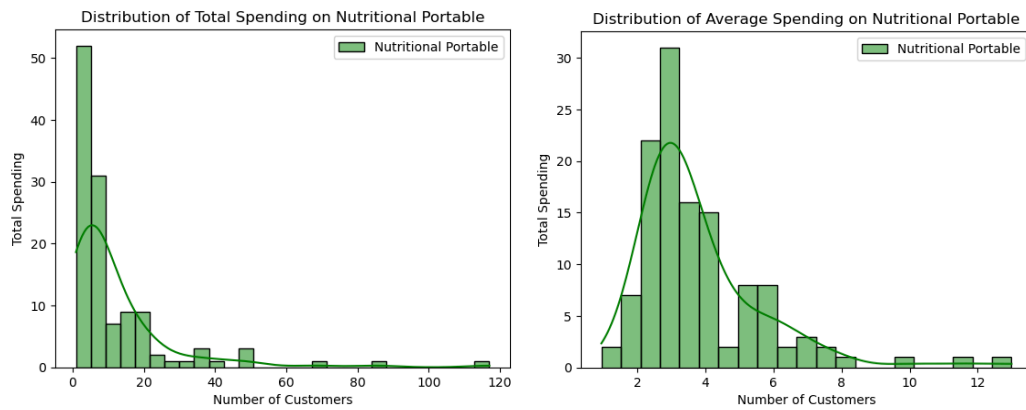Confidence Level of Kellogg Subcategory Purchases by General Mills Customers Using Association Rules

| Kellogg_subcategory | Confidence |
|---|---|
| Cereal Rte | 55.50% |
| Nutritional Portable | 29.58% |
| Breakfast | 27.04% |
| Crackers/Health Cake | 4.81% |
| Halloween | 1.14% |
| Coating Mixes | 0.71% |

The distribution of 3 engineered features of "Cereal Rte" "Nutritional Portable" subcategory

The distribution of 3 engineered features of "Nutritional Portable" subcategory

## Distribution of sales statistics