

Essays on the use of probabilistic machine learning for estimating
customer preferences with limited information

Nicolas Padilla

Submitted in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy
under the Executive Committee
of the Graduate School of Arts and Sciences

COLUMBIA UNIVERSITY

2021

©2021
Nicolas Padilla
All Rights Reserved

Abstract

Essays on the use of probabilistic machine learning for estimating customer preferences with limited information

Nicolas Padilla

In this thesis, I explore in two essays how to augment thin historical purchase data with other sources of information using Bayesian and probabilistic machine learning frameworks to better infer customers' preferences and their future behavior. In the first essay, I posit that firms can better manage recently-acquired customers by using the information from acquisition to inform future demand preferences for those customers. I develop a probabilistic machine learning model based on deep exponential families to relate multiple acquisition characteristics with individual level demand parameters, and I show that the model is able to capture flexibly non-linear relationships between acquisition behaviors and demand parameters. I estimate the model using data from a retail context and show that firms can better identify which new customers are the most valuable. In the second essay, I explore how to combine the information collected through the customer journey — search queries, clicks and purchases; both within-journeys and across journeys — to infer the customer's preferences and likelihood of buying, in settings in which there is thin purchase history and where preferences might change from one purchase journey to another. I propose a non-parametric Bayesian model that combines these different sources of information and accounts for what I call *context heterogeneity*, which are journey-specific preferences that depend on the context of the specific journey. I apply the model in the context of airline ticket purchases using data from one of the largest travel search websites

and show that the model is able to accurately infer preferences and predict choice in an environment characterized by very thin historical data. I find strong context heterogeneity across journeys, reinforcing the idea that treating all journeys as stemming from the same set of preferences may lead to erroneous inferences.

Table of Contents

List of Figures	iv
List of Tables	vi
Acknowledgments	vii
Introduction	1
Chapter 1: Overcoming the Cold Start Problem of CRM using a Probabilistic Machine	
Learning Approach	5
1.1 Introduction	8
1.2 Previous literature	14
1.3 The “cold start” problem of CRM	17
1.3.1 The “cold start” problem	17
1.3.2 Augmenting cold start data with acquisition characteristics	19
1.3.3 Predictive power of augmented data	21
1.3.4 Modeling challenges	24
1.4 Modeling framework	27
1.4.1 Model development	27
1.4.2 Estimation and identification	40
1.4.3 Model inferences for newly acquired customers	41
1.4.4 Model performance	43
1.5 Empirical application	48
1.5.1 Data and model specification	48
1.5.2 Estimation	57
1.5.3 Results	59
1.5.4 Overcoming the cold start problem	61
1.6 Conclusion	68
Chapter 2: The Customer Journey as a Source of Information	72
2.1 Introduction	75
2.2 Relevant literature	80
2.3 Empirical setting	84
2.3.1 The customer purchase journey of airline tickets	84
2.3.2 Extracting information from the data	90
2.3.3 Inferring preferences from purchase journey data	98
2.4 Model	100
2.4.1 Model intuition	101
2.4.2 Model development	106

2.4.3	Specification of query variables, covariates, and sensitivities	120
2.4.4	Model estimation and prediction for partially observed journeys	122
2.5	Results	123
2.5.1	Model estimates	123
2.5.2	Contexts in the data	125
2.5.3	Prediction of purchase incidence and product choice	132
2.5.4	Illustration of how the model infers contexts and preferences along the journey	137
2.6	Conclusion and discussion	142
References		145
Appendix A: Appendix to Essay 1 - Overcoming the Cold Start Problem of CRM using a Probabilistic Machine Learning Approach		
		152
A.1	Augmenting the acquisition characteristics via product embeddings	152
A.1.1	Data processing	152
A.1.2	Word2vec algorithm	153
A.1.3	Interpreting the product dimensions	155
A.1.4	Product mapping for first purchase data	156
A.2	Brief description of DEFs	157
A.3	Model priors and automatic relevance determination component	159
A.3.1	Automatic relevance determination	159
A.3.2	Model priors	160
A.4	Rotation of traits	162
A.5	Algorithm for newly-acquired customers	165
A.6	Further details about the simulation analyses	166
A.6.1	Simulation design	166
A.6.2	Data generation process	167
A.6.3	Estimated models	173
A.6.4	Assessing model performance	177
A.6.5	Interpreting the model parameters and results	181
A.6.6	Why is the model giving superior performance?	185
A.6.7	Exploring the number of dimensions per layer	186
A.6.8	Model performance “at scale”	191
A.7	Empirical application: Additional results	198
A.7.1	Possible sources of endogeneity in the model components	198
A.7.2	Latent attrition benchmarks models	200
A.7.3	Interpreting the latent traits	201
A.7.4	FIM predictive accuracy using in-sample customers	204
A.7.5	Population distribution and individual-level posterior distributions	205
A.7.6	Exploring the latent factors	206
A.7.7	Details on the (Machine Learning) benchmark models	208
Appendix B: Appendix to Essay 2 - The Customer Journey as a Source of Information 211		
B.1	Model priors	211

B.2	Blocked-Gibbs sampler algorithm	212
B.3	Parameter estimates per context	219
B.3.1	Query context location parameters	219
B.3.2	Click and purchase context location parameters	220

List of Figures

1.1	Transactions versus acquisition characteristics	22
1.2	Transactions versus interaction of acquisition characteristics	22
1.3	Graphical model of first impressions	39
1.4	Model performance for Scenario 3	47
1.5	Acquisition characteristics for customers with top/middle/low CLV.	64
1.6	Acquisition characteristics for customers with top/middle/low sensitivity to Email.	65
1.7	Acquisition characteristics for customers with top/middle/low sensitivity to DM.	66
1.8	Demand parameters vs. binary acquisition characteristics.	67
1.9	Demand parameters vs. continuous acquisition characteristics.	71
2.1	Flow of the customer purchase journey for roundtrip flights	85
2.2	Mock-up of purchase journey steps	88
2.3	The data generating process	102
2.4	Model intuition	104
2.5	DAG of customer journey model	107
2.6	Expected number of clusters from a Dirichlet Process vs. a Pitman-Yor process.	118
2.7	Example of a context distribution drawn from a Pitman-Yor process prior	119
2.8	Posterior mean and 95% CPI of contexts probabilities, π_c	125
2.9	Number of contexts per customer.	126
2.10	Posterior mean of context location parameters θ_c	128
2.11	Top 50 routes per context	131
2.12	Queries of two holdout journeys from the same customer	138
2.13	Posterior of context for each journey example.	139
2.14	Posterior of price coefficient for each journey example.	141
A1	Model selection for Word2vec: Perplexity when varying the number of dimensions from 2 to 10.	154
A2	Visual representation of the product embeddings	157
A3	Visualization of the benchmark models	177
A4	In sample individual posterior mean vs. true intercepts of the demand model.	178
A5	Out of sample individual posterior mean vs. true intercepts of the demand model.	180
A6	Posterior distribution of α	182
A7	Posterior mean of α as a function of number of dimensions in lower layer and upper layer.	188
A8	Posterior distribution of pseudo- α^1 (Linear scenario).	189

A9	Posterior distribution of pseudo- α^1 (Interactions scenario).	191
A10	Square correlation between simulated and predicted β for Covariate 1 in Scenario 2: Interaction	192
A11	Population distribution and individual-level posterior distribution for customers in the <i>Test</i> sample.	207
A12	Posterior distribution of α	208
A13	Convergence of α	209
A14	Posterior distribution of pseudo- α^1	210

List of Tables

1.1	Accuracy of predictions of demand parameters for (out-of-sample) customers	45
1.2	Summary of time-varying marketing actions.	51
1.3	Summary statistics of selected acquisition characteristics.	55
1.4	Correlations among selected acquisition characteristics.	55
1.5	Parameter estimates of FIM.	59
1.6	Comparison with benchmark models (<i>Validation</i> sample).	61
1.7	Identifying valuable customers using <i>Test</i> customers.	62
2.1	Summary statistics of query variables	91
2.2	Summary statistics of product attributes in page results	96
2.3	Data summaries, per customer and per journey.	99
2.4	Parameter estimates of click and purchase models. We show the average across customers and contexts ($\bar{\mathbf{b}}$) and standard deviation across customers (σ_b).	124
2.5	Estimated models	133
2.6	AUC for purchase incidence using each piece of information from the customer journey.	134
2.7	Hitrate and RMSE for product choice per attribute using each piece of information from the customer journey.	136
A1	Top 5 products per dimension of the product embeddings.	156
A2	True values for factors f_{i1} and f_{i2} impact on acquisition parameters (B_{1p} and B_{2p}).	168
A3	Simulated values for ω_k^1 in the Linear scenario	169
A4	Simulated values for ω_k^1 and Ω_k^2 in the Quadratic/Interaction scenario	171
A5	Simulated values for ω_k^1 in the Positive part scenario	172
A6	Posterior mean of lower layer weights (\mathbf{W}^y and \mathbf{W}^a) for FIM.	183
A7	True associated effects of factors on demand and acquisition variables.	184
A8	Squared correlation (true vs predicted) for Covariate 1; Quadratic/Interaction Scenario.	186
A9	Model at scale results	197
A10	Latent attrition benchmarks models.	201
A11	Posterior mean of correlations across customers of individual lower level traits \mathbf{z}_i^1	202
A12	Rotated traits weights' on acquisition and demand variables	203
A13	Model fit and prediction accuracy for the <i>Training</i> sample	205
B1	Posterior mean of query location parameters per context	219
B2	Posterior mean of location click and purchase parameters	220

Acknowledgments

Going through my doctoral studies would not have been possible without the unconditional love and support of my wife Jacqueline. I would like to thank her for her patience and encouragement during all these wonderful years in New York.

During my time at Columbia, I met brilliant people that helped me during the first steps of my academic journey. First and foremost, I would like to thank my advisors Oded Netzer and Eva Ascarza from whom I had learned tremendously as a researcher and as a person. I could not have been luckier to be their student, and I will always aspire to guide students in the future the way they have guided me. Second, I would also like to thank the rest of my committee, Asim Ansari, Vicki Morwitz, and Eric Schwartz, who have provided me with invaluable feedback for the research contained in this dissertation. Third, I would also like to thank my friends in the doctoral program: Malek Ben Sliman, Alain Lemaire, Ma'ayan Malter, Ryan Dew, Verena Schoenmueller; I will always cherish all the conversations we had about life, research, Real Madrid, and many other fun topics over a good cup of coffee.

I am also particularly thankful to Ricardo Montoya, my advisor in Chile, who encourage me to start a career in academia and guided me during my undergraduate and

post-graduate days at Universidad de Chile. I also thank Marcelo Olivares, Marcel Goic, Richard Weber, and Andrés Musalem, for their contribution in my formation years in Chile.

I would also like to acknowledge the Marketing Science Institute for the financial support to develop this dissertation.

Many other friends have helped me through these years in New York and Santiago. I would like to thank especially David Aranda, Francisco Liebbe, Isidora Valdés, Gabriel Guggisberg, Fernanda Abarzúa, Cristian Urbina, Miguel Biron, Laura Torres, Francisco Castro, and Mauro Escobar.

Finally, completing this doctoral program would not have been possible without the warm support of my parents, Begoña Pérez and Luis Padilla; my siblings, Felipe and Begoña; and my extended family, Ruth, Sergio, Pablo, Daniela, and Constanza.

To my beloved wife, Jacqueline.

Introduction

In the data-rich environment, firms and researcher aim at inferring customer preferences from their history of past purchases, to predict whether customers will buy again, what product they will buy, and how they may respond to marketing actions. Firms are often pressured to understand customers at the time they make decisions, right after they are acquired, or when they are still interacting with the firm at a time where such decisions are most effective. However, traditional approaches to understand customer heterogeneous preferences often rely on long history of past purchases. There are many reasons why long-history of purchases by consumers may not be available. First, many product categories may have a long product cycle (e.g. cars, mortgage), or the customer may purchase very infrequently in the category (e.g. flights, hotel stays). Second, the firm may be particularly interested in understanding specific customers with short purchase history, because they have recently purchased for the first time, or the firm may want to understand the needs of an infrequent customer at the moment of interaction. Third, rising concerns regarding consumer privacy has resulted, and may continue to result, in regulatory changes that limit firms' ability to store long historical data at the individual level.

The solution to thin history about customer purchases may lie in different sources of data beyond purchases that can be collected by the firm. Firms do not only store whether a customer transacts with the firm, but they can also register several other relevant pieces of information on the interaction between the firm and the customer.

In this dissertation I explore in two essays how to incorporate these other sources of information to better predict customers' future behavior. In both contexts, I employ a similar methodological strategy. I assume that these different sources of information are outcomes from individual parameters that correlate with those of interest. Thus, even if these sources of information may be different in nature from purchase outcomes, they carry valuable information about the customer's underlying preferences, and are hence useful to better predict purchase outcomes.

In the first essay "Overcoming the Cold Start Problem of CRM using a Probabilistic Machine Learning Approach" I posit that firms can better manage recently-acquired customers by using the information from acquisition to inform future demand preferences for those customers. In this essay, I develop a probabilistic machine learning model based on Deep Exponential Families to relate multiple acquisition characteristics with individual level demand parameters, and I show that the model is able to capture flexibly non-linear relationships between acquisition behaviors and demand parameters. I estimate the model using data from a retail context and show that firms can better identify which new customers are the most valuable.

In the second essay "The Customer Journey as a Source of Information," I explore how the information along the journey that the customer undertakes carries valuable information about the purchase that may take place. This information is particularly

valuable for high involvement purchases, such as flights, insurance, and hotel stays, where the firm observe at most only a handful of purchases during a customer lifetime. Moreover, customers in these industries often look for products that satisfy different needs depending on the context of the purchase (e.g., flights for a family vacation vs. flights for a business trip), further complicating the task to understand what a customer might prefer in the next purchase occasion. To overcome those challenges, I propose a non-parametric Bayesian model that combines different sources of information from the customer journey — search queries, clicks and purchases; both within-journeys and across journeys — to infer the customer’s preferences. The model accounts for what I call *context heterogeneity*, which are journey-specific preferences that depend on the context of the specific journey. I apply the model in the context of airline ticket purchases using data from one of the largest travel search websites and show that the model is able to accurately infer preferences and predict choice in an environment characterized by very thin historical data. I find strong context heterogeneity across journeys, reinforcing the idea that treating all journeys as stemming from the same set of preferences may lead to erroneous inferences.

Beyond the main substantive question, these two essays share a common methodological approach. Both models are developed using a flexible probabilistic framework and relate to the use of probabilistic machine learning and Bayesian methods marketing contexts. In the first essay I develop a deep probabilistic model of demand and acquisition characteristics where the individual-level parameters of each of these sub-models are projected into a lower-dimension space using a two-layered deep exponential family (DEF) component. This flexible component allows the model to capture potential non-linear relationships between these set of parameters, while reducing the dimensionality of a large

set of potentially correlated acquisition characteristics. In the second essay, I use Bayesian nonparametrics methods to uncover the number of (unobserved) purchase contexts that shift individual preferences for each customer journey differently. These methods used in the two essays of this dissertation illustrate how marketers can flexibly capture customer preferences and their future demand propensities by leveraging other sources of data.

Chapter 1

Overcoming the Cold Start Problem of CRM using a Probabilistic Machine Learning Approach

This essay forms the basis of a paper of the same name jointly authored with Eva Ascarza which is under third round review at the *Journal of Marketing Research*.

Abstract

The success of Customer Relationship Management (CRM) programs ultimately depends on the firm's ability to understand consumers' preferences and precisely capture how these preferences may differ across customers. Only by understanding customer heterogeneity, firms can tailor their activities towards the right customers, therefore increasing the value of customers while maximizing the return on the marketing efforts. However, identifying differences across customers is a very difficult task when firms attempt to manage new customers, for whom only the first purchase has been observed. For those customers, the lack of repeated observations poses a structural challenge to infer unobserved differences across them. This is what we call the "cold start" problem of CRM, whereby companies have difficulties leveraging existing data when they attempt to make inferences about customers at the beginning of their relationship.

In this research we propose a solution to the cold start problem by developing a modeling framework that leverages the information collected at the moment of acquisition. The main aspect of the model is that it flexibly captures latent dimensions that govern both the behaviors observed at acquisition as well as future propensities to buy and to respond to marketing actions. Using probabilistic machine learning, we combine deep exponential families with the demand model, relating behaviors observed in the first purchase with consequent customer behavior. The model can be integrated with a variety of demand specifications and is flexible enough to capture a wide range of heterogeneity structures (both linear and non-linear), thus being applicable to a variety of behaviors and contexts. We validate our approach in a retail context and illustrate how the focal firm can overcome

the cold start problem by augmenting the (thin) historical data for new customers using the firm’s transactional database and applying the proposed modeling framework to those data. We empirically demonstrate the model’s ability at identifying high-value customers as well as those most sensitive to marketing actions, right after their first purchase. Leveraging the model predictions, the firm can also identify the most relevant variables—transaction characteristics or products being purchased at the moment of acquisition—that are predictive of behaviors of interest (e.g., sensitivity to email communications).

Keywords: Customer Relationship Management (CRM), Deep Exponential Families, Probabilistic Machine Learning, Cold Start Problem.

1.1 Introduction

Customers are different, not only in their preferences for products and services, but also in the way they respond to marketing actions. Understanding customer heterogeneity is at the heart of Customer Relationship Management (CRM) programs—from obtaining accurate estimates of the value of current and future customers, to deciding which customers should be targeted in the next marketing campaign. Over the last three decades, the marketing literature has provided researchers and analysts with methods to empirically identify unobserved differences across customers using their past history—e.g., customers with higher versus lower expected lifetime value (e.g., Schmittlein et al., 1987; Fader et al., 2005, 2010), those who are less sensitive to a price increase (e.g., Rossi et al., 1996; Allenby and Rossi, 1998), or those who are more receptive to marketing communications (e.g., Ansari and Mela, 2003). However, when firms attempt to implement CRM programs on customers who have been acquired recently, they only observe these customers’ first purchase. This lack of repeated observations presents a structural challenge for estimating unobserved differences across recently-acquired customers, precluding firms from leveraging such heterogeneity. We call this the “cold start” problem of CRM; that is, the challenge that firms face when trying to make inferences about customers at the outset of the relationship, for whom data is limited.

Firms have traditionally relied on demographics (e.g., age, gender) and/or recency metrics (e.g., how many weeks since your last transaction) to target marketing efforts with limited data (Shaffer and Zhang, 1995). These approaches, however, face practical

limitations: Recency metrics, for example, do not differentiate among recently acquired customers (as they all were acquired at the same time), and relevant personal information is generally hard to collect or poses data privacy challenges. Although, thanks to technological advances, firms can now increasingly observe a wider range of behaviors on each customer touch. What in the past might have been considered simply a transaction added to a customer base is now a collection of behaviors that a customer incurs while making a first purchase (e.g., is the purchase online or offline, did they buy a new product or an old best-seller, did they buy on discount or at full price). While some of these characteristics may be purely coincidental with the moment in which the customer made their first purchase, others may carry important information as they reflect latent customer preferences/attitudes. Thus, whereas firms only observe a just-acquired customer in one occasion, they now have many more cues to form a “first impression” of who this customer is, which can be used to understand heterogeneity across recently acquired customers. We present a solution to the cold start problem that is flexible, scalable, and general. Specifically, we augment transactional data with information collected when a customer makes their first purchase—information already available in the firm’s database—and propose a probabilistic machine learning modeling framework that extracts information relevant to making inferences about the customer’s future behavior. The model, which we term the “First Impression Model” (FIM), reflects the premise that behaviors and choices observed in newly-acquired customers can be informative about underlying traits that are, in turn, predictive of their future behavior. We operationalize these customer traits via a finite set of latent factors that enable the model to reduce the dimensionality of, while extracting

relevant signals from, the data, and assume those traits to drive, at least partially, customer behaviors observed both at the moment of acquisition and in the future.

In essence, the FIM is a deep probabilistic model of demand (main outcome of interest to the firm) and acquisition characteristics (customer outcomes that are observed to the firm at the moment of acquisition) where the individual-level parameters of each of these sub-models are projected into a lower-dimension space using a two-layered deep exponential family (DEF) component. The lower layer of the DEF component captures the relevant correlations among the individual-level parameters. We incorporate automatic relevance determination priors (ARD) for this layer, enforcing sparsity and automatically reducing the dimensionality of the individual-level parameters, similarly as in a Bayesian PCA model and modern applications of “supervised” factor models. The model departs from the aforementioned models by allowing non-linear relationships among the factors in the lower layer, through the upper layer.

First among four notable aspects of the proposed modeling approach is that the model is able to capture a wide range of relationships between observed behaviors and variables of interest, for example, the interaction effects between two (or more) acquisition variables and the outcomes of interest. As the model will recover them from the data, those (linear or non-linear) relationships do not need to be pre-specified. Second, unlike traditional dimensionality reduction methods, the number of latent factors do not need to be specified a priori. The model infers the number of relevant dimensions from the data through automatic relevance determination. Third, the model is scalable, being applicable to datasets with large numbers of customers and many acquisition characteristics, some of which might contain missing observations. When present, these missing observations are easily handled by the

FIM, which models them as outcomes using a Bayesian estimation framework. Lastly, the proposed modeling framework is general in the sense that can be integrated with any demand specification, from simple linear specifications to more complex model structures that incorporate a latent attrition component (a.k.a., “buy-till-you-die” models) or other forms of customer dynamics (e.g., hidden Markov models). This desirable feature implies that marketers across business settings, contractual and non-contractual, can use this framework by making minor adjustments to the demand/transactional model.

Using a set of simulation analyses, we demonstrate the FIM inferences for newly-acquired customers’ to be more accurate than those generated by multiple tested benchmarks. Unlike other models, our approach accommodates flexible relationships among relevant behaviors, enabling the model to make accurate inferences about newly-acquired customers when the relationships between acquisition characteristics and demand parameters are unknown to the firm or researcher.

We then apply the FIM to a retail context and demonstrate how the focal firm can overcome the cold start problem by augmenting the (thin) historical data using their transactional database and employing the proposed modeling framework that extracts the relevant information from the augmented customer data. First, we use the transactional data to extract the characteristics of every customer’s first purchase (namely price paid, number of products purchased, etc.) as well as observed product characteristics such as category purchased, package size, etc. Second, we leverage the transactional data from customers outside our sample to create a continuous multidimensional representation of products (or product embeddings). Specifically, we use the word2vec algorithm — a machine learning approach originally developed to analyze textual data — to model the co-occurrence of

products in customer baskets. This yields a set of product embeddings that can be used to augment data on customers' first transactions based on the specific products they bought. We then estimate the FIM to the augmented cold start data and make individual-level predictions for newly-acquired customers outside the calibration sample.

We empirically demonstrate the superiority of the FIM at distinguishing, immediately after they make their first purchase, heavy spenders from those expected to yield less value. The model can be also used to highlight the set of acquisition characteristics most predictive of future behavior. For example, we find the predicted Top 10% heavy spenders to be less likely to be acquired during the holiday period and more likely to be acquired offline, and their first purchases to tend to include expensive and discounted products. The model also captures differences in customer responsiveness to marketing actions, enabling firms to identify and characterize those most (or least) sensitive to specific marketing communications. For example, we find that customers most sensitive to email marketing are more likely to be acquired online and buy less expensive products, and their first purchases to include fewer units. We also find non-linear relationships between acquisition characteristics and customer responsiveness to marketing actions. For example, the differences in email sensitivities across customers that received discounts on their first purchase only exist for those who also purchased a recently introduced product.

The present research develops a modeling framework that overcomes the cold start problem by linking customers' early observed behaviors and choices with future purchase behavior, enabling firms to make meaningful predictions about customers just acquired. Methodologically, our paper contributes to the CRM literature by being the first to incorporate in a general, flexible, and scalable way information obtained at the moment of

acquisition (generally discarded due to an inability to use it effectively). Substantively, our research is relevant to marketers faced with the challenge of managing customers soon after acquisition. We show how the proposed modeling framework enables firms to identify and characterize, from information collected at the moment of acquisition, high-value customers and those most sensitive to marketing communications. From a practical perspective, our research guides firms in the use of cold start data to augment information already in their databases. To that end, we employ recent developments in machine learning and natural language processing to create a matrix of product “embeddings” that enable firms to characterize (even recently acquired) customers based on the products they purchase. We believe this approach to customer segmentation to be highly promising, enabling firms to obtain rich information about individual customers without recourse to customer-provided data or external sources that might pose privacy concerns.

The remainder of the paper is organized as follows. Following a brief review of the literature related to our work, we introduce the cold start problem and illustrate the main challenges to solving it in practice. We next present our modeling framework, discuss its components, and evaluate its performance vis-à-vis existing approaches that could be used to solve the cold start problem. We then apply our model in the context of an international beauty and cosmetic retailer. We conclude with a discussion of the implications, managerial relevance, and future directions of our research.

1.2 Previous literature

Our research relates to the broad literature on customer-base analysis that has provided managers and analysts with tools for understanding, forecasting, and managing the (heterogeneous) behavior of customers. It relates particularly to work that has incorporated the effect of marketing variables or, more generally, time-varying covariates in customer lifetime value (CLV) models. Notable work in this area includes Schweidel and Knox (2013) and Schweidel et al. (2014) who, building on the foundations of the Beta-Geometric/Beta-Binomial (BG/BB) model (Fader et al., 2010), incorporate the effect of direct marketing activity and past customer activity on the latent attrition process and the customer's purchase propensity while alive, and Knox and van Oest (2014) and Braun et al. (2015) , who incorporate the effect of the customer service experience and customer complaints on the latent attrition process of the Beta-Geometric/NBD (BG/NBD) model (Fader et al., 2005). Our research and methodological objectives differ in two main ways. Whereas the main purpose of the aforementioned studies is to capture the effect of time-varying marketing variables (e.g., direct marketing activities, customer complaints) on customer behavior, we extract as much information as possible from cold start data. The referenced models, although they could be used to incorporate a handful of pre-specified acquisition variables, are not well suited to extract relevant information from noisy and redundant variables, the case with cold start data. Second, we do not build on a specific demand specification tied to a business context, but rather provide a modeling framework that can incorporate any of the models of behavior presented in the foregoing papers.

On a substantive level, our work relates to Gopalakrishnan et al. (2016), who propose a framework for multi-cohort data able to predict the behavior of new cohorts of customers for whom little transactional data is available. Gopalakrishnan and colleagues build a model that allows customers to be inherently different depending on when they were acquired (i.e., *which cohort* they belong to), while capturing the underlying dynamics across cohorts. We posit that such inherent heterogeneity can be explained (at least partially) by individual-level observed characteristics collected when customers make their first purchase. This is consistent with Anderson et al. (2020) who document the existence of “harbinger products.” These are products that, when purchased by a customer in their first transaction, are an indicator of the customer being less likely to purchase again, and hence, provide less value to the firm. Our work also relates to Loupos et al. (2019), who use social network data for recently acquired customers to explain heterogeneity in their future value to the firm. To the best of our knowledge, our approach is the first to integrate several types of information collected at the moment of acquisition, and to differentiate responsiveness to marketing actions — not only individual propensity to transact — on the basis of customers’ first purchases. The latter aspect is crucial in cases in which targeting occurs soon after the customer is acquired or when securing a second purchase is challenging.

The premise that behaviors observed at the moment of acquisition can help firms explain heterogeneity in future behavior is consistent with empirical findings in the CRM literature (e.g., Fader et al., 2007; Voigt and Hinz, 2016), specifically, work on customer acquisition that has investigated the relationship between acquisition-related information — e.g., channel of acquisition — and subsequent customer lifetime value (e.g., Verhoef and Donkers, 2005; Lewis, 2006; Villanueva et al., 2008; Chan et al., 2011; Steffes

et al., 2011; Schmitt et al., 2011; Uncles et al., 2013; Datta et al., 2015). Our work, although it investigates relationships between acquisition-related variables and subsequent customer behavior, differs in two important ways. First, our end goal is to inform decisions related to the management of already acquired customers (e.g., whom to target in the next campaign) rather than the design of optimal strategies for customer acquisition (e.g., free trials to increase customer acquisition). The goal of our modeling framework is to extract as much observed heterogeneity as possible from initial behaviors while controlling for firms' acquisition activities rather than estimate the casual impact of these acquisition variables on future behavior. Second, this literature suggests that customers are inherently different depending on how they have been acquired. We broaden the range of acquisition-related behaviors by looking not only at *how* a customer was acquired (e.g., online vs. offline, trial vs. regular), but also *what* they did when they were acquired (e.g., what kind of product did they buy? how much did they pay?), hence extracting more information from the initial transaction. The latter is especially relevant for managers and analysts in large retail and hospitality businesses, among others, such information not only being easily observed, but typically already residing in their databases.

From a methodological perspective, we contribute to the literature on applying probabilistic machine learning methods to marketing (Jacobs et al., 2016; Dew and Ansari, 2018; Dew et al., 2020). More specifically, our work relates to the literature on applying deep exponential families (Ranganath et al., 2015) as building blocks of more complex models (Ranganath et al., 2016; Wang and Blei, 2019), and other generative models such as Bayesian Principal Component Analysis (Bishop, 1999; Mohamed et al., 2008).

1.3 The “cold start” problem of CRM

We turn to a retail context to illustrate the cold start problem, and to motivate and validate our modeling framework. Retail is a good context to examine this phenomenon for several reasons. First, firms in this sector increasingly collect transactional data and rely on analytics to better manage their customers (Forbes, 2015). Second, retail represents a large proportion of the total economy, with revenues accounting for 31% for the global GDP (Research and Markets, 2016). Finally, the data structure in most retail settings—in particular, the one used in this research—resembles that in many other industries such as hospitality, entertainment business, or non-for-profit organizations, that face similar data challenges when implementing CRM programs.

1.3.1 The “cold start” problem

Consider a retailer that sells cosmetic/beauty products both via online and offline channels.¹ Like most other companies, it records the transactions of all individual customers since the moment they were acquired, including the time of purchase, the products purchased in each particular transaction, their price and discounts (if any), along with information about the CRM activities that the company engaged with, such as email marketing activities. With these transactional data at hand, the focal company could apply some of the aforementioned models and be able to predict, with a good degree of accuracy, the number of transactions that customers with different transaction patterns would make in future periods (e.g., Fader

¹This will be the specific context of our empirical application. The full set of details about the focal firm and the data will be presented in Section 1.5; in this section we only present the relevant information to motivate the business problem and the modeling challenges.

et al., 2010). The marketer can also incorporate the historical marketing actions to capture how those variables affected transaction propensities and customer value (e.g., Schweidel and Knox, 2013; Schweidel et al., 2014). However, when making these types of inferences for recently acquired customers, for whom the firm has no transactional history nor past marketing interventions, the “best guess” that the marketer can get is the population average. This is what we call the “cold start problem of CRM” whereby firms cannot make individual-level inferences about newly-acquired customers that differentiates them, therefore diminishing the effectiveness of future CRM activities.

The premise of this research is that, while it is the lack of (historical) data that causes the cold start problem, firms nowadays have access to other data sources that, properly leveraged, can help them overcome the cold start problem. Granted, if firms only observed that the customer made “a transaction” it would be very difficult to overcome the cold start problem. However, most firms not only know when a customer made their first transaction but also record the details such as the channel/store used, the exact product the customer purchased, the price paid, whether they bought in discount, the time of the day, and so forth.² We propose leveraging those (already existing) data and extract what we call “acquisition characteristics” from each customer’s first transaction.³ We contend that these

²Note that the amount of data collected by firms also include data *prior* to the moment of acquisition. For example, e-retailers collect information via cookies, which could identify which customers have visited the website previously (yet, not making a purchase). When available, those data can be included in the exact same fashion as the acquisition characteristics. For simplicity, we denote “acquisition” data to all information available to the firm at the moment of acquisition, acknowledging that such data could also incorporate actions the customer performed before their first transaction.

³In theory, the data could also be augmented with characteristics of the second, or third transaction, for customers who are repeat buyers. However, we only use the first transaction because that is the data that *every* customer — just acquired and existing users — have in common, which will be the key to make inferences about recently-acquired customers. Adding information about each later transactions might add precision to the individual-level inferences of repeat users, but not necessarily to the inferences of recently-acquired customers, which is the main focus of this paper.

acquisition characteristics/choices can be informative about underlying customer differences which can be predictive of customers behavior in the future. Because these data are also available for customers with longer tenure with the company, the firm would be able to uncover the (subtle) relationships between the choices observed at the moment of acquisition and the customer behavior down the road.

1.3.2 Augmenting cold start data with acquisition characteristics

Considering the retailer introduced above, who is trying to make inferences about its customers right after they have been acquired. A natural first step for the analyst would be to select a handful of variables collected at the acquisition moment (e.g., channel of acquisition) and use existing models to relate those characteristics to future demand (e.g., Chan et al., 2011). The caveat of doing so is that merely few variables might not fully capture the richness of the acquisition data, and the level of personalization would likely be limited as these few variables only capture a coarse representation of customers' heterogeneity. We propose to fully augment the acquisition data to broaden the amount of information that would (potentially) be linked to future behavior, therefore increasing the chance to solve the cold start problem.

Specifically, using the (existing) data from each first transaction, we propose to augment cold start data with three types of acquisition variables: *transaction characteristics* (e.g., channel, price paid, holiday season) and *product characteristics* (e.g., product category, package size), which are easily extracted from the transactional database, and *shopping basket (latent) representation*. The latter type of data aims to capture the “nature” of

products that the customer purchased, above and beyond what the standard (observed) product categories represent. Our premise is that the nature of products purchased can signal the type of customer who purchases those. For example, in the market of cosmetics, certain ingredients or aroma characterize lines of products. It is possible that customers who discover the brand by buying products of certain “nature” are similar in the way they behave in the future. Because such information is not readily available from the firm’s database, we need a method to encode the information embedded in each product, to then aggregate it at the basket level.⁴

Previous literature has used different methods to encode such information, from human coding based on full description of the product, to machine learning approaches that apply textual analyses to the description of products, or that leverage co-occurrence of products in basket data to create measures of similarity across products (e.g., Jacobs et al., 2016; Ruiz et al., 2017; Kumar et al., 2020; Chen et al., 2020). We take the latter approach and leverage the transaction data from anonymous customers to create continuous multidimensional representations of products, called product embeddings, that capture the nature of the product. Specifically, we create a co-occurrence matrix based on the composition of shopping baskets — i.e., which SKUs are purchased together — and implement *word2vec* (Mikolov et al., 2013), a machine learning approach widely used for natural language processing, to map each item to a multi-dimensional vector that captures similarities across products. This exercise is similar to creating a perceptual map from

⁴One alternative to this solution would be to include a dummy variable per (available) SKU. This approach would be straight forward in business contexts where the product space is small. However, when the firm offers a large selection of items or SKUs — as it is the case for most retailers — the vector of dummy variables would be too sparse to capture similarities among baskets and thus would prevent any model to learn across customers. For those cases, we recommend using a lower-dimensional vector representing the product space, as we do in this research.

association data (Netzer et al., 2012) in which the co-occurrence of products in a basket is used as proxy of association between two products. (See Appendix A.1 for all the details about how we process the transaction data and create the product embeddings using the *word2vec* algorithm.) Once we represent each product by a continuous vector, we can easily characterize the first purchase of any customer by computing moments of the product vectors in that basket.

In sum, using the transactional data already collected by the firm, one can easily augment each customer’s data with a high-dimensional vector that captures a wide variety of acquisition characteristics including details about the first transaction as well as the type of products purchased.⁵

1.3.3 Predictive power of augmented data

A natural question to ask is: Do acquisition characteristics carry information about future behavior? While this is an empirical question, we present preliminary evidence from our empirical application that these augmented acquisition characteristics in turn explain differences in subsequent demand behavior across customers. To do so, we select customers who have been with the company for at least 15 months and relate their total number of repeat purchases during those 15 months with their (augmented) acquisition characteristics. We explore the relationship between individual acquisition characteristics and future transactions (Figure 1.1), as well as possible interactions among acquisition variables in their correlation with future demand (Figure 1.2).

⁵In our empirical application this vector has 31 dimensions. Further details are presented in Section 1.5.

Indeed, acquisition characteristics are predictive of customers future transactions. Consistent with common belief in the industry (e.g., Artun, 2014; RJMetrics, 2016), customers that were acquired during the holiday season are less valuable to the firm, as we find that they are less likely to transact in the future. On the other hand, customers who

Figure 1.1: Transactions versus acquisition characteristics. Observed repeated transactions as a function of a sample of augmented acquisition characteristics. All acquisition variables are constructed from the first transaction of each customer. Repeated transactions do not include the first transaction.

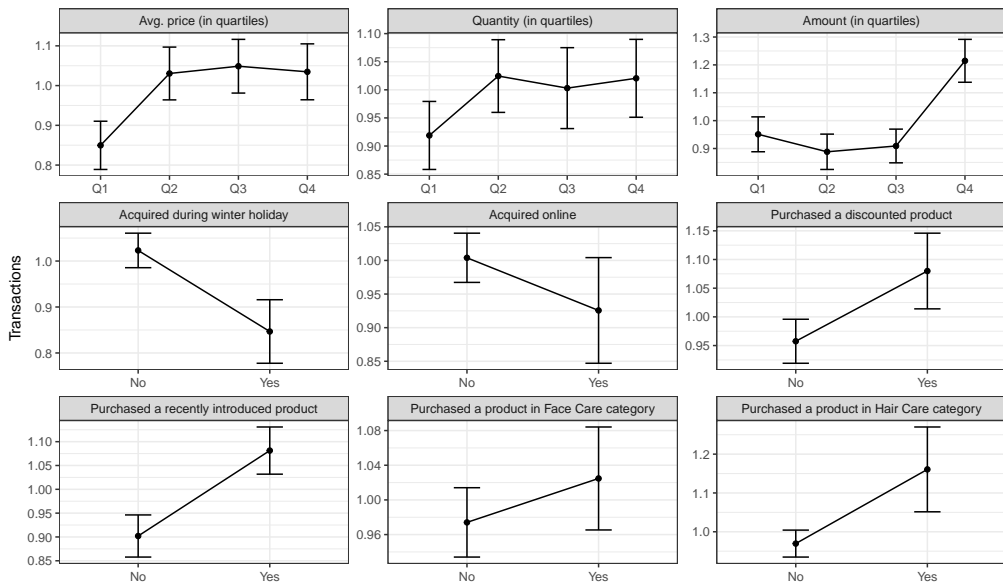
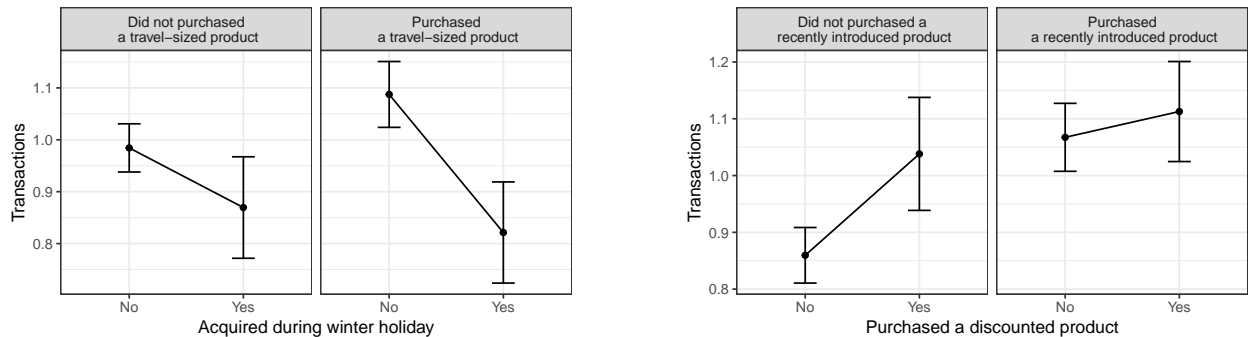


Figure 1.2: Transactions versus interaction of acquisition characteristics. Observed repeated transactions as a function of interactions among acquisition characteristics. All acquisition variables are constructed from the first transaction of each customer. Repeated transactions do not include the first transaction.



bought using discounts on their first transaction generally buy more during the next 15 months than customers who did not. A similar pattern exists for customers who bought a recently-introduced product on their first transaction, and those who purchased products from the hair care category. Interestingly, this model-free analysis also suggest that some of these relationships are likely to be non-linear. For example, looking at average price paid per item, customers that bought more expensive products in their first transaction tend to buy more frequently in the future. Noteworthy, this relationship is not linear. Customers in the lowest quartile (Q1) tend to buy less frequently in their first 15 periods than all other customers. Similar non-linear relationships appear for the number of units and the total amount of the ticket.

Interesting patterns also emerge in Figure 1.2. On the left, we group customers on whether they were acquired during the winter holiday season, coupled with whether they purchased travel-size products. We find that purchasing travel-size products moderates the relationship between being acquired during the holidays and the future number of transactions. Turning to the figure on the right, we observe that purchasing a discounted product on the first transaction signals lower value *only* if such a purchase did not include a new product. Taken together, these results present evidence of a relationship between acquisition characteristics and future transactions, confirming that augmenting cold start data with acquisition characteristics incorporates relevant information to infer customers' differences.

Nevertheless, this simple analysis is insufficient for solving the cold start problem of CRM as would likely miss useful information from the data. First, it can only be performed for sub-sample of customers — those for whom we observe for relatively long period of time

(e.g., 15 months) — in order to have a fair comparison across customers over the same number of periods. Second, this type of analysis examines each variable independently (Figures 1.1), at most allowing for single interactions (Figure 1.2). Given that the goal is to extract relevant correlations in high-dimension cold start data, it will be more effective (and efficient) to examine these correlations collectively, while allowing for flexible relationships among the variables. Furthermore, the model-free analysis does not shed any light about customers’ response to marketing actions. These results indicate that “holiday” customers are less likely to transact again. However, are they more/less sensitive to the firm’s communication? How strongly will they react product introductions? A model would be certainly necessary to effectively extract the information from the acquisition characteristics to predict differences in transaction propensities *as well as* in responsiveness to marketing actions. Before presenting our modeling framework, we describe the methodological challenges that such a model should overcome.

1.3.4 Modeling challenges

Our solution to overcome the cold start problem ultimately depends on the ability of the model to extract the information hidden in the augmented data that is predictive of future behavior. Naturally, increasing the dimensionality of the acquisition data increases the chances of adding (at least potentially) information that will be relevant to infer customer differences down the road. However, expanding the dimensionality of the acquisition data also adds methodological challenges.

First, several of those augmented variables are likely to be irrelevant. Many of the behaviors observed in the first purchase are likely to be random and not systematically related with how customers will behave in the future. Second, some of these augmented data are multiple signals from the same underlying behaviors, implying that much of those data would be redundant. For example, a price-conscious customer may purchase a set of travel-sized, cheap products that are discounted. Although, the variables price and discount capture different types of information (e.g., a discounted product may still be an expensive one), these variables are clearly correlated as they are both signals of this customer's preferences for inexpensive products. Moreover, if one also were to include latent representations of the products bought, these representations may also correlate with the prices that these products are sold and how frequently they are discounted; adding to the total correlation present among augmented variables. Taken together, these characteristics suggest that it is likely that cold start data would have low "signal-to-noise" ratio, increasing the difficulty of recovering the relationships between acquisition characteristics and future behavior.

Importantly, the underlying relationships between acquisition variables and future demand is unknown. As indicated by the early exploration of the data (Figures 1.1 and 1.2), those relationships are unlikely to be linear. It is unrealistic to recommend that a firm would explore all possible interactions and non-linear specifications among their augmented acquisition characteristics, and is especially cumbersome when also interested in customers' response to marketing actions. Moreover, increasing the dimensionality of the augmented data only emphasizes this challenge as it would increase the number of potential non-linear relationships and interactions among acquisition variables. Another potential limitation of

increasing the dimensionality of the acquisition variables is that some variables might be missing for some customers. Missing observations present challenges to estimate models that use those missing variables as covariates as they require imputation methods—cumbersome for high-dimensional spaces—or deletion of customers (or variables) from the data—which directly reduces the amount of information, defeating the purpose of the data augmentation step.

In this research, we propose a modeling framework that overcomes all these issues at once. We combine a flexible demand specification (such that can be applicable to a wide range of marketing contexts) with state-of-the-art machine learning methods (addressing nonlinearities and data redundancy) within a Bayesian framework (that extract signals from the acquisition characteristics while handling missing data). The resulting modeling framework is a flexible probabilistic machine learning model that links the individual-level parameters governing customer’s future behavior (e.g., transaction propensities, sensitivity to marketing actions) with a latent representation of the behaviors/choices observed at the moment of acquisition. This modeling approach seamlessly captures flexible relationships among variables (linear and non-linear) without the need to pre-specify those relationships a priori. Moreover, the model explicitly accounts for correlations in the acquisition data which helps regularize the flexible model avoiding overfitting.

These benefits will become clear as we build and validate the model in the next section, where we also show how this approach dominates existing alternatives that addressed some (but not all) modeling challenges. For example, we compare it with a standard hierarchical Bayesian model with acquisition characteristics are included as covariates; a fully hierarchical model where acquisition characteristics and demand are jointly

correlated using a multivariate Gaussian distribution; or a (supervised) Bayesian PCA that aims to reduce dimensionality of acquisition characteristics as well as demand parameters.

Finally, as we show in our empirical application that, if we simplify the task and only consider the model’s ability to predict future transactions, our modeling approach performs at the level of traditional machine learning (ML) approaches such as a random forest and a deep neural network (proven to capture non-linear relationships very well). Our model stands out in comparison with these ML benchmarks in two main ways. Methodologically, it can be easily be combined with multiple demand specifications, as well as allows for missing observations in acquisition characteristics without relying on data imputation. Practically, our model provides inferences beyond predictions of future transactions, enabling marketers to get insights about customer heterogeneity in preferences and in sensitivity to marketing actions.

1.4 Modeling framework

1.4.1 Model development

Our modeling framework — which we call “First Impression Model” (FIM) — comprises three main components: (1) the *demand model*, main outcome of interest to the firm, which could include customers transactions, purchase volume, etc., (2) the *acquisition model*, capturing all customer outcomes that are observed to the firm at the moment of acquisition, and (3) the *probabilistic model* that links the underlying customer parameters influencing these two types of behaviors through hidden traits.

1.4.1.1 Demand model

We start by assuming a general model for demand, suitable for different specifications, and parametrized using individual-level parameters and population-level parameters. Specifically, for customer i at period t , we denote

$$p(y_{it}|\tilde{\mathbf{x}}_{it}^y, \boldsymbol{\beta}_i^y, \boldsymbol{\sigma}^y) = f^y(y_{it}|\tilde{\mathbf{x}}_{it}^y, \boldsymbol{\beta}_i^y, \boldsymbol{\sigma}^y) \quad i \in \{1, \dots, I\}, t \in \{1, \dots, T_i\}, \quad (1.1)$$

where I represents the total number of customers, T_i denotes the number of periods since the customer was acquired, $\boldsymbol{\beta}_i^y$ is a vector containing customer i 's individual-level parameters, the vector $\boldsymbol{\sigma}^y$ contains the parameters that are common across customers, and $\tilde{\mathbf{x}}_{it}^y$ includes the observed covariates for customer i at period t . Finally, $f^y(\cdot)$ is the pdf/pmf for outcome y_{it} ; for example, if the outcome of interest is purchase incidence, we would specify

$$p(y_{it} = 1) = \text{logit}^{-1} [\mathbf{x}_{it}^{y'} \cdot \boldsymbol{\beta}_i^y].^6$$

1.4.1.2 Acquisition model

We denote A_i the vector of characteristics that are collected at the moment of acquisition, and a_{ik} the k 'th component/behavior (e.g., did the customer purchase a discounted product on their first transaction?). These acquisition characteristics are likely to be influenced by individual-level parameters (e.g., does this customer have the tendency to buy on discount?)

⁶The model can easily be adapted to other forms of demand (e.g., continuous demand, count) and extended to dynamic specifications such as latent attrition models. For the latter, one could define (1.1) as a state-space model (e.g., a hidden Markov model) with state variable s_{it} and $p(y_{it}, s_{it}|y_{i1:t-1}, s_{i1:t-1}) = p(y_{it}|s_{it}) \cdot p(s_{it}|s_{i1:t-1})$. We would implement such a model by having two individual level vectors, $\boldsymbol{\beta}_i^{yq}$ and $\boldsymbol{\beta}_i^{ye}$, as well as two population level vectors, $\boldsymbol{\sigma}^{yq}$ and $\boldsymbol{\sigma}^{ye}$, that would govern transitions among the hidden states and emissions in a state, respectively. We would substitute (1.11) for $p(y_{it}, s_{it}|y_{i1:t-1}, s_{i1:t-1}, \mathbf{x}_{it}^y, \boldsymbol{\beta}_i^y, \boldsymbol{\sigma}^y) = p(y_{it}|s_{it}, \mathbf{x}_{it}^y, \boldsymbol{\beta}_i^{yq}, \boldsymbol{\sigma}^{yq}) \cdot p(s_{it}|s_{i1:t-1}, \mathbf{x}_{it}^y, \boldsymbol{\beta}_i^{ye}, \boldsymbol{\sigma}^{ye})$, where $\boldsymbol{\beta}_i^y = [\boldsymbol{\beta}_i^{yq} \quad \boldsymbol{\beta}_i^{ye}]$, and $\boldsymbol{\sigma}^y = [\boldsymbol{\sigma}^{yq} \quad \boldsymbol{\sigma}^{ye}]$ be the parameters of the demand model.

but also by the market conditions at the moment of acquisition (e.g., was the company running heavy discounts during that period?). We account for these effects by modeling the acquisition characteristics as a probabilistic outcome, rather than as an input/covariate to the model. Note that we do not model acquisition per se, i.e., *whether* the customer is acquired or not. Rather, we model the characteristics of their first purchase, given that the customer was acquired.

Modeling the acquisition characteristics as an output not only allows us to control for the time-varying factors that shift demand at the moment of acquisition, but also allows for a flexible modeling specification of the latent traits that overcome challenges such as redundancy, irrelevance of variables, and missing data commonly encountered in the firm’s database. (We discuss these challenges in Section 1.4.1.3). Specifically, we denote

$$p(a_{ip}|\beta_{ip}^a, \boldsymbol{\sigma}_p^a, \mathbf{x}_{m(i)\tau(i)}^a) = f_p^a(a_{ip}|\beta_{ip}^a, \boldsymbol{\sigma}_p^a, \mathbf{x}_{m(i)\tau(i)}^a) \quad i \in \{1, \dots, I\}, p \in \{1, \dots, P\}, \quad (1.2)$$

where P is the number of different types of behaviors collected at acquisition, β_{ip}^a is an individual level parameter that reflects tendency to observe such a behavior when customer i is acquired, $\boldsymbol{\sigma}_p^a$ denotes a vector of parameters that are common across customers, and $\mathbf{x}_{m(i)\tau(i)}^a$ comprises the set of market-level covariates, with $m(i)$ indicating the market customer i belongs to, and $\tau(i)$ denoting the time period at which the customer was acquired.

The term $f_p^a(\cdot | \cdot)$ is the pdf/pmf of a distribution to model acquisition behavior p . Note that some of these behaviors will likely be discrete (e.g., whether the customer was

acquired online), in which case we specify $\boldsymbol{\sigma}_p^a = [\mathbf{b}_p^a]$ and model p as

$$p(a_{ip} = 1) = \text{logit}^{-1} [\beta_{ip}^a + \mathbf{x}_{m(i)\tau(i)}^a \cdot \mathbf{b}_p^a]. \quad (1.3)$$

For continuous acquisition variables (e.g., total amount spent in the first transaction) we define $\boldsymbol{\sigma}_p^a = [\mathbf{b}_p^a, \sigma_p^a]$ and model p as

$$p(a_{ip}) = \mathcal{N}(\beta_{ip}^a + \mathbf{x}_{m(i)\tau(i)}^a \cdot \mathbf{b}_p^a, \sigma_p^a), \quad (1.4)$$

specification that can be easily adjusted for multivariate outcomes as we do with some acquisition variables in our empirical application.

All of these types of variables are easily incorporated by adjusting the acquisition model accordingly. We define $\boldsymbol{\beta}_i^a = \begin{bmatrix} \beta_{i1}^a & \dots & \beta_{iP}^a \end{bmatrix}$ and $\boldsymbol{\sigma}^a = \begin{bmatrix} \sigma_1^a & \dots & \sigma_P^a \end{bmatrix}$ as the full set of individual- and population-level vectors of acquisition parameters, respectively.

Note that we only have one observation per individual and behavior. Hence, in theory, having an individual-level parameter β_{ip}^a could completely capture the residual variance of a_{ip} that is not systematically explained by the market-level factors (as in a regression with individual random effects but only one observation per individual). However, because we model demand and acquisition jointly, our model will balance fitting each acquisition behavior a_{ip} with fitting the other acquisition characteristics, as well as fitting demand, with a reduced set of individual factors or traits. Therefore, the individual level parameters β_{ip}^a will not have full flexibility to accommodate perfectly to the behavior a_{ip} . Rather, these parameters will capture the residual variance that is correlated with the rest of

the acquisition variables and with the demand model. This remark will become clearer when we specify the relationship between the individual-level demand and acquisition parameters, β_i^y and β_i^a , as we do in the next section.

Finally, the term $\mathbf{x}_{m(i)\tau(i)}^a$ controls for the overall marketing intensity that a yet-to-be-acquired customer might have been exposed to in a particular market at the moment of acquisition. For example, if there is a strong promotional activity in market m in period t , one would likely observe a higher-than-usual share of discounted products among the acquisition characteristics, not only driven by the customers’ propensity to buy on discount, but also by the fact that the majority of products were discounted.⁷ Accordingly, we want to capture this systematic shift in the acquisition characteristics as a market-related shift and not as a customer-driven shift, and therefore set \mathbf{b}_p^a common across customers.

1.4.1.3 Linking acquisition and future demand: Deep probabilistic model

We use a deep exponential family (DEF) component (Ranganath et al., 2015) to relate demand and acquisition parameters hierarchically, through hidden layers. We chose such specification because of its hierarchical nature — allowing the model to identify/extract individual-level traits that affect both acquisition and future demand — and because the presence of multiple layers facilitates the reduction of dimensionality while accommodating a wide range of possible relationships between acquisition and demand variables. Furthermore, one important characteristic of DEFs is that the latent variables are distributed according to distributions that belong to the exponential family (e.g., Gaussian, Poisson, Gamma),

⁷If the model did not control for these market-level conditions and the firm managed acquisition and retention efforts strategically, the correlations between acquisition characteristics and demand parameters obtained by the model could be spurious in the sense that they could be driven by the firm’s actions and not by customers’ underlying preferences.

making them a good candidate to model the wide range of data types encountered in the firm’s database. Finally, DEFs also enjoy the flexibility of probabilistic models, allowing them to be easily incorporated in more complex model structures, as we do in this research. (See Appendix A.2 for more details on DEFs.)

Turning our attention to our modeling challenge, the primary goal of our model is to infer the individual-level parameters β_i^y . Therefore, we specify the DEF component such that the lowest level captures the individual-level traits that affect both the acquisition characteristics and future demand. Specifically, we define

$$\beta_i^y = \boldsymbol{\mu}^y + \mathbf{W}^y \cdot \mathbf{z}_i^1 \quad (1.5)$$

$$\beta_i^a = \boldsymbol{\mu}^a + \mathbf{W}^a \cdot \mathbf{z}_i^1 \quad (1.6)$$

such that the individual level parameters, β_i^y and β_i^a are a (deterministic) function of mean parameters, $\boldsymbol{\mu}^y$ and $\boldsymbol{\mu}^a$, and individual deviations from this mean which are a function of the lower layer vector \mathbf{z}_i^1 , and weight matrices \mathbf{W}^y and \mathbf{W}^a . Similarly as in a Bayesian Principal Components Analysis (Bayesian PCA) model (Bishop, 1999), the vector \mathbf{z}_i^1 captures the individual level traits that explain jointly demand and acquisition behavior. The weight matrices \mathbf{W}^y and \mathbf{W}^a capture how each one of these traits manifests in both demand and acquisition characteristics respectively.

We assume that each component k of the lower layer, z_{ik}^1 , is distributed Gaussian with mean $g(-\mathbf{w}_k^{1'} \cdot \mathbf{z}_i^2)$, and variance 1,

$$p(z_{i,k}^1 | \mathbf{z}_i^2, \mathbf{W}^1) = \mathcal{N}\left(z_{i,k}^1 | g\left(-\mathbf{w}_k^{1'} \cdot \mathbf{z}_i^2\right), 1\right) \quad k \in \{1, \dots, N_1\}, \quad (1.7)$$

where N_1 is the dimension of the lower layer, $g(x) = \log(\log(1 + \exp(x)))$ is the log-softplus function (Ranganath et al., 2015),⁸ and \mathbf{W}^1 is the weight matrix that links the upper and lower layers. The upper layer captures higher-level traits (resembling the structure of neural networks), while allowing for non-linear correlations between the traits in the lower level \mathbf{z}_i^1 . The correlations among the lower layer components are induced by reducing the dimensionality of top layer (\mathbf{z}_i^2 is a vector of length N_2 , with $N_2 < N_1$)⁹ whereas the non-linear relationships are captured by the non-linear link function $g(\cdot)$, which relates the higher-level traits with the lower-level traits that manifest in demand and acquisition. Finally, we model the upper layer using a standard Gaussian distribution,

$$p(z_{i,k}^2) = \mathcal{N}(z_{i,k}^2 | 0, 1) \quad k \in \{1, \dots, N_2\}. \quad (1.8)$$

To sum, we link the individual-level demand and acquisition parameters using a DEF component of two Gaussian layers, \mathbf{z}_i^1 and \mathbf{z}_i^2 . The model could easily accommodate more layers (e.g., Ranganath et al., 2015, use up to 3 layers, $L \leq 3$, in their empirical applications).¹⁰

⁸In Stan, the softplus function, defined as $f(x) = \log(1 + \exp(x))$, can be computed using `log1p_exp(·)`.

⁹In theory, N_2 could be larger than N_1 but such a model would not necessarily reflect patterns in data as information would be lost going from the upper layers of the DEF to the lower layers of the DEF. Ranganath et al. (2015) only estimate models with decreasing dimensions of upper layers.

¹⁰We follow the specifications from Ranganath et al. (2015), where the model is estimated using, at most, 3 layers ($L \leq 3$). In that paper, the model is trained on two large text corpora (5.9K and 8K terms), two matrix factorization tasks on a movie ratings dataset (50K users and 17.7K movies), and a click dataset (18K users and 20K documents). All of these datasets are considerably larger than our data (both in the simulations and in the empirical application). Furthermore, Tables 2 and 3 from Ranganath et al. (2015) do not show consistently whether $L = 3$ is better than $L = 2$. As a result, we use $L = 2$ as it is the smallest configuration that allows for non-linear relationships.

1.4.1.4 Dimensionality of the DEF component

At first glance, the choice of the layers dimensions N_1 and N_2 may seem cumbersome. On the one hand, high values of N_1 and N_2 increase the computational burden of the inference procedure, which is not desirable. On the other hand, a model with low values for N_1 and N_2 may miss relevant correlations that are needed to infer customers' parameters. In the extreme, if the number of components of the lower layer, N_1 , is set to one, the model would only learn a single trait to describe the variation across all parameters, which will fail to capture the heterogeneity in the demand parameters, and their (potentially non-linear) relationships with acquisition characteristics. Similarly, if the number of components of the higher layer, N_2 , is set to zero, the model would be stripped away from the non-linear function $g(\cdot)$ that allows the model to capture non-linear relationships between demand and acquisition parameters.

Similar to other latent-space models, one could test all possible combinations of N_1 and N_2 (increasing in magnitude) and choose the optimal values using cross-validation. Such exercise is certainly required when using Maximum Likelihood Estimation, as more flexibility in a model leads to over-fitting following the classical bias-variance trade-off, and therefore poor performance in holdout samples. However, when using Bayesian inference, this exercise would not only be computationally very costly, but also unnecessary, provided that adequate priors such as spike-and-slab or sparse-gamma (Karaletsos and Rättsch, 2015; MacKay, 1995; Neal, 2012) are used to induce regularization in the parameters governing the weights that activate the traits. Using such priors ensures that a trait only manifests in a particular

variable if the improvement in fit is significant; otherwise, that trait is “shut down” by the prior (Ranganath et al., 2015).

Therefore, our approach to specifying the dimensionality of the model is to set a “large enough” number of traits to ensure that all relevant traits are recovered, while using sparse priors to ensure that the model only activates the relevant traits, thus avoiding over-fit the data. Specifically, we use sparse Gamma priors for \mathbf{W}^1 and hierarchical Gaussian automatic relevance determination (ARD) priors for \mathbf{W}^y and \mathbf{W}^a , both of which are spike-and-slap-like priors that have shown to perform well on feature selection (e.g., Bishop, 2006; Kucukelbir et al., 2017). These priors ensure that once a trait is “shut down,” adding more traits (i.e., increasing N_1 or N_2) would just add irrelevant traits with weights all being close to zero, not affecting the performance of the model. (See Appendix A.3.1 for details about these priors.)

The added benefit of inducing regularization through the priors is that we can look at the posterior estimates of the variances of the weights (\mathbf{W}^y , \mathbf{W}^a , and \mathbf{W}^1) to evaluate whether the number of dimensions (N_1 and N_2) are sufficient to represent the data.

Examining N_1 is straightforward as the model parameter α^1 captures the variance of the lower layer traits. Regarding N_2 , while there is not one specific parameter capturing the relevance of the upper layer traits, we can compute a pseudo- α_m^1 for each upper trait m using the components of the weight matrix \mathbf{W}^1 that map to relevant lower level traits (see Appendix A.6.7 for details). Finally, examining the posterior estimates of α^1 and pseudo- α_m^1 — and observing that some traits have been “shut down” by the model — we corroborate whether N_1 and N_2 are “large enough” for any specific dataset.

These insights are further developed in Appendix A.6.7 where we explore the dimensionality of the DEF component by analyzing the results of estimating the FIM on

simulated data, where we know how many traits are needed. There we show how the performance of the model remains largely unchanged by the additional dimensions (on either N_1 or N_2) after the relevant number of traits are accounted for. We also show how the posterior estimates of the variances of the weights (α^1 and pseudo- α_m^1) are diagnostic of relevant and non-relevant traits.¹¹

To sum, we take a hybrid approach to model selection in which we make sure that the number of pre-specified dimensions is large enough — phenomenon that can be validated from the model parameters — while we rely on the priors of the model to ensure regularization.

1.4.1.5 Bringing it all together

We briefly discuss how each part of the model contributes to the desired goals and how the FIM compares with alternative approaches to overcome the cold start problem. In essence, the model comprises a demand and an acquisition model, whose individual-level parameters are projected into a lower-dimensional space through a two-layered DEF component. The lower layer of the DEF captures the relevant correlations among the individual-level parameters while reducing the dimensionality of those vectors. An alternative approach to link the acquisition and demand parameters could be through using traditional full hierarchical Bayesian priors (e.g., multivariate Gaussian). Such an approach would assume that all individual-level parameters (β_i^y and β_i^a) are distributed jointly according to a flexible multivariate distribution which parameters capture all the potential correlations among the variables. However, this full hierarchical approach would require the model to

¹¹The posterior distribution of α and \mathbf{W}^1 from real world data sets would not display as clear cut distinction between those traits that are meaningful and those that are not compared to our simulation analyses. We come back to this point when discussing the specification of the FIM for our empirical application.

estimate a very high-dimensional correlation matrix which can become computationally expensive, especially as the number of acquisition variables increases. On the contrary, because the FIM includes ARD priors for the lower layer of the DEF, the model only allows for “relevant” correlations to emerge, automatically reducing the dimensionality of the individual-level parameters. This is a desirable feature not only because the number of acquisition variables could be large, but also because some of the acquisition variables are likely to be correlated among each other.¹²

The upper layer of the DEF, and in particular, the non-linear link function $g(x)$ that relates the higher-level traits with the lower-level traits allows the model to capture a wide range of relationships — linear and non-linear — among the variables of interests. A simpler specification of the FIM would be one that does not incorporate the second layer and therefore imposes linear relationships among the individual parameters. Such a nested version of the FIM would be equivalent to a “supervised” factor analysis or Bayesian PCA where the latent traits are extracted from the acquisition variables as well as from the demand model. The limitation of such a (nested) approach is that the model would lose its accuracy at forming first impressions the moment the assumption of linearity does not hold, either because acquisition variables relate to demand parameters in a non-linear way, or when two (or more) acquisition variables interact in their relationship with the demand parameters. As we show in Section 1.4.4, our FIM specification (that includes the second

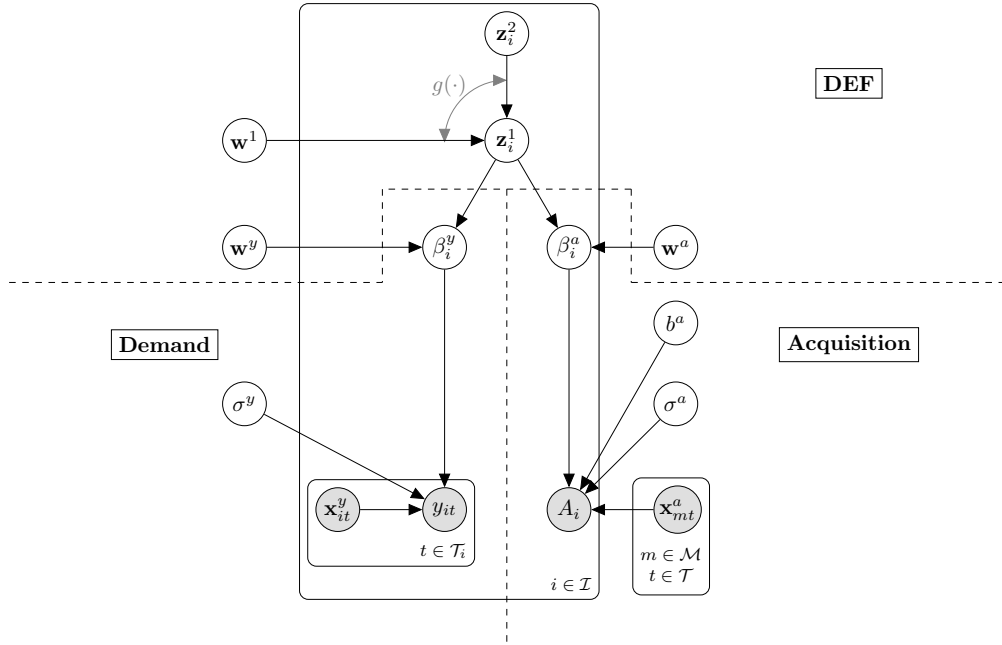
¹²An alternative but similar specification for the model could be a two-step approach that first reduces dimensionality among the acquisition variables (i.e., connecting z_i^1 to β_i^a) and then connects those factors with future demand. We choose to connect the lower level of the DEF model with both components jointly in order to be robust to the possibility that the residual variance of the acquisition variables not explained by the main factors of the first step is predictive of demand behavior; and to inform the choice of factors that are predictive of demand behavior, as in supervised topic models (Mcauliffe and Blei, 2008), and therefore, to overcome redundancy and irrelevance of acquisition variables simultaneously.

layer) captures several forms of relationships (including linear, interaction effects, and maximum function) without the need for specifying those relationships a priori. This is a very desirable property of the model because managers/researchers/data scientists generally do not know the exact form of the relationships among the variables of interest.

Finally, a different approach to overcome the cold start problem could be to simply specify the individual-level demand parameters (β_i^y) as a direct function of the acquisition variables (A_i^y). Such a specification would resemble a typical demand model with interactions, or a multi-level (hierarchical) model in which β_i^y are a function of the observed A_i and some population distribution (Rossi et al., 1996; Allenby and Rossi, 1998; Ansari and Mela, 2003; Chan et al., 2011). While a linear model is attractive for its simplicity and ease of interpretation, such an approach becomes intractable when the parameter space for the acquisition variables increases. Moreover, if the underlying relationships between the acquisition variables were not linear (or did not follow the specified relationship, due to variable transformation), the model will fail at inferring individual-level demand parameters for newly-acquired customers with certain level of accuracy. In addition, specifying the demand parameters as a direct function of the acquisition characteristics prevents the researcher from using acquisition characteristics that have missing observations. This is a key benefit that provides modeling acquisition characteristics as an outcome, as opposed to a direct function (we show this benefit by incorporating acquisition characteristics with missing observations in our empirical application in Section 1.5).

To conclude, Figure 1.3 shows the graphical model for the FIM, connecting all the individual components. We propose a model of demand and acquisition characteristics where the individual-level parameters of each of these sub-models are projected into a

Figure 1.3: Graphical model of first impressions



lower-dimension space via a DEF component. The specification of the demand sub-model is general such that the modeling framework can be applicable to a wide range of business contexts. The sub-model for acquisition characteristics enables the model to control for market conditions or firm-initiated actions that can potentially shift the type of customers that are acquired over time. If these shifts were not captured, the model would not be able to differentiate market conditions from customer underlying preferences. Regarding the DEF component, there are three main benefits of using a two-layered DEF to connect both types of individual-level parameters. First, the model provides dimensionality reduction, avoiding the curse of redundancy and irrelevance of variables among the acquisition variables. Second, the model allows for flexible relationships (e.g., non-linear relationships) among the model components. Third, the model can incorporate acquisition characteristics with missing observations, as these are modeled as outcomes which are easily handled using a Bayesian estimation framework. These benefits will become clearer in Sections 1.4.4 through 1.5,

when we compare the predictive accuracy of the FIM with that of several alternative specifications.

1.4.2 Estimation and identification

We estimate the model using full Bayesian statistical inference with MCMC sampling. We sample the parameters from the posterior distribution which is proportional to the joint,¹³

$$\begin{aligned}
p(\{\mathbf{z}_i^1, \mathbf{z}_i^2\}_{i=1}^I, \mathbf{W}^y, \mathbf{W}^a, \mathbf{W}^1, \boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}_a, \{y_{i1:T}, A_i\}_i) = \\
\left[\prod_{i=1}^I \prod_{t=1}^{T_i} p(y_{it} | \mathbf{x}_{it}^y, \mathbf{z}_i^1, \mathbf{W}^y, \boldsymbol{\mu}^y, \boldsymbol{\sigma}^y) \right] \cdot \left[\prod_{i=1}^I p(A_i | \mathbf{x}_i^a, \mathbf{z}_i^1, \mathbf{W}^a, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^a, \mathbf{b}_a) \right] \\
\cdot \left[\prod_{i=1}^I p(\mathbf{z}_i^1 | \mathbf{z}_i^2, \mathbf{W}^1) \right] \cdot \left[\prod_{i=1}^I p(\mathbf{z}_i^2) \right] \\
\cdot p(\mathbf{W}^y, \mathbf{W}^a, \mathbf{W}^1, \boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}_a). \tag{1.9}
\end{aligned}$$

In particular, we use the No U-Turn Sampling (NUTS) Hamiltonian Monte Carlo algorithm, implemented in the Stan probabilistic programming language (Carpenter et al., 2016; Hoffman and Gelman, 2014), which is freely available, and facilitates the use of this model among researchers and practitioners.¹⁴

Regarding the identification of the model parameters, the demand and acquisition parameters ($\boldsymbol{\beta}_i^y$, σ^y , $\boldsymbol{\beta}_i^a$ and σ^a) are identified, provided the functional forms described in (1.1) and (1.2) are well specified. On the contrary, not every single parameter of the DEF component is fully identified. [*Lower layer*] The parameters that link the lower layer of the

¹³All details about the prior distribution $p(\mathbf{W}^y, \mathbf{W}^a, \mathbf{W}^1, \boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}_a)$ are presented in Appendix A.3.2.

¹⁴The code is available from the authors.

DEF with β_i^y and β_i^a are identified up to a rotation, similar to a traditional factor analysis model. Specifically, the scales of the lower layer trait (\mathbf{z}_i^1) and weights (\mathbf{w}^y and \mathbf{w}^a) are identified through the priors scales. Small rotations are identified by the sparsity of the ARD priors (see Appendix A.3 for details) — these priors favor the activation of fewer traits, avoiding the rotation of a large trait into smaller ones. Orthogonal rotations are not fully identified due to possible sign change in traits and label switching. However, we can obtain behavioral insights from the lower layer of model — e.g., what trait(s) are most predictive of specific behaviors — by carefully rotating the lower layer traits and weights parameters across draws to maintain a consistent interpretation of these parameters (see Appendix A.4 for details). [Top layer] The top layer of the DEF and the parameters that link the top and lower layer are not identified. This is similar to deep neural networks, in which the lower layer is a combination of the values of the upper layer and the weights linking them. In our model specification, this translates to the value of the top layer (\mathbf{z}_i^2) not being identified as different combinations of \mathbf{z}_i^2 and \mathbf{w}^1 could generate the same value for \mathbf{z}_i^1 . Most importantly, this lack of identification in the DEF component does not preclude the model from uniquely identifying the individual-level demand parameters β_i^y (as corroborated in Sections 1.4.4 and 1.5), which is the main goal when overcoming the cold start problem.

1.4.3 Model inferences for newly acquired customers

Recall that the main purpose of the model is to assist firms in the task of making inferences about how individual customers will behave in the future (e.g., how they will respond to marketing interventions), based on the observed behaviors at the moment of acquisition.

Intuitively, that process would work as follows: A new customer is acquired and the firm observes their behaviors at the moment of acquisition. At that point, and given the firm's prior knowledge of the market (i.e., the model parameters and market conditions), the firm makes an inference about that particular customer's latent traits, which are then used to infer the individual-level parameters that will determine their demand (e.g., how likely is it that the customer will purchase in the future, their responsiveness to marketing interventions).

More formally, we want to infer $p(\boldsymbol{\beta}_j^y | A_j, \mathcal{D})$ for customer j who was not in the training sample, for whom we observe acquisition characteristics A_j , and where $\mathcal{D} = \{y_{i1:T_i}, A_i\}_{i=1}^I$ comprises the calibration data. Denoting $\Theta = \{\boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \mathbf{W}^y, \mathbf{W}^a, \mathbf{W}^1, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}^a\}$ the population parameters and $\mathbf{Z}_j = \{\mathbf{z}_j^1, \mathbf{z}_j^2\}$, we can write $p(\boldsymbol{\beta}_j^y | A_j, \mathcal{D})$ by integrating out over the parameters Θ and \mathbf{Z}_j , and using the conditional independence properties of our model. That is,

$$\begin{aligned}
p(\boldsymbol{\beta}_j^y | A_j, \mathcal{D}) &= \int p(\boldsymbol{\beta}_j^y, \mathbf{Z}_j, \Theta | A_j, \mathcal{D}) \cdot d\mathbf{Z}_j \cdot d\Theta \\
&= \int p(\boldsymbol{\beta}_j^y | \mathbf{Z}_j, \Theta, A_j) \cdot p(\mathbf{Z}_j | \Theta, A_j) \cdot p(\Theta | A_j, \mathcal{D}) \cdot d\mathbf{Z}_j \cdot d\Theta \\
&= \int_{\Theta} \left[\int_{\mathbf{Z}_j} p(\boldsymbol{\beta}_j^y | \mathbf{Z}_j, \Theta, A_j) \cdot p(\mathbf{Z}_j | \Theta, A_j) \cdot d\mathbf{Z}_j \right] \cdot p(\Theta | A_j, \mathcal{D}) \cdot d\Theta \\
&\approx \int_{\Theta} \left[\int_{\mathbf{Z}_j} p(\boldsymbol{\beta}_j^y | \mathbf{Z}_j, \Theta, A_j) \cdot p(\mathbf{Z}_j | \Theta, A_j) \cdot d\mathbf{Z}_j \right] \cdot p(\Theta | \mathcal{D}) \cdot d\Theta. \quad (1.10)
\end{aligned}$$

The last equation suggests that if the number of customers in the calibration data is large, we can approximate the posterior of the population parameter with focal customer j by the posterior distribution obtained without the focal customer j . In other words, adding one

more customer would not significantly change the posterior of the population parameters. This approximation is very useful in practice because it allows us to draw from $p(\Theta|\mathcal{D})$ using the calibration sample, and draw the individual parameters of the focal customer j once this customer has been acquired, without the need to re-estimate the model to incorporate A_j . (See Appendix A.5 for a description of the corresponding algorithm.)

1.4.4 Model performance

Before applying the new modeling framework to the empirical context, we need to demonstrate the accuracy of the model at inferring the individual-level parameters for newly-acquired customers. Because individual-level parameters are, by definition, unobserved, we perform this task using a simulation analysis in which we know the exact values of β_j^y and can therefore evaluate the model's ability at recovering the true parameters using (1.10). Unlike other simulation exercises, the goal of this analysis is *not* to confirm that the model can recover the (population) parameters. Rather, we use simulations to demonstrate that the proposed model is able to recover customers' individual-level parameters accurately, even when the data generating process for those individual-level parameters is not known, and possibly different from the modeling assumptions. In reality, marketers (and researchers) never know the exact relationship between acquisition characteristics and future demand parameters, therefore, having a flexible model that performs well in a variety of contexts is of critical importance. (We briefly describe the main aspects of the simulation design while including all details in Appendix A.6.)

We generate three scenarios for the underlying relationship between acquisition variables and demand parameters. In each scenario, customers are “endowed” with a set of demand parameters that follow a specific relationship with their observed acquisition characteristics, namely (1) *linear*, (2) *quadratic/interactions* (allowing the relationship between one acquisition variable and the demand parameters to vary depending on the value of other acquisition characteristics), and (3) *positive-part* (forcing the relationship between acquisition characteristics and demand parameters being zero for low values of the acquisition characteristic). Given those individual-level demand parameters, customer transaction history is simulated for 2,200 customers. We use 2,000 customers to estimate the model, and the remaining 200 customers to evaluate the accuracy of the model at inferring demand parameters for newly-acquired customers. Specifically, only using the acquisition characteristics for these 200 customers, we use the model to infer their individual-level demand parameters, and compare those estimates with the true values.

We compare the performance of the FIM with that of three other specifications: (i) a HB-linear model, where individual demand parameters are specified as a linear function of the acquisition characteristics (this corresponds to the simulated data under the *linear* scenario), (ii) a full hierarchical model, where demand and acquisition parameters are jointly distributed according to a multivariate Gaussian distribution with a flexible covariance matrix, and (iii) a Bayesian PCA model. As discussed in Section 1.4.1.5, the Bayesian PCA model is a nested specification of the proposed FIM (in which the second layer does not exist) whereas the full hierarchical model and HB-linear specifications reflect alternative (simpler) ways in which past research has modeled these types of data. To measure the accuracy of each model, we compare the predicted posterior mean vs. the actual values for

the demand parameters (both the intercept and the effect of the covariates) of the 200 out-of-sample customers. Table 1.1 includes the results for all models across all scenarios.¹⁵ We also include the results of estimating a hierarchical Bayesian (HB) demand-only model in which acquisition characteristics are not incorporated, to have a reference of how much error one would obtain by simply predicting the population mean.

Table 1.1: Accuracy of predictions of demand parameters for (out-of-sample) customers

	Scenario 1		Scenario 2		Scenario 3	
	Linear		Quadratic/interactions		Positive part	
	R-squared	RMSE	R-squared	RMSE	R-squared	RMSE
<i>Intercept</i>						
HB demand-only	0.001	6.703	0.020	7.624	0.007	8.514
Linear HB	0.988	0.734	0.711	4.113	0.783	4.056
Full hierarchical	0.988	0.735	0.704	4.164	0.781	4.091
Bayesian PCA	0.988	0.736	0.706	4.484	0.780	4.329
FIM	0.988	0.738	0.888	2.661	0.928	2.987
<i>Effect of covariates</i>						
HB demand-only	0.005	2.562	0.004	4.589	0.001	4.604
Linear HB	0.986	0.303	0.258	3.969	0.736	2.363
Full hierarchical	0.986	0.303	0.258	3.970	0.733	2.378
Bayesian PCA	0.986	0.301	0.245	4.364	0.738	2.752
FIM	0.986	0.302	0.515	3.229	0.745	2.325

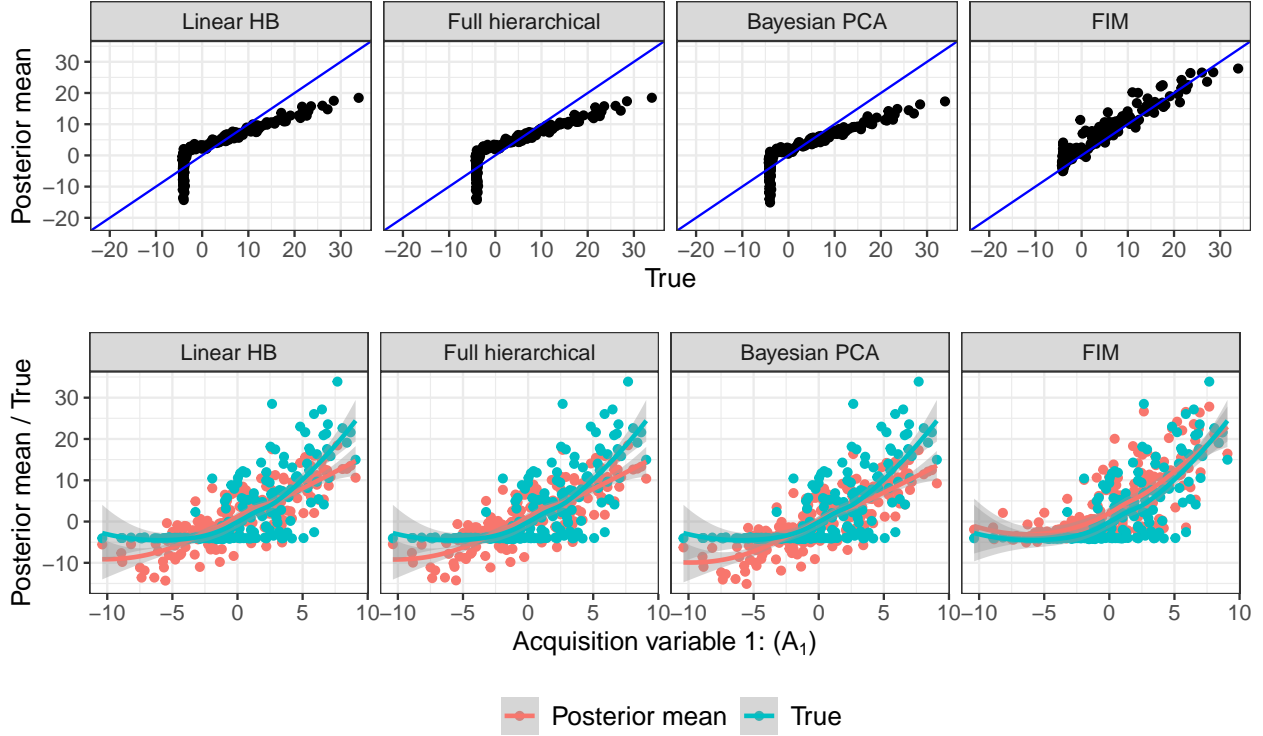
First, under a true linear relationship (Scenario 1), the FIM predicts the individual parameters as good as the benchmark models. The RMSE of the FIM is comparable to the benchmark models, and the R-squared is equal to the benchmark models. This result verifies that the FIM does not overfit the training data or, in other words, that the additional model complexity—even when not needed—does not hurt the accuracy of predictions for customers outside the calibration sample. Second, when the relationship among the model parameters is not perfectly linear (Scenarios 2 and 3), the FIM significantly outperforms the benchmark models in all dimensions. In particular, the R-squared of the FIM is higher than that of the benchmarks, demonstrating that the model is superior at sorting customers based

¹⁵See Appendix A.6.3 for more details about the specification of the benchmark models and Appendix A.6.4 for details on the performance metrics.

on their demand parameters. Moreover, the RMSE for the FIM is substantially lower than that of the benchmarks, indicating that the proposed model predicts the exact magnitude of customer parameters (e.g., purchase probability, sensitivity to marketing actions) more accurately than any of the benchmarks. These results hold when we examine the model “at scale”, when we significantly increase the amount of data collected by the firm and also add standard regularization techniques (e.g., LASSO) to the benchmark models. (Please see Appendix A.6.8 for details.)

To help understand what drives the greater accuracy of these predictions, we further explore the results for Scenario 3 (when the true relationship is positive-part). The first row of Figure 1.4 shows the scatter plot of the predicted ($\hat{\beta}_{j1}^y$) versus actual (β_{j1}^y) individual demand intercepts from each model, which displays the superior performance of the FIM, as detailed in Table 1.1. The second row of Figure 1.4 shows the predicted and actual demand intercepts as function of the first acquisition variable for each model. The blue dots show the true relationship between these two variables (i.e., positive-part) whereas the red dots correspond to the relationship estimated by the model. These plots evidence that the FIM can better recover the positive-part relationship between the acquisition variables and the demand parameters.

Figure 1.4: Model performance for Scenario 3: positive-part individual results of intercept. The first row shows the scatter plot of the individual true vs. posterior mean for each model. The second row shows the individual posterior mean (red) and true (blue) as a function of acquisition variable 1 (A_1).



Finally, to better understand which aspect of the model is responsible for this accuracy of predictions, we compare the BPCA and the FIM model more closely, allowing both specifications to vary the dimensionality of their latent components. Such an analysis indicates that the presence of the second layer of the DEF component is contributing significantly to the improvement in accuracy for scenarios where the relationship is not linear. The results suggest that incorporating that second layer, even if specified with low dimensionality, allows the model to flexibly capture the non-linear relationship between acquisition and demand parameters. (Please see full details in Appendix A.6.6.)

To sum, these analyses demonstrate the effectiveness of the FIM at overcoming the cold start problem. We have shown that the FIM can accurately infer customer parameters using only acquisition data, even when such a model is not used to simulate the true parameters. While the benchmark models fail to form accurate inferences of newly-acquired customers when the underlying relationships among variables are not perfectly linear, the FIM is flexible enough to reasonably recover those parameters. This latter point is of great importance because in reality the researcher/analyst never knows the underlying relationships among variables. Therefore, having a flexible model able to accommodate multiple forms of relationships is crucial to accurately infer customers' parameters.

1.5 Empirical application

1.5.1 Data and model specification

Our focal firm is an international retailer that sells its own brand of beauty and cosmetic products (e.g., skincare, fragrance, haircare).¹⁶ Customers can only purchase the company's products via owned stores, either offline (the company owns "brick and mortar" stores across many countries) or online (with one online store per country). While the company is present in many countries, most marketing functions (e.g., promotional campaigns, product introductions) are centralized and therefore operations are very consistent across markets. Like most other companies, the focal firm records the transactions of all individual customers, along with other information about the CRM activities, such as direct marketing campaigns and email marketing activities.

¹⁶The authors thank the Wharton Customer Analytics Initiative (WCAI) for providing this data set.

1.5.1.1 Transactional data

We obtain individual-level transactions for registered customers in the six major markets—USA, UK, Germany, France, Italy, and Spain. We observe customers from the moment they make their first purchase (starting in November of 2010). At the point of purchase, customers are asked to provide their name, email, and address so that they can receive promotions and other marketing communications from the firm. We track their behavior up to 4 years after that date (ending in November of 2014). We have 13,473 customers, with a minimum of 3 and a maximum of 51 periods of individual observations, resulting in 287,584 observations.¹⁷ During this time, we observe a total of 15,985 repeated transactions (i.e., the average number of transactions per customer is 2.19; or 1.19 repeated transactions). In addition to the behavior of the 13,473 registered customers, we collect data on all purchases made by “anonymous” customers in all six markets—i.e., those who never shared their identity with the firm. While their behavior is not included in our main analysis (the firm can neither track their future behavior nor communicate with them via email or mail), we use these anonymous transactional data to extract product-level information which will be used to augment the cold start data and to control for shocks in distribution channels that affect the timing of the introduction of new products in specific markets.

We specify demand as a logistic regression where $y_{it} = 1$ if customer i transacts at period t , and $y_{it} = 0$ otherwise. Specifically, $f^y(\cdot | \cdot)$ from (1.1) is defined as

$$p(y_{it} = 1) = \text{logit}^{-1} \left[\mathbf{x}_{it}^{y'} \cdot \boldsymbol{\beta}_i^y + \delta_{rec} \cdot \text{Recency}_{it} + \alpha_m \right], \quad (1.11)$$

¹⁷A period corresponds to exactly 28 days. We do not use a calendar month as our unit of analysis because we want to have the same number of days in all periods.

where we control for latent attrition using recency as a covariate (Neslin et al., 2013)¹⁸ and include market-level fixed effects to capture differences in purchase frequencies across countries (i.e., in this case $\tilde{\mathbf{x}}_{it}^y = [\mathbf{x}_{it}^y, \text{Recency}_{it}]$ and $\boldsymbol{\sigma}^y = \{\delta_{rec}, \alpha_1, \dots, \alpha_{M-1}\}$, with M the number of markets).

1.5.1.2 Marketing actions

The firm regularly sends emails and direct marketing to registered customers. The content of these promotional activities is set globally (i.e., the same promotional materials are used across countries, translated to the local language), though their intensity is set by market (e.g., the USA tend to send more emails than France).¹⁹ In addition to promotional activity, the company uses product innovation as a marketing tool. Like other major brands in this category, the focal retailer regularly adds extensions and/or replacements to their product lines. The sense among the company managers is that such an activity not only helps in acquiring new customers but also keeps current customers more engaged with the brand. When the company introduces a new product, it does so in all markets simultaneously. There is, however, some variation across markets regarding when new products were introduced. Conversations with the company confirmed that such variation is due to differences (and random shocks) in the local distribution channels.

While direct and email marketing are observed at the individual level (we denote them by **DM** and **Email**, respectively), the availability of new products is not observed at a granular

¹⁸As discussed in Section 1.4.1.1, the proposed FIM can accommodate different demand specifications such as “buy-til-you-die” models or HMMs. For our empirical application, we corroborate that adding recency is sufficient to control for latent attrition, which reduces the estimation time when compared with adding a probabilistic latent absorbing state (e.g., Chan et al., 2011).

¹⁹We only observe email activity sent after September 2012. Therefore, we will only consider customers acquired after that date for the estimation of the model.

level. We create a new product introduction variable (**Introd**) by combining point-of-sale data (at the SKU level) with a firm-provided SKU list of new products. Specifically, we obtain the list of all new products introduced during the period of our study. We identify the SKUs for all products in that list and infer inventory in each market from *all* purchases observed in that particular market (including all 304,497 transactions from “anonymous” customers). We assume that a new product was introduced in a market at the time the first unit of that SKU was sold. We then create a period/market-level variable representing the number of new products that were introduced in each market in each time period.

Table 1.2: Summary of time-varying marketing actions.

Marketing action	Statistic	Mean	SD	N
Email	Across observations	3.267	4.686	287,584
	Indiv. average	4.272	3.612	13,473
	Indiv. st. dev.	3.404	1.790	13,473
	Indiv. coeff. of variation	1.425	1.082	13,336
Direct Marketing	Across observations	1.006	1.889	287,584
	Indiv. average	1.329	1.018	13,473
	Indiv. st. dev.	1.731	0.769	13,473
	Indiv. coeff. of variation	2.031	1.205	13,455
Products introduced	Across observations	0.923	1.264	287,584
	Indiv. average	0.657	0.532	13,473
	Indiv. st. dev.	0.755	0.534	13,473
	Indiv. coeff. of variation	1.354	0.478	11,927

Table 1.2 shows the summary statistics for the marketing actions summarized across observations and across individuals. For the latter, we summarize individual average, individual standard deviation, and the individual coefficient of variation. The variation in these data is very rich both across customers and within customers.

We define the vector of demand time-variant covariates \mathbf{x}_{it}^y as the intercept, firm-initiated marketing actions, and seasonal factors such as holiday periods,

$$\mathbf{x}_{it}^{y'} = \left[1, \text{Email}_{it}, \text{DM}_{it}, \text{Introd}_{m(i)t}, \text{Season}_{m(i)t} \right]',$$

where **Email**, **DM**, and **Introd** are the marketing actions, and **Season** is a dummy variable that equals 1 for the winter holiday, and 0 otherwise.²⁰

Given the business nature of our application, the information provided by the firm about how the managers conduct their marketing actions, the rich longitudinal and cross-sectional variation in our data (Table 1.2), and our model specification, we argue that the potential endogenous nature of the marketing actions is not a main concern in this research (see Appendix A.7.1 for details). Nevertheless, in situations where these conditions do not hold (due to different strategic behavior by the firm or for data limitations), the demand model should be adjusted to account for the firm’s targeting decisions. Given the flexibility of our modeling framework, those adjustments would merely involve extending the demand model to capture unobserved shocks between firm’s actions and individual-level responsiveness (Manchanda et al., 2004) or adding correlations between firm decisions and unobserved demand shocks through copulas (Park and Gupta, 2012), depending on how these actions are determined by the firm. Those changes would only affect the demand (sub)model and not the overall specification of the FIM.

²⁰We compute such a variable for each market separately because the exact calendar time for the holiday period varies across countries. For example, in the USA the holiday “shopping” period covers Thanksgiving week until the last week of December (i.e., the end of Christmas), whereas in Spain the only holiday season corresponds to Christmas, which starts at the end of December and ends after the first week of January.

1.5.1.3 (Augmented) acquisition characteristics

Transaction characteristics: We compute **Avg.Price** as the total amount in euros of the ticket divided by the number of units bought at the first transaction; **Quantity** is the total number of units bought at the first transaction; **Amount** is the total amount in euros of the ticket at the first transaction;²¹ **Discount** is a dummy variable that equals 1 if the customer received discounts in the first transaction, and 0 otherwise; **Online** is a dummy variable that equals 1 if the first transaction was made online, and 0 otherwise. We also create a **Holiday** dummy variable that equals 1 if customer made their first transaction during the winter holiday period and 0 otherwise (analogously as the time-varying covariate **Season**).

Product characteristics: Directly from the observed product characteristics, we create a 10-dimensional vector that indicates whether the basket includes a product from a **Category**, including Body care, Face care, Hair care, Toiletries, etc., as defined by the focal company. Moreover, given that product innovation is very important in markets of beauty and cosmetic products, we create a **NewProduct** dummy variable that equals 1 if the customer bought a product that had been introduced in the 30 days prior to the purchase, and 0 otherwise. We also include the average **Size** of the packages in the basket, operationalized as relative size with respect to other products in the same sub-category, and a **Travel** dummy which equals 1 if the basket includes products on travel size, and 0 otherwise.

Latent representation of shopping baskets: As described in Section 1.3.2, we characterize each customer’s first purchase by computing moments of the products included in their shopping basket. The resulting product embeddings in our empirical application is a 6-dimensional

²¹We transform the **Avg.Price** and **Amount** variables using a log function, and the **Quantity** variable with a log-log function.

vector that represents the position of each product in a similarity space, which we call the “nature” of a product. Once those product embeddings are created, we create `BasketNature`, computed as the “average” product purchased, and `BasketDispersion`, computed as the element-wise standard deviation across products in the same basket, with missing values when the first purchase includes only one product.²²

Formally, the vector of acquisition characteristics is specified as follows,

$$A_i = [\text{Avg.Price}_i, \text{Quantity}_i, \text{Amount}_i, \text{Discount}_i, \text{Online}_i, \text{Holiday}_i, \\ \text{Category}_i, \text{NewProduct}_i, \text{Travel}_i, \text{Size}_i, \\ \text{BasketNature}_i, \text{BasketDispersion}_i].$$

The variation in the acquisition data is very rich (Table 1.3). For example, 22% of the sample was acquired over the holiday period, and 30% of first transactions included at least one discounted product, 35% included products in the face care category. The standard deviations of price, number of items purchased, amount, relative size, and basket dispersion are large, reflecting the heterogeneous behavior of customers across the six markets. Note that several of these acquisition characteristics are missing for some customers—for example, products for which the package size could not be retrieved from the data have missing `Package Size` observations, baskets that include single items have missing `BasketDispersion` observations, and so forth. These missing observations do not present a challenge in the estimation of the FIM— i.e., there is no need to eliminate observations or

²²In addition, if a first transaction of a customer includes only SKUs of products that were not purchased in any transaction of those anonymous customers’ transactions used for generating the product embeddings, then both `BasketNature` and `BasketDispersion` will have missing values as well.

to input population averages — because of the way the acquisition characteristics enter the probabilistic model in (1.2).

Table 1.3: Summary statistics of selected acquisition characteristics.

Variable	Description	Mean	SD	N
Avg. price (€)	Average price per unit, in euros	11.642	10.237	13,473
Quantity	Total number of units purchased	4.934	5.298	13,473
Amount (€)	Total ticket amount, in euros	39.567	38.433	13,473
Holiday	Whether customer was acquired during the Holiday	0.220	--	13,473
Discount	Whether discounts were applied in transaction	0.302	--	13,473
Online	Whether the transaction was online	0.176	--	13,473
New product	Whether a new product was purchased	0.431	--	13,473
Travel	Whether a travel-size product was purchased	0.397	--	13,473
Package Size	Average size of products (relative to its subcategory)	1.080	0.701	13,352
Avg. BasketDispersion	Average basket dispersion across all dimensions	1.338	0.660	9,928
Face Care	Whether a product in the Face Care category was purchased	0.352	--	13,473
Hair Care	Whether a product in the Hair Care category was purchased	0.120	--	13,473

Note: For the sake of simplicity, we omit the descriptive statistics for the 6 BasketNature variables and 8 remaining product categories. We also aggregate the BasketDispersion variables, by averaging across all dimensions of the *word2vec* representations. Missing values correspond to first purchases including products with missing information and for the case of BasketDispersion, those with only one item in the basket.

Table 1.4: Correlations among selected acquisition characteristics.

	Avg. price	Quantity	Amount	Size	Holiday	Discount	Online	New product	Travel	Face care
Avg. price	1.000									
Quantity	-0.330	1.000								
Amount	0.251	0.594	1.000							
Size	0.396	-0.238	0.038	1.000						
Holiday	-0.082	0.179	0.090	-0.027	1.000					
Discount	-0.200	0.285	0.184	-0.160	0.055	1.000				
Online	-0.241	0.411	0.168	-0.097	0.056	-0.049	1.000			
New product	-0.036	0.250	0.248	-0.055	0.068	0.066	0.106	1.000		
Travel	-0.350	0.347	0.122	-0.348	0.088	0.289	0.009	0.149	1.000	
Face care	-0.066	0.366	0.298	-0.113	0.051	0.096	0.483	0.177	0.083	1.000
Hair care	-0.124	0.261	0.121	-0.091	-0.016	0.084	0.266	0.139	0.063	0.155

Note: We dropped missing values in pairwise computations only.

Consistent with the challenges mentioned in Section 1.3.4, some acquisition characteristics are correlated with each other (Table 1.4) — e.g., customers who purchased many items paid less per item (correlation= -0.330), and those who bought on discount also paid slightly lower than those who paid full price when they were first acquired (correlation= -0.200). Online first purchases tend to include more items in the basket (correlation= 0.411) and contain products in the face care category (correlation= 0.483).

While it is to be expected that some of these variables will be correlated, as they capture different behaviors incurred by the *same* customer, some of these correlations might also arise from the market conditions at the moment in which a customer was acquired (e.g., if the company introduces all of its new products during the holiday, customers with `Holiday= 1` will also have `NewProduct= 1` and vice versa).²³ As discussed in Section 1.4.1.2, our modeling framework separates these two types of correlations by incorporating firm’s market-level actions, $\mathbf{x}_{m(i)\tau(i)}^a$, that potentially affect these acquisition behaviors.

Specifically, we include market-level CRM activities such as number of emails (`MarketEmail`), DMs (`MarketDM`),²⁴ and the number of products introduced by the firm (`Introd`) in that period.²⁵ That is,

$$\mathbf{x}_{m(i)\tau(i)}^a = \left[\text{MarketEmail}_{m(i)\tau(i)}, \text{MarketDM}_{m(i)\tau(i)}, \text{Introd}_{m(i)\tau(i)} \right]'$$

Because the span of the acquisition data covers 4 years from 6 different markets, we have substantial variation (longitudinal and cross-sectional) to separate any firm-related

²³If not accounted for, the latter case could be potentially problematic because the model would not be able to separate the predictive power of being a “holiday customer” from that of being a “new product customer.” And, if the company were to change its policy in the future (e.g., introducing new products in June), our model inferences about just-acquired customers could be misleading.

²⁴We calculate market-level number of emails and DMs as the average number of emails and DMs sent in a particular period to customers in that market. Note that the focal customer i cannot receive these marketing communications before being acquired, thus these variables are computed using the set of already existing customers at that time.

²⁵Note that the number of products introduced in a particular period enters both the demand and the acquisition model (\mathbf{x}_{it}^y and $\mathbf{x}_{m(i)\tau(i)}^a$, respectively). This is not problematic because the objective is different on each component. In the demand model, this variable captures the effect of introducing products at a particular period on the purchasing behavior of an existing customer for that particular period. In the acquisition model, this variable serves as a control for extracting the component of the acquisition variables that reflects individuals’ traits. For example, the fact that a customer bought a new product on their first transaction could be a signal of customers traits, and/or a consequence of more products being introduced by the firm when the customer was acquired.

systematic relationship among acquisition characteristics from correlations induced by customers' underlying preferences.

1.5.2 Estimation

We apply our modelling framework to this retail context to show how a firm can make meaningful inferences about newly acquired customers. The firm would do so by calibrating the FIM using historical data from its existing customers and making inferences about newly acquired customers for whom only the acquisition characteristics are observed.

We restrict our analysis to periods in which the firm was engaging in marketing activities, which span from October 2012 to November 2014 ($N = 8,985$ customers). In order to mimic the problem faced by the firm, we estimate the model with the transactional behavior of (existing) customers up to April 2014 and use those estimates to form first impressions for customers acquired after April 2014, using only their acquisition variables.²⁶ Specifically, we split all customers into three groups: *Training*, *Validation*, and *Test*. We randomly select customers that were acquired before April 2014 to use in our *Training* sample ($N = 5,000$) and use their behavior prior to April 2014 to train the models. Regarding the dimensionality of the FIM, and following the approach discussed in Section 1.4.1.4, we find that $N_1 = 13$ and $N_2 = 5$ are enough to recover the meaningful correlations present in our data. The posterior distribution of α is concentrated close to the origin for a set of lower level traits, indicating that $N_1 = 13$ is high enough to capture the traits that directly affect the demand and acquisition parameters. Similarly, the posterior

²⁶We chose this date to reasonably balance the amount of data we need to estimate the model, with the sample size remaining for the prediction analysis.

distribution of the computed pseudo- α shows that at least one upper level trait is not relevant for impacting the lower level traits, suggesting that $N_2 = 5$ is enough to capture the upper level traits.²⁷ (For further details see Appendix A.7.6.)

We also select another set of customers acquired during the same period for our *Validation* sample, which we will use to compare the predictive accuracy of the models at estimating demand ($N = 1,000$). Finally, we use the remaining customers acquired before April 2014, and combine them with those acquired after April 2014 to form our *Test* sample, which we will use to identify valuable customers and to inform our targeting policy ($N = 2,985$).²⁸

Similarly as in Section 1.4.4, we estimate all models (linear HB, Bayesian PCA and FIM) using NUTS in Stan.²⁹ We also estimate a set of probability models (also estimated with Stan) that have been proposed in the literature to model these type of data as they explicitly account for latent attrition (e.g., Chan et al., 2011; Schweidel and Knox, 2013; Schweidel et al., 2014). For completeness, we test multiple specifications varying the inclusion of time-varying covariates in the transaction process and time-invariant covariates in the attrition process, namely (1) Linear model with marketing actions + logistic attrition process (without acquisition covariates), (2) Linear model (without marketing actions) + logistic attrition with acquisition covariates, and (3) Linear model with marketing actions + logistic attrition with acquisition covariates (see details in Appendix A.7.2). Finally, we estimate two

²⁷For robustness, we estimate another FIM specification with $N_2 = 2$ instead, and we find that all upper traits are relevant, suggesting that $N_2 = 2$ may not be enough to capture the non-linear relationships present in the data.

²⁸Ideally, we would like to test our targeting policies using only customers acquired after the calibration period. However, given the low incidence of purchases in this empirical context, we would not observe such a group of customers for a long enough period to have reliable data to validate our predictions.

²⁹We do not show the Full hierarchical model given its similar performance to the linear-HB specification.

Machine Learning (ML) methods widely used for supervised learning (i.e., whether a customer transact) namely a feed-forward deep neural network (DNN) and a random forest (RF). Both ML models include time-varying covariates, acquisition characteristics, and market-conditions at the moment of acquisition. (See details in Appendix A.7.7 for details about the packages used for estimation of the ML methods and related model specifications.)

1.5.3 Results

1.5.3.1 Parameter estimates

Table 1.5 shows the population mean and standard deviation of each of the demand parameters. Customers in the sample have a low propensity to transact on average ($\beta_{intercept}^y = -3.110$). Email and direct marketing communications have a positive average impact on purchase ($\beta_{email}^y = 0.111$ and $\beta_{dm}^y = 0.121$, respectively), whereas product introduction effects are not significant on average. Finally, customers return to transact more on holiday periods ($\beta_{season}^y = 0.361$). In Section 1.5.4 we explore the observed heterogeneity in these components (captured by the FIM) as well as the implications for the managers of the firm.

Table 1.5: Parameter estimates of FIM.

Demand parameter		Posterior statistics			
		Post. mean	Post. sd	PCI 2.5%	PCI 97.5%
Intercept	Pop. mean	-3.110	0.051	-3.205	-3.024
	Pop. std. dev.	0.364	0.086	0.245	0.549
Email	Pop. mean	0.111	0.026	0.061	0.163
	Pop. std. dev.	0.167	0.031	0.110	0.235
DM	Pop. mean	0.121	0.028	0.067	0.174
	Pop. std. dev.	0.137	0.023	0.094	0.182
Product introductions	Pop. mean	-0.058	0.048	-0.164	0.024
	Pop. std. dev.	0.213	0.046	0.128	0.310
Season	Pop. mean	0.361	0.072	0.235	0.502
	Pop. std. dev.	0.362	0.065	0.245	0.505

Another set of interpretable parameters of the FIM are the posterior estimates of the lower layer of the DEF component. Properly rotated, these parameters could be used to interpret the latent factors that connect acquisition characteristics and demand parameters. For the sake of brevity, in this section we focus on the model performance at solving the cold start problem and include those interpretable results in Appendix A.7.3.

1.5.3.2 Comparison with the benchmark models

Unlike the simulation exercise, in the empirical application we do not know the true value of the demand parameters (β_i^y), and therefore have to rely on the model predictions to evaluate the quality of the model. We compare the (out-of-sample) accuracy of the FIM predictions with those of the benchmark models in Table 1.6.³⁰ (For completeness, the performance of all models on the *Training* sample is presented in Appendix A.7.4.) The FIM outperforms all the nested and latent attrition benchmarks in out-of-sample fit (i.e., Log-Like) as well as at making predictions at the observation, customer, and period level. This results not only corroborate the results presented in Section 1.4.4, now on a real-world setting, but also indicate that in this application, the traditional CLV models that explicitly model attrition do not outperform the Linear HB model with recency, even when including the acquisition variables as time-invarying covariates (e.g., Chan et al., 2011). Not surprisingly, the DNN method provide the most accurate results when looking at in observation level RMSEs, with

³⁰Arguably one should test these performance metrics on a different set of customers for which we selected the FIM specification. However, most FIM specifications deliver a similar performance on this Validation sample, and thus, would perform similarly well against the benchmark models. More importantly, the main performance test of the FIM is whether it can better identify valuable customers, which we perform using the Test sample in Section 1.5.4.

the FIM doing as well as the RF. However, when looking at customer- and period-level RMSE, the FIM outperforms all of the above models.

Table 1.6: Comparison with benchmark models (*Validation* sample).

Model	Log-Like	RMSE		
		Observation	Customer	Period
Linear HB	-2134.6	0.247	1.307	4.570
Latent Attrition w/ Acq.	-2367.4	0.249	1.403	4.951
Latent Attrition w/ Mktg. Actions	-2194.1	0.250	1.361	4.499
Latent Attrition w/ Acq.+Mktg. Actions	-2384.5	0.253	1.421	4.722
Bayesian PCA	-2010.0	0.240	1.184	4.240
Feed-Forward DNN	--	0.235	1.095	7.468
Random Forest	--	0.236	1.118	6.783
FIM	-1927.0	0.236	1.046	4.058

These analyses demonstrate that the FIM outperforms the benchmark models at accurately inferring individual-level demand parameters when only acquisition characteristics are available. The benefits of the proposed model are most salient when the underlying relationship between the acquisition characteristics and the parameters governing future demand are not linear, as it is the case for many empirical applications. In the next section we illustrate the managerial value of these predictions and discuss other insights (provided by the model) that are of managerial relevance.

1.5.4 Overcoming the cold start problem

First, we investigate how accurately the firm can identify “heavy spenders” using only the data from their first transaction. We do so by leveraging the information from customers in the *Test* sample. Specifically, we combine the estimates of the models (calibrated with the *Training* sample) and the acquisition characteristics observed for customers in the *Test* sample, and infer their individual-level demand parameters (see Appendix A.7.5) to predict each individual’s expected number of transactions. We then compare these inferences with

their actual behavior using two sets of prediction metrics (Table 1.7). First, we compute the RMSE on the individual-level average number of transactions per period.³¹ Second, based on each individual’s expected number of transactions, we flag whether a customer belongs to the top 10% and top 20% of highest average number of transactions and report the proportion customers correctly identified/classified in each group.³² For reference, we compare those figures with what a random classifier would predict (shown in the last row).

Table 1.7: Identifying valuable customers using *Test* customers.

Model	RMSE	% customers correctly classified	
		Top 10%	Top 20%
Linear HB	0.157	0.151	0.253
Latent Attrition w/ Acq.	0.520	0.113	0.207
Latent Attrition w/ Mktg. Actions	0.303	0.213	0.248
Latent Attrition w/ Acq.+Mktg. Actions	0.242	0.090	0.191
Bayesian PCA	0.138	0.208	0.313
Feed-Forward DNN	0.098	0.349	0.450
Random Forest	0.106	0.193	0.310
FIM	0.131	0.401	0.477
Baseline (random)	–	0.100	0.200
	–	(0.067,0.127)	(0.170,0.230)

Note: The proportion of top spenders is computed by predicting over the observed periods, computing the average number of transactions per period, and selecting customers with highest predicted values.

As Table 1.7 shows, the FIM can predict reasonably well the value of customers: the FIM has a lower RMSE than the Linear HB and the Bayesian PCA models, only outperformed by the RF and the DNN. Moreover, Linear HB and BPCA are significantly better than the baseline at identifying valuable customers, which proves that acquisition characteristics carry valuable information to predict the value of customers. Nevertheless, the FIM significantly improves the identification of valuable customers over the benchmark models, including the DNN, being able to correctly identify 40.5% of customers in the Top 10% and 47.7% of customers in the Top 20%. These results are consistent with the notion that, because the FIM captures the non-linearities in the relationship between acquisition

³¹Using our notation, the individual level average number of transactions per period is $\bar{Y}_i = \frac{1}{T_i} \sum_{t=1}^{T_i} y_{it}$.

³²We make predictions and compute recovery rates for each draw of the posterior distribution and report posterior means and 95% CPI.

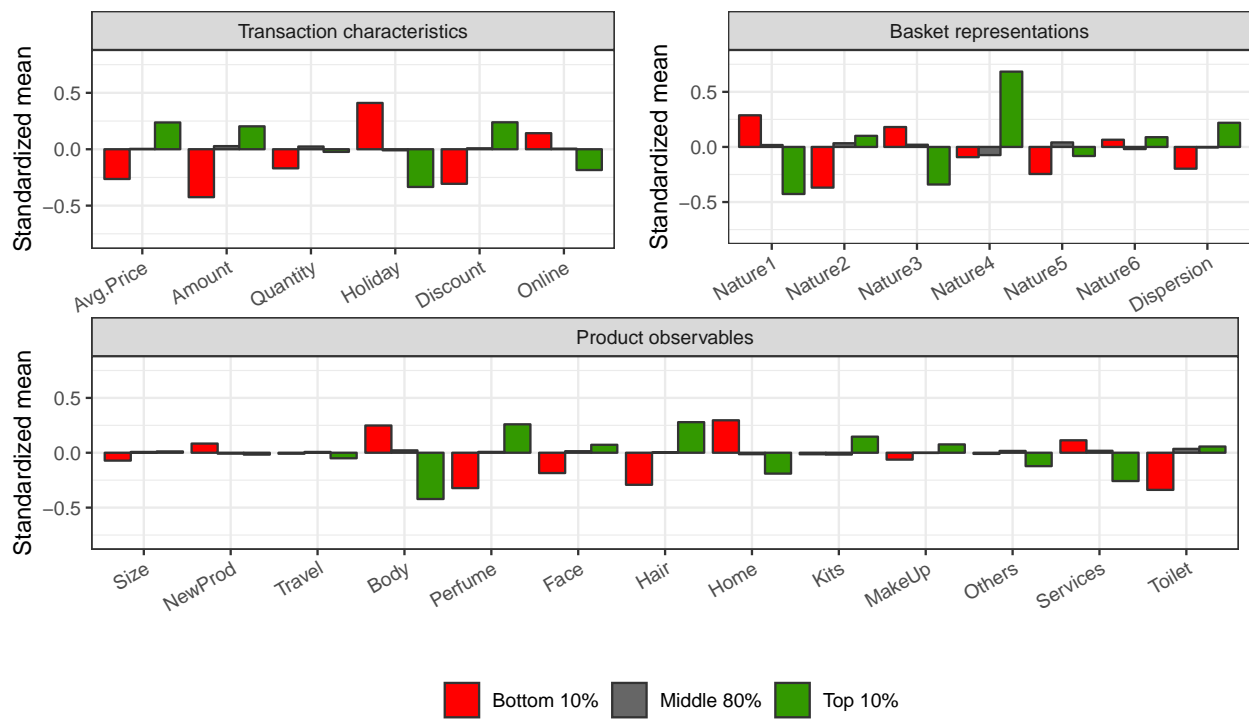
characteristics and future demand parameters, it does an excellent job—significantly better than the benchmarks—at sorting customers based on their expected value inferred from their acquisition characteristics.

Similarly, a firm would use the FIM to identify which customers are the most sensitive (or least sensitive) to marketing interventions; information that will be instrumental in increasing the effectiveness of its marketing actions (e.g., Ascarza, 2018). Unfortunately, our data does not enable us to quantify the exact value that the focal firm could extract from a FIM-based targeting approach—ideally, one would run a field experiment to test the effectiveness of targeting policies based on the predictions of the FIM. Nevertheless, combining the results from Section 1.4.4, where we demonstrate the model’s ability to predict the (individual-level) demand intercept as well as the sensitivity to the covariates, with the results in Table 1.7, where we corroborate some of those findings in our empirical application, we are confident that implementing targeting policies based on predictions of the FIM would generate incremental revenues to the firm. We trust that future research will be able to quantify these benefits empirically.

Second, we use the FIM results to explore the acquisition variables that better characterize “heavy spenders” (separately from light users), customers with “high sensitivity to email” (from those who are better left out in the email campaigns), and those who are “most sensitive to direct marketing” campaigns. Based on the model predictions, we split customers from the *Test* sample in three groups: Top 10%, Middle 80% and Bottom 10% for each of the three categories and summarize the average value of each of the (standardized) variables observed at the moment of acquisition. Figure 1.5 shows the results when sorting customers on the basis of expected future value. Several interesting findings emerge:

Consistent with the patterns observed when exploring the predictive power of the acquisition variables (Figure 1.1) we find that the Top 10% heavy spenders are less likely to be acquired during the holiday period, more likely to being acquired offline, and tend to buy expensive and discounted products in their first purchase, compared to those at the Bottom 10%. They are also characterized to buy certain types of products, as indicated by the high chance to include Perfume and Hair products in their first transaction (less likely to contain products in the Body Care, Home and Services categories), as well as by a high score in dimension 4 of the product embeddings.³³

Figure 1.5: Acquisition characteristics for customers with top/middle/low CLV.



We repeat the analysis now sorting customers based on their predicted sensitivity to email (Figure 1.6) and predicted sensitivity to DM (Figure 1.7). Consistent with the

³³This dimension is related to products such as “Grape Line Showers” and “Olive Harvest Conditioner,” see Table A1 in Appendix A.1.

previous findings, several acquisition characteristics exhibit a non-linear relationship with the sensitivities to marketing actions. Both the Top 10% and Bottom 10% email sensitivity groups are less likely to buy in the Body Care category during their first transaction, compared with the remaining 80% of customers in between. Customers who are the most sensitive to email marketing are more likely to be acquired online, buy less expensive products, and fewer units at their first purchase. With respect to DM, low sensitive customers buy fewer units and more expensive products in their first transaction, while high sensitive customers are more likely to buy relatively small sized products, recently introduced products, and products in the Perfume Category at their first purchase.

Figure 1.6: Acquisition characteristics for customers with top/middle/low sensitivity to Email.

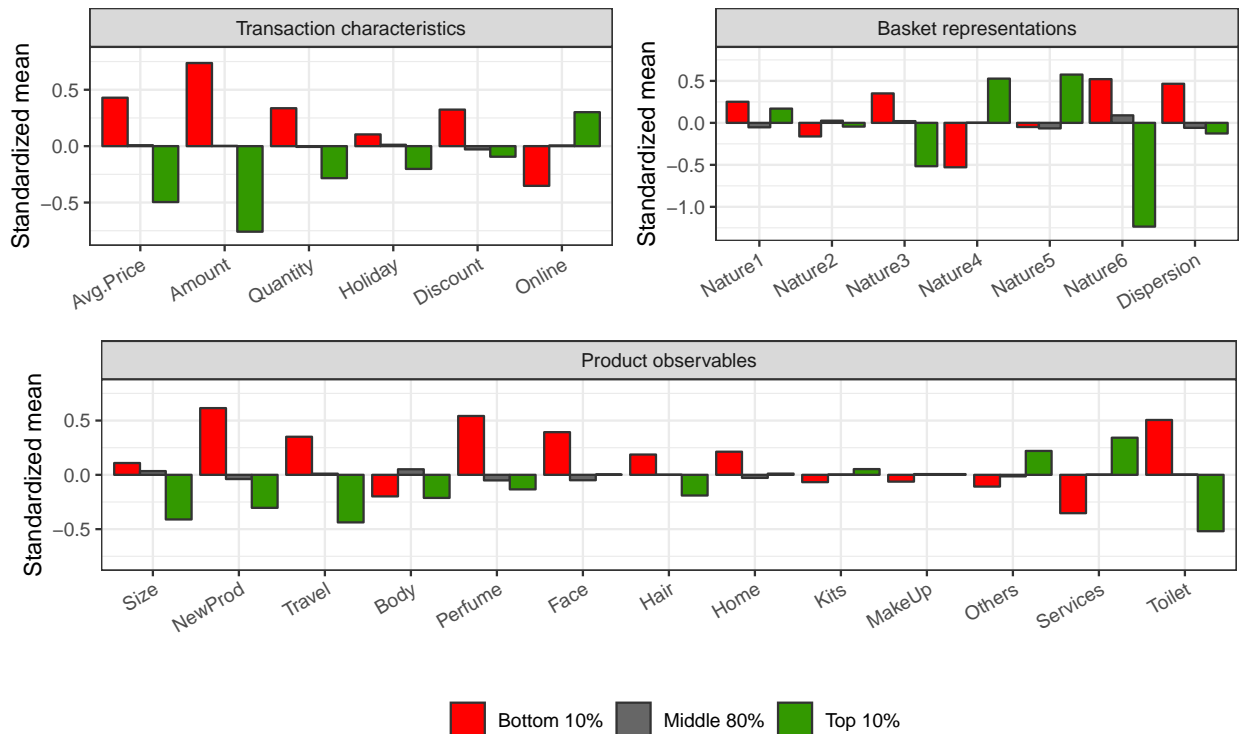
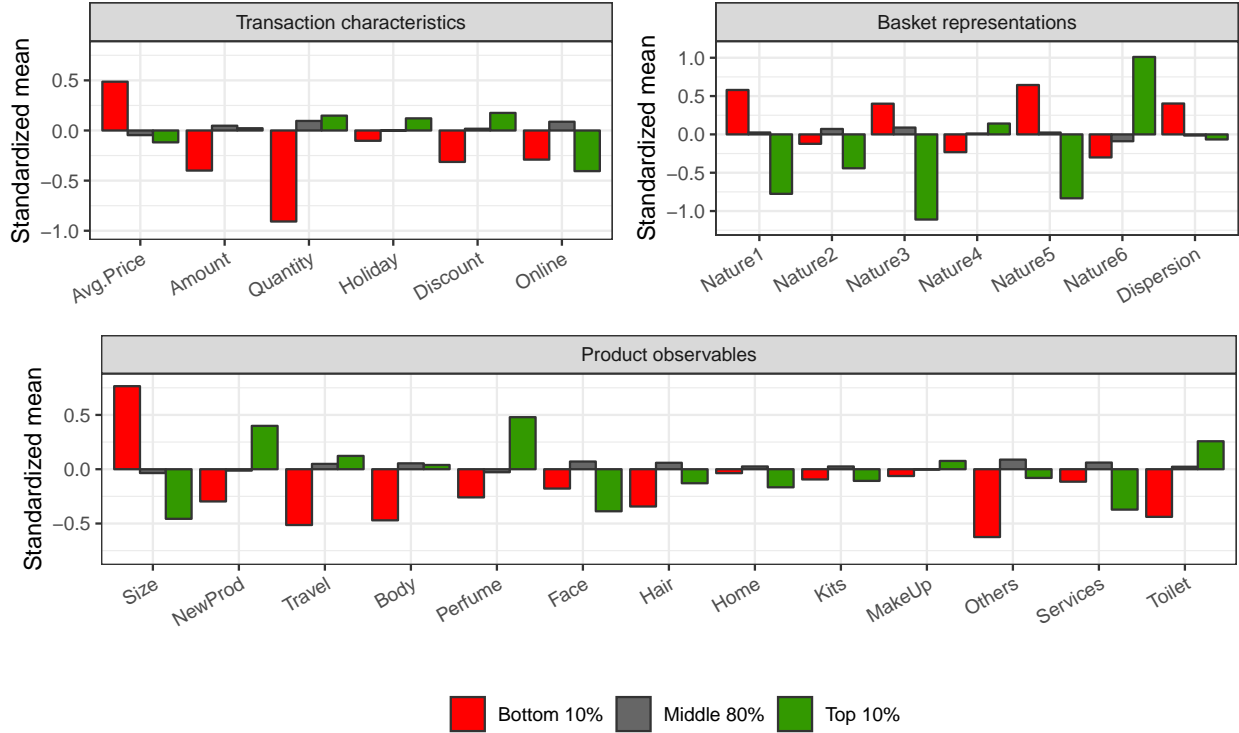


Figure 1.7: Acquisition characteristics for customers with top/middle/low sensitivity to DM.

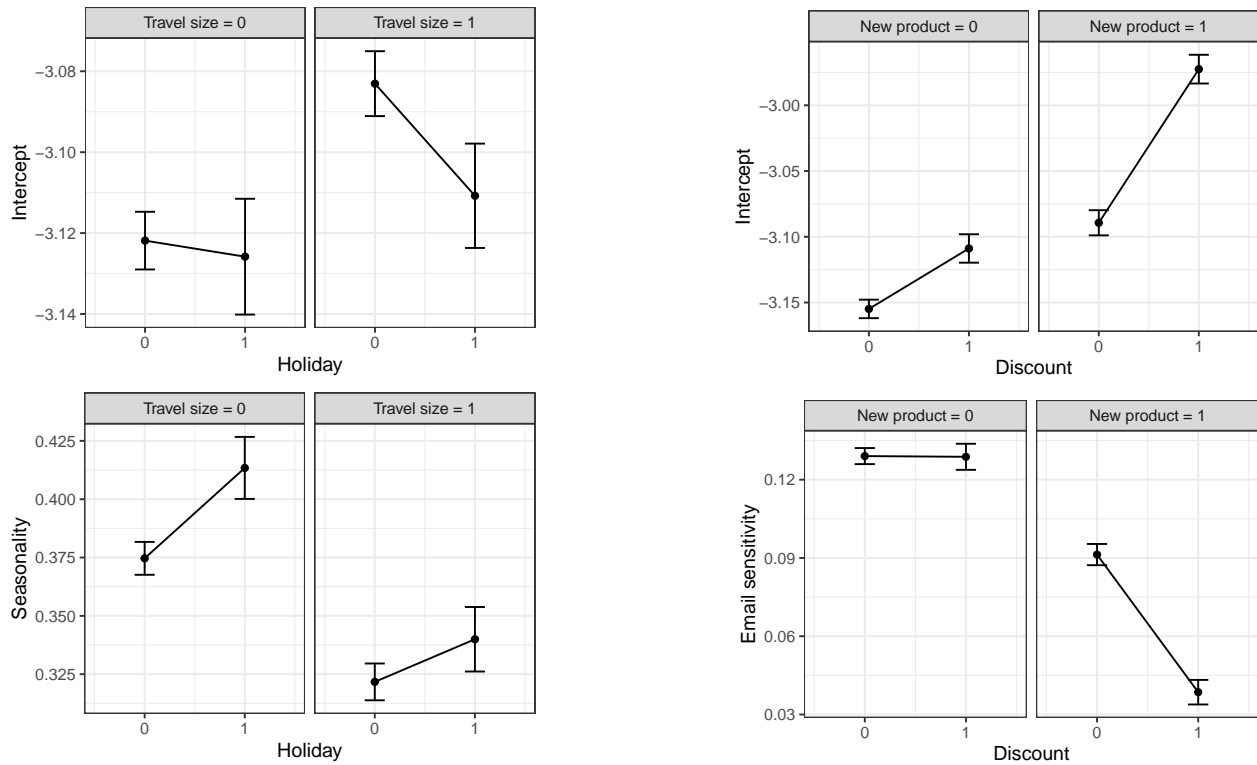


Finally, we use the inferred demand parameters from these test customers to explore the relationships between the magnitude of the demand parameters and the acquisition characteristics. Figure 1.9 shows the individual level posterior mean of the demand parameter vs. the acquisition characteristics for a set of demand parameters and acquisition characteristics. In particular, we find that these plots corroborate that there are non-linear relationships that the model allows to uncover.³⁴ Figure 1.8 explore possible interactions by presenting box plots of individual level posterior mean demand parameters and pairs of discrete acquisition characteristics. The model replicates the model-free insights shown in Figure 1.2: (1) the relationship between the intercept and whether the customer was acquired during the winter holiday season (`Holiday`) depends on whether the customer

³⁴Note, that these plots show marginal relationships of demand parameters and acquisition characteristics (i.e., one at a time) where indeed the model cover relationships accounting for all acquisition characteristics.

purchased a travel-sized product (**Travel Size**), and (2) the relationship between the intercept and whether the customer purchased discounted products at acquisition (**Discount**) depends on whether the customer purchased a recently introduced product (**New Product**). Moreover, the model not only captures these relationships for the intercept but also for other demand parameters. For instance, the holiday season lift is higher for customers that were acquired during a past holiday season compared to those that were not, but this difference is considerably larger for those that did not purchased a travel-sized product when acquired. Also, the differences in email sensitivities across customers that received discounts on their first purchase only exist for those who purchased a recently introduced product at acquisition.

Figure 1.8: Demand parameters (posterior mean) vs. some binary acquisition characteristics.



1.6 Conclusion

We have developed a modeling framework (FIM) that, leveraging information collected when customers are acquired, enables firms to overcome the cold start problem of CRM. Using a probabilistic machine learning approach, the model connects underlying acquisition and demand parameters using a set of hidden factors modeled via deep exponential families. The multi-layer structure with flexible relationships among layers enables the researcher or analyst to be agnostic about the (assumed) underlying relationship among variables. The hidden factors automatically extract relevant information from existing data — i.e., identify the traits that relate acquisition characteristics with future outcomes — overcoming the challenge (commonly faced by firms) of maintaining significant amounts of redundant and irrelevant data in their customer databases.

We have illustrated the benefits of using the FIM in a retail setting. First we have shown how the focal firm can further leverage its existing database to augment the cold start data using readily-available techniques. We have further demonstrated how subtle signals extracted from the augmented data by the FIM enables the focal firm to make individual-level inferences about just-acquired customers, for example, distinguish high-value customers from those unlikely to purchase again and those most and least sensitive to marketing interventions, such as email campaigns or direct marketing. We leverage the model predictions to identify characteristics of first transactions that are predictive of customer behavior in future periods. For example, compared to the rest, Top 10% heavy spenders are more likely to be acquired online and their first purchases to be expensive and

discounted products, and customers identified as most sensitive to email marketing to also be more likely to be acquired online but buy less expensive products, and their first purchases to be of fewer units.

These findings suggest that firms can meaningfully categorize customers based on characteristics of their first transactions. We believe this approach to customer segmentation to be promising in relying neither on sometimes difficult to obtain customer-provided data (Dubé and Misra, 2017) and nor on external sources of data that could pose privacy concerns. The resulting insights can be used both to prune acquisition data and inform decisions about the types of variables worth collecting from customers that make a first transaction or first visit a company’s website. Our research shows that firms leave value on the table by not fully leveraging the multiple behaviors observed when a customer makes a first transaction, and provides a general framework for extracting meaningful but hard-to-pinpoint relationships imprinted in subtle ways in “cold start” data.

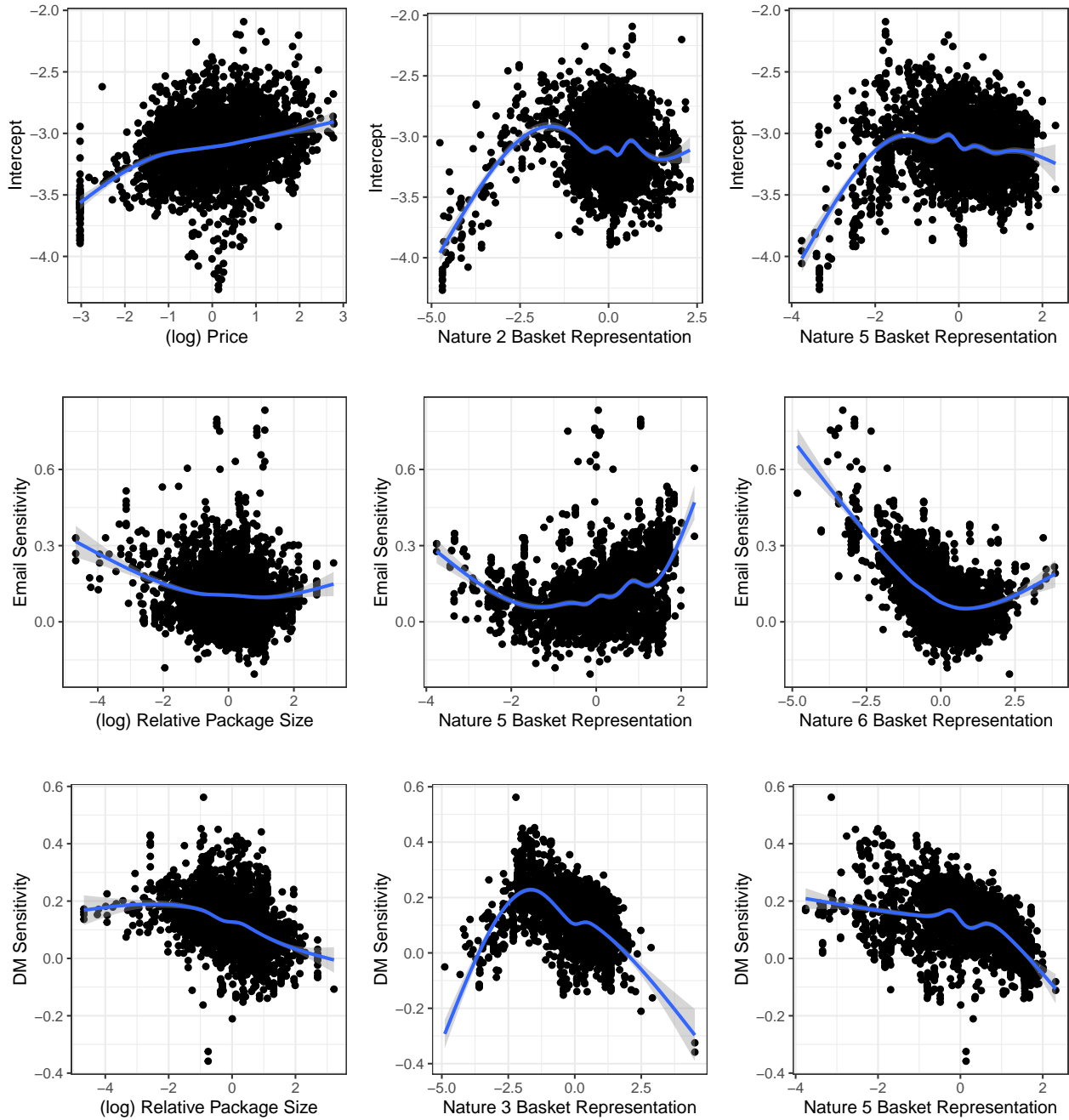
While this research highlights the value of using the FIM to tackle the cold start problem of CRM, it is also important to acknowledge some limitations of the present research. The simulation analyses enabled us to validate the accuracy of the model at inferring individual-level parameters, but doing so in an empirical setting, in which only realized purchases are observed, is more difficult. We leave it to future research to examine and quantify the effectiveness of targeting policies based on the predictions of the FIM. Regarding the model specification, we investigated model performance using linear and logistic specifications for the demand and acquisition models. Although the proposed FIM is extremely flexible so as to be adaptable to other modeling frameworks, we have not empirically tested the model’s performance in more complex structures. The current model

estimation is computationally feasible for datasets with thousands of customers, dozens of time periods, and a handful of variables (as in our empirical application). Although the model scales readily to situations with more acquisition variables, increasing the sample size to, for example, millions of customers will increase estimation time substantially, constraining the ability to gauge customers' first impressions in a timely manner. For such cases, variational inference might be a better way to estimate and use the model.

A natural extension to this research would be to investigate a wider range of acquisition characteristics and the relevance thereof to customers' first impressions in different contexts. The results of our empirical application could be built on to further augment the data from first purchases and incorporate other acquisition characteristics that, although not currently collected (e.g., whether the customer visited the store alone or with family), could be valuable in identifying which marketing actions are most likely to increase future sales. We encourage further research to investigate these research settings and identify additional drivers and methods that might help companies overcome the cold start problem.

The main goal of this work being to provide a flexible model that overcomes the cold start problem, we have not formally investigated the latent traits that drive all the observed behaviors. It would be relevant for researchers and marketers to identify individual traits that characterize shopper behavior, to which end customer behavior in a variety of contexts could be measured and estimated in a unifying FIM framework. We hope that this research opens up new avenues for understanding "universal" shopping traits and identifies the behaviors that best relate to those generalizable findings.

Figure 1.9: Demand parameters (posterior mean) vs. some continuous acquisition characteristics.



Chapter 2

The Customer Journey as a Source of Information

This essay forms the basis of a working paper of the same name jointly authored with Eva Ascarza and Oded Netzer.

Abstract

In high involvement purchases such as flights, insurance, and hotel stays, firms often observe at most only a handful of purchases during a customer lifetime. The lack of multiple past purchases presents a challenge for firms to infer individual preferences. Moreover, customers in these industries often look for products that satisfy different needs depending on the context of the purchase (e.g., flights for a family vacation vs. flights for a business trip), further complicating the task to understand what a customer might prefer in the next purchase occasion. Fortunately, in such high involvement purchases, these settings also collect other pieces of information; prior to a purchase, firms often have access to rich information on the customer journey, over the course of which, customers reveal their journey-specific preferences as they search and click on products prior to making a purchase. The objective of this essay is to study how firms can combine the information collected through the customer journey — search queries, clicks and purchases; both within-journeys and across journeys — to infer the customer’s preferences and likelihood of buying, in settings in which there is thin purchase history and where preferences might change from one purchase journey to another. We build a non-parametric Bayesian model that links the customer clicks over the course of a journey, and across journeys, with the customer’s history of purchases. The model accounts for what we call *context heterogeneity*, which are journey-specific preferences that depend on the context in which the journey is undertaken. We apply our model in the context of airline ticket purchases using data from one of the largest travel search websites. We show that our model is able to accurately infer preferences and predict choice in an environment characterized by very thin historical data. We find

strong context heterogeneity across journeys, reinforcing the idea that treating historical journeys as reflecting the same set of preferences may lead to erroneous inferences.

Keywords: Customer Journey, Bayesian Non-parametrics, Clickstream Data, Customer Search.

2.1 Introduction

In purchase of high involvement products such as flights and hotels stays, firms often observe at most only a handful of purchases during a customer lifetime. The short history of past purchases presents a challenge for firms who want to infer customer preferences; in particular, whether the customer will buy, and if so, what product will they buy. Moreover, in these settings, even when firms observe multiple purchases per customer, it is often the case that different purchase occasions are aimed at satisfying different customer needs (e.g., leisure versus business travel). As a result, it is not obvious how to aggregate information across purchase occasions in a meaningful way. To mitigate the thin historical data on the customer's past purchases, firms in these contexts often have access to rich information about the customer journey prior to purchase. In particular, firms not only observe the initial search query, but also the clicking steps the customer makes until making purchase.

We argue and demonstrate that, in these settings, firms can use the customer purchase journey — search queries, clicks and purchase the customer makes while in the market for a product — as a source of information to overcome both the lack of historical purchase data and to account for context changing preferences. During the course of a purchase journey, the customer reveals information in two ways. First, the customer types the search query, hence allowing the firm to infer the particular need the customer is looking for in the current purchase journey. Second, as the customer clicks on certain products but not on others, the customer starts to reveal his/her more stable as well as purchase journey-specific preferences. For example, if Adam is searching for a flight from Chicago to

Orlando, observing that he is adding children to the query may provide information that this flight will be purchased for family vacation, which may inform about Adam's stronger preferences for non-stop flights as he wants to avoid making connections with the kids. Then, as Adam clicks on a non-stop American Airlines flight the firm may infer Adam's preferences towards American Airlines, especially if Adam's purchase history is short. The firm can use this information to continue showing relevant products to Adam. Moreover, even if Adam decides to wait and not buy in that moment, the firm can use this information to recommend certain flights through re-targeting efforts, and/or to show these flights at the top of the page the next time Adam searches for the flight from Chicago to Orlando.

Accordingly, the objective of this paper is to study how firms can use the customer journey path from search to transaction as a source of information and combine it with, possibly one or a few, past journeys to infer the customer's preferences and likelihood of buying. To do so, we build a non-parametric Bayesian model that links the customer query with the clicks over the course of a journey, and integrates that information across journeys, and across the customer's history of purchases.

The model accounts for what we call *context heterogeneity*, which are journey-specific preferences that depend on the context in which the journey is undertaken. We model the journey decisions on what to search, what to click and what to buy to be both a function of customers' stable preferences and the unique needs of the context of the trip. Intuitively, contexts are unobserved segments that capture need-specific preferences that are shared for different journeys across customers. We uncover those contexts non-parametrically using a Pitman-Yor process as prior for the distribution of contexts in the population. Our model allows for creation of new contexts that have not been previously observed as new journey

observations arrive. Our model also allows for preferences over products to express differently when customers click early on in the process versus when they choose to buy.

To infer journey-specific preferences, our model leverages three sources of information to estimate the customer’s preferences in each particular purchase journey: (1) within journey’s behavior (e.g., the customer query and what the customer search for and clicked on in the focal journey); (2) past journeys’ behavior (e.g., what the customer clicked on and purchased in past journeys); and (3) across customers’ behavior (e.g., what other customers with similar search behavior clicked on and purchased). The within journey information, particularly click information, allows us to identify the unique journey-specific preferences. The past journeys’ information (not only past purchases but also searches and clicks) allows us to inform the customer’s stable preferences. Finally, the information across customer augments the, possibly thin, information from the two other sources with data from a host of customers with similar context-specific preferences. Thus, in an environment with infrequent purchases, but with observed interactions throughout the customer journey, our model allows us to leverage across customer information with within journey information to augment the relatively thin or non-existent historical data.

We estimate our model in the context of airline ticket search and purchases, using data from one of the largest online travel website. For each journey we observe the query search, and, if made, clicks and a purchased flight for a sample of active customers who searched for flight tickets from May 2017 to November 2017. These customers have on average fewer than 6 journeys each, result in 0.81 purchases per customer on average. Note that even when the journey does not end up in a purchase, these journeys contain valuable information as customers click on products, providing signals of their preferences.

We find strong context heterogeneity across journeys such that over time, customers search for multiple contexts, reinforcing the idea that a model that treats historical journeys as reflecting the same set of preference may lead to erroneous inferences. While we find that customers have negative sensitivity to price, as well as to the number of stops or the length of the flights, the extent to which customers care about these attributes heavily depends on the context. We uncover 19 different contexts for customers’ travel, with an average of 3.30 contexts per customer for customers with more than one journey. Those contexts vary in terms of their search queries (e.g. who is the customer flying with, where and when is the customer flying, and when is the customer searching), as well as their preferences for product attributes (e.g. price, number of stops, length of flight and departing and arrival times). The different contexts capture different trip purposes, and customers exhibit different preferences depending on the purpose of their trip. For example, a customer who searches for short domestic business trips—characterized by a flight that is less likely to include a weekend, without children or multiple adults, and not searching far in advance of the departing date—is less price sensitive and prefer return flights that arrive in the evening. In a different context, customers who search for long distance vacations with their families are more price sensitive, have stronger preferences for non-stops, and strongly prefer avoiding a return flight that departs between midnight and sunrise.

We compare our full model that accounts for context heterogeneity, past journeys and the information collected during the current journey to a host of nested models that do not consider some of these components. We show that leveraging the customer journey as a source of information helps the firm predict more accurately whether the customer is going to purchase. Moreover, for customers who buy, our model also outperforms the benchmark

models at predicting the type of flight the customer is likely to purchase. We also show the benefits of our model in cases in which the identity of customers cannot be tracked/stored, for example due to privacy reasons. We show that when queries and clicks are observed, even a model that cannot identify a customer identity (i.e. treating every customer as if s/he were a new customer) can alleviate the lack of purchase history to estimate preferences using context information inferred from search queries and clicks. This benefit of our model is particularly relevant given the recent concerns regarding consumer privacy, which often limits firms' ability to store historical data at the individual level.

Beyond our empirical application (i.e., travel websites), our model can be useful in other industries, with high involvement purchase and involved customer journeys such as cars and durable goods, which also exhibit thin individual purchase history. Experiential purchases such as hotel stays, restaurants reservations, food delivery and media consumption often involve purchase journeys with varying contexts and needs. We believe that our empirical setting is one of the complex situations as it contains *all* these aforementioned challenges one needs to address to infer preferences accurately at the journey level. Sub components of our model could be generalizable to other settings in which only some of the challenges take place.

The rest of the paper is organized as follows. We start in Section 2.2 by providing a review of relevant literature. We describe our empirical context and the data in Section 2.3. Then, in Section 2.4 we develop our modeling framework to integrate the customer journey as a source of information to infer preferences. We show the insights from our model as well as its prediction performance in Section 2.5. Finally, we conclude discussing the generalizability of our modeling approach, as well as potential limitations and future directions.

2.2 Relevant literature

The current work contributes to the rich literature in marketing on using transactional and search data to estimate customers' preferences at the individual level. On the one hand, the early days of scanner panel data saw the marketing literature developing methods that allow researchers to use past-purchase data to infer individual-level preferences (Rossi et al., 1996; Allenby and Rossi, 1998; Duvvuri et al., 2007; Fiebig et al., 2010). These panel data models have been widely used by researchers and practitioners in settings where individual transactions are available. However, there are many business contexts in which observing several purchases by the same customer is rare (e.g., purchasing a car, or booking a week-long vacation), preventing these models from estimating individual-level preferences in a reliable manner, and therefore limiting the manager's ability to understand and leverage customer heterogeneity. We extend this literature by incorporating the information from the customer journey (e.g., search and purchase processes observed from clickstream data) from current and past purchase occasions; even from those that did not end up in a purchase.

There is a rich literature on consumer search and the use of clickstream data (e.g., Montgomery et al., 2004; Kim et al., 2010, 2011, 2017; Ghose et al., 2012, 2014, 2019; Seiler, 2013; Koulayev, 2014; Honka, 2014; Bronnenberg et al., 2016; Honka and Chintagunta, 2017; Chen and Yao, 2017; De los Santos and Koulayev, 2017; Ursu, 2018) that uses within-journey information to infer customer preferences and to predict purchase. For example, Montgomery et al. (2004) find that customers' browsing behavior can predict future steps in the browsing process as well as conversion. De los Santos and Koulayev (2017) shows how

firms can use data on the current visit to optimize click through rates. Nevertheless, these studies analyze only the focal journey ignoring the information provided by previous journeys (for exception see Dong et al., 2019), limiting the model’s ability to capture rich heterogeneity in customer preferences. In turn, in most of the aforementioned search models, and partially due to model complexity of these models, unobserved heterogeneity is often taken into account in a fairly limited manner and is mainly used to unbiasedly account for substitution patterns in the market that are caused by heterogeneous tastes, rather than capturing the rich customer heterogeneity both within and across journeys. We extend this literature by providing a method that integrates within-journey information (i.e., search queries and clicks from the focal journey), cross-journey information from multiple journeys by the same customer (i.e., past search queries, clicks and purchases), and journeys from other customers (i.e., search queries, clicks and purchases from other customers).

There is a fundamental challenge that arises when combining information across journeys; it is often the case that different purchase occasions are aimed at satisfying different needs, resulting in customers exhibiting situational-based preferences (Belk, 1975; Dickson, 1982; Holbrook, 1984; Bucklin and Lattin, 1991; Jacobs et al., 2016; DeSarbo et al., 2008; Liu and Dzyabura, 2017; Thomadsen et al., 2018).¹ For example, a customer looking for a business trip might exhibit different preferences than when searching for a family vacation. As a result, when inferring customer preferences for a focal journey, it is not clear how we should integrate the information from within-journey and across journeys in a meaningful way. At first glance, one would argue that within-journey clicks are more informative than

¹Note that situational-based is distinct from dynamics in preferences (e.g., Erdem and Keane, 1996; Netzer et al., 2008) or from learning models (e.g., Dzyabura and Hauser (2019)) in which customer preferences may change longitudinally in a systematic way. In our case, situational-based preferences are affected by the context of each purchase occasion with no particular (or systematic) longitudinal pattern.

clicks from past journeys; and that past purchases are more relevant than past clicks.

However, how exactly a model should combine all this information is far from obvious.

This paper combines these three research streams and contributes to the literature by proposing a method to infer customer preferences in settings where there is thin purchase history—i.e., most customers have not purchased multiple times, some customers have not even purchased yet—and where preferences might change from one purchase occasion to another. We do so by jointly modeling information on the full customer purchase journey: search queries, clicks and purchases; both within-journeys and across journeys.

Our work also relates to the literature on context-dependent product recommendations (e.g., Sarwar et al., 2001; Hidasi et al., 2016; ?; Yoganarasimhan, 2019). This growing literature in the areas of computer science as well as in marketing has proposed diverse machine learning approaches—including item-to-item recommendation approaches using similarities across products, Recursive Neural Networks, or topic modeling—to recommend products when there is lack of historical individual-level data. Most of these methods require the observation of several individuals interacting (e.g., clicking or buying) with the same set of products, as well as each individual interacting with several products. Our approach relaxes this stringent requirement, as we extract preferences for attributes and not only for “entire” products. As a result, our model can be applicable when the product space is large and includes non-purchased items, and when the number of available products is growing over time. Additionally, in our application, the context is determined not only by the product but also by the search occasion environment (e.g., a journey in which a customer is buying a flight for tomorrow might differ from a journey in which s/he is looking for the same destination, but purchasing 2 months in advance).

Methodologically, our work builds on and contributes to the literature on Bayesian non-parametric models in marketing (Ansari and Mela, 2003; Kim et al., 2004; Braun and Schweidel, 2011; Bruce, 2019), and particularly, on models to capture multiple sources of heterogeneity. For example, Dew et al. (2020) uses hierarchical Gaussian processes to capture dynamic heterogeneity; and Boughanmi et al. (2019) uses a hierarchical Dirichlet processes to uncover themes of musical albums that are predictive of success. In this paper, we introducing to the marketing literature the Pitman-Yor process for inferring heterogeneous discrete distributions with unknown number of components (e.g., number of contexts). The Pitman-Yor process (Pitman and Yor, 1997) generalizes Dirichlet processes by introducing an additional parameter that allows for more flexible patterns of the drawn discrete distributions.

More generally, this paper relates to previous work that has incorporated other sources of information when data on the main behavior of interest is thin — for example, by leveraging preferences from other product categories (Iyengar et al., 2003), by semantically linking web pages content and clicking to text-based search queries in search engines (Liu and Toubia, 2018), or by leveraging detailed acquisition data to infer future purchase behavior (Padilla and Ascarza, 2019). The present research contributes to this stream of literature by highlighting the value of extracting information from current and previous customer’s journeys to infer current customer preferences.

2.3 Empirical setting

We demonstrate the use of the customer journey as a source of information in the context of airline ticket purchases. We use data on flight search and purchases, from one of the largest worldwide online travel agencies. The dataset contains each query search, click and purchased flight for a sample of 5,000 active customers that searched for flight tickets between May 2017 and November 2017.² For each web page shown to those customers, we observe the customer id, the timestamp of when the customer accessed the page, the parameters of the search query associated with that page, and the list of results, including the flight attributes (price, length, airline carrier, etc) observed by the customer after entering the query. We observe a total of 5,285,770 flight offers displayed in 133,012 results pages, which resulted in 4,053 flight itineraries purchased.

We start by describing how the website works, and how we construct the “customer journey” in this context, and in particular, the query variables, the click occasions, and the purchase occasions. We discuss the trade-off in our data between thin historical purchase data, but rich with journey search observations.

2.3.1 The customer purchase journey of airline tickets

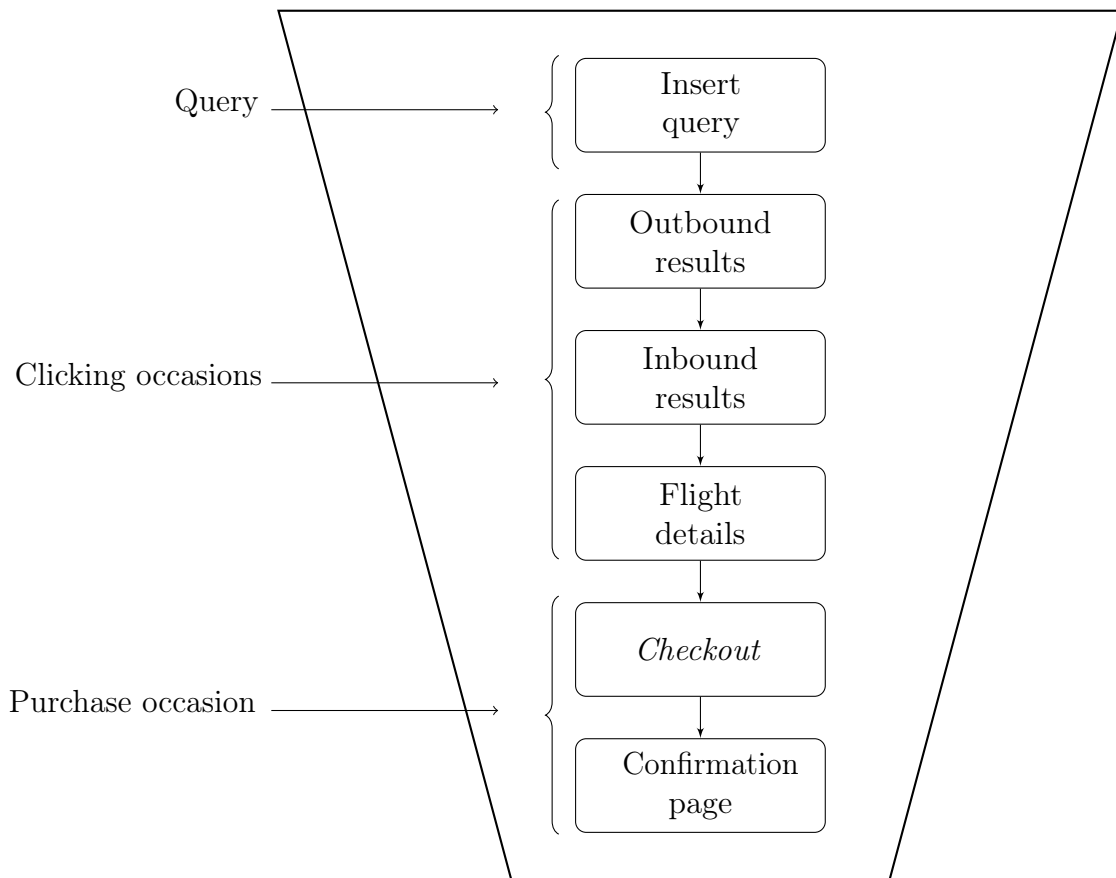
We describe the flow of a typical purchase journey on the website. The purchase journey starts when the customer lands at the homepage of the website to search for a flight.³ There

²We remove one customer from the analysis for which prices were unobserved for some of this customer’s journeys.

³Our data does not contain searches on packages (e.g., flight + hotel), and therefore we focus on purchase journeys over flights only. The model could be easily adapted for trip packages.

are two types of trips that the customer can choose from: (1) Roundtrip, and (2) One-way.⁴ For roundtrips the customer includes an origin and a destination, a departure date for the portion of the trip from the origin to the destination, known as the *outbound leg*, and a returning date for the portion of the trip from the destination back to the origin, known as the *inbound leg*. Each leg of the trip is composed by either one non-stop flight or multiple connecting flights. One-ways itineraries have only one direction of travel. Next we describe the flow of the roundtrip purchase journey as one-way is a nested version of the roundtrip purchase journey.

Figure 2.1: Flow of the customer purchase journey for roundtrip flights



⁴We drop from our analysis the the third type of trip multi-cities trips, as they constitute a very small portion of the trips.

Figure 2.1 shows the flow of steps of the roundtrip purchase journey for a roundtrip. The customer journey has 6 major steps: (1) Insert query, (2) Outbound results, (3) Inbound results, (4) Flight details, (5) Checkout pages (unobserved in our dataset), and (6) Confirmation page. At homepage, the customer starts specifying the search query (see Figure 2.2a) by selecting the type of trip to search for (e.g., *roundtrip*), and filling multiple fields (all of them required): origin and destination cities/airports, outbound and inbound departing dates (i.e., “departing” and “returning dates” in Figure 2.2a, respectively), and number of travelers. The customer then clicks on the “Search” button, which triggers the website to search the flight results that match the information from the query. After the search query is performed, the website displays the set of results for the outbound itineraries (see Figure 2.2b). Each of these itineraries are fully described by a path of flights that start at the origin airport and finish at the destination airport. The website clearly displays all relevant information of the outbound legs of the product search results, including price,⁵ the total duration of leg, the marketing airline carrier,⁶ the number of stops, departing and arrival times. Except for price, note that at this stage the website does not display information from the inbound leg.

If the customer clicks on the “Select” button of one of the outbound offers, the website displays the set of corresponding inbound results for the clicked outbound leg (see Figure 2.2c),⁷ For those resulting inbound offers, the website displays the same level of information displayed for the outbound offers (see Figure 2.2c), including the extra price of

⁵The price display corresponds to the price of the complete roundtrip itinerary, including the price of the outbound leg and the cheapest inbound leg that corresponds to the outbound leg

⁶The operating airline could differ.

⁷When the website queries the airlines servers, they may return offers for the whole outbound-inbound itinerary or generate a combination of multiple airlines separately for outbound and inbound legs to construct outbound-inbound itineraries, in order to find better prices.

each alternative compared to the minimum price (i.e., the price displayed in the outbound page of results). Once the customer clicks on the “Select” button of one of the inbound offers, the website shows a page with the details of all the information mentioned before from both the outbound and the inbound legs (see Figure 2.2d), as well as the full breakdown of the price (taxes and fees clearly displayed).

After the customer clicks on “Continue Booking”, the customer fills information about the passengers and proceeds with the payment steps.⁸ Finally, after finalizing the purchase the customer is shown a confirmation page. The one-way purchase journey is very similar, with the exception that instead of clicking through two set of results (outbound and inbound), the customer is displayed only one page of results, “One-way results”.

⁸While we do not observe the customer’s activity in the checkout page, we can track whether s/he clicked to the next page. In other words, if the customer follows all of these steps, the customer is shown a confirmation step, which we do observe. This means, if the customer bought the tickets, we do observe that outcome; but if s/he did not, we would not know at which specific page the customer decided not to purchase the flight.

Figure 2.2: Mock-up of purchase journey steps

(a) Examples of query page

(b) Example of outbound page results

Select your departure to Los Angeles Mon, Nov 18

Prices are roundtrip per person, include all taxes and fees, but do not include baggage fees

6:00pm - 9:33pm American Airlines JFK - LAX	6h 33m (Nonstop)	\$394 roundtrip	Select
7:00am - 10:07am United EWR - LAX	6h 7m (Nonstop)	\$397 roundtrip	Select
7:30pm - 10:40am Alaska Airlines JFK - LAX	6h 10m (Nonstop)	\$397 roundtrip	Select
9:20pm - 12:35am+1 Delta JFK - LAX	6h 15m (Nonstop)	\$397 roundtrip	Select

(c) Example of inbound page results

Your selected departure Mon, Nov 18 | [Change](#)

7:00am - 10:07am 6h 7m (Nonstop) from \$397
United EWR - LAX roundtrip

Select your return to New York Sat, Nov 23

Prices are roundtrip per person, include all taxes and fees, but do not include baggage fees

11:30pm - 1:30pm+1 American Airlines LAX - 3h 58m in BOS - JFK	11h 0m (1 stop)	+ \$0 roundtrip	Select
11:15pm - 7:55am+1 Alaska Airlines LAX - JFK	5h 40m (Nonstop)	+ \$44 roundtrip	Select
11:15pm - 9:05am+1 United LAX - 57m in ORD - EWR	6h 50m (1 stop)	+ \$65 roundtrip	Select

(d) Example of flight details results

Review your trip

<p>Mon, Nov 18</p> <p>From Liberty Int. (EWR) To Los Angeles Intl. (LAX)</p> <p>United 7:00am 10:07am 6h 7m, Nonstop EWR LAX</p> <p><small>Fare rules and Restrictions:</small></p> <ul style="list-style-type: none"> Pay to choose your seat Carry-on bag not allowed 	<p>Trip Summary</p> <p>Return: Arrives on 11/24/2019</p> <p>Traveler 1: Adult \$431.29 Booking Fee \$0.00</p> <p>Trip total: \$431.29</p> <p><small>Rates are quoted in US dollars</small></p> <p>Continue Booking</p>
<p>Sat, Nov 23</p> <p>From Los Angeles Intl. (LAX) To Liberty Int. (EWR)</p> <p>United 11:15pm 9:05am 6h 50m, 1 stop LAX EWR</p> <p><small>Fare rules and Restrictions:</small></p> <ul style="list-style-type: none"> Pay to choose your seat Carry-on bag not allowed <p><small>Arrives Sat, Nov 24</small></p>	

[< Change flights](#)

Note that, while Figure 2.1 shows a linear purchase funnel, in reality the journey can be highly non-linear. That is, the customer may go back from each step to enter a new/revised query, to click on alternative outbound or inbound results, etc. Moreover, this process does not need to occur during the same internet session, but can occur over the

course of multiple days (Lee et al., 2018). We create a flexible definition of the customer journey by combining pages/sessions that belong to the same trip need. This process is not straightforward because customers sometimes modify the search query aiming to obtain a new set of products that potentially would satisfy the same need. To allow for a flexible definition of a journey we combine into a journey session with similar queries that: (1) occur at different points in time sometime over days and weeks; (2) have departing or arrival dates within up to 4 days; and (3) have origin or destination to close-by airports and cities within a 140 miles range (approx. 225 kms.). Accordingly, we use the queries in our data to construct the journeys. To construct the journey, we combine the different pages each customer saw within the same journey, sort them by timestamp, and remove all pages within a journey after a purchase is made, to remove the infrequent behavior of customers checking prices of the same itinerary *after* purchase. This process resulted in a total of 28,025 journeys, corresponding to an average of 5.6 journeys per customer.

We believe that our conceptualization of journeys (instead of simply using individual search queries) better captures the nature of the purchase journey; because searches included in the rules described above are, most likely, aimed to satisfy the same need for a trip. Moreover, this broader definition of journey allows us to seamlessly integrate behavior across sessions that are aimed at covering the same need.

2.3.2 Extracting information from the data

2.3.2.1 Search queries

Using the query information, we construct several variables that aim to capture in more details the context of this trip. While some pieces of information are directly provided by the customer (e.g., destination), others could be indirectly determined. For example, whether the trip includes weekends can be extracted from the dates, or the trip distance can be inferred from the origin and destination airports. We combine these variables into a set of “query variables” that aim to capture information about the journey in four different dimensions: (1) who is traveling, (2) which market this flight belongs to (origin-destination), (3) when is the trip, (4) when was the search made. These variables will help inform preferences capturing journey-specific needs, even for different journeys of the same customer. Table 2.1 shows these variables and their corresponding summary statistics.

Overall, we observe a great variety of trip characteristics. Among all journeys customers undertake in our sample, 66% of them are roundtrip (vs. one-ways); in 28% of cases, customers are searching for more than one adult, whereas in 8% they search for trips with kids. With respect to the dates of the trip, the average stay for roundtrips is 11.80 days, 37% of journeys are searching for flights during the summer season,¹¹ 3% for the holiday season¹², and 66% of flight searches include stays over weekends. In terms of origin-destination of the trip, the average trip distance is 3,548 kilometers or 2,205 miles (e.g., approx. New York to Las Vegas); 59% of journeys are domestic (including US-Canada,

¹¹We define the summer season from June 30th, to September 4th.

¹²We mark as holiday season stays that include either Thanksgiving, Christmas or New Year’s holidays.

Table 2.1: Summary statistics of query variables

Query variable	Mean	SD	Quantiles		
			5%	50%	95%
Continuous					
Trip distance (kms)	3,584.16	3,465.07	448	2,269	11,529
Time in advance to buy (days)	50.73	59.82	1	29	182
Length of stay (only RT) (days)	11.80	21.25	2	6	37
Binary					
Is it roundtrip?	0.66	.	0	1	1
Traveling with kids?	0.08	.	0	0	1
More than one adult?	0.28	.	0	0	1
Is it domestic? ⁹	0.59	.	0	1	1
Is it summer season?	0.37	.	0	0	1
Holiday season?	0.03	.	0	0	0
Does stay include a weekend?	0.66	.	0	1	1
Flying from international airport?	0.74	.	0	1	1
Searching on weekend?	0.21	.	0	0	1
Searching during work hours?	0.49	.	0	0	1
Categorical					
Market					
US Domestic	0.51	.	0	1	1
US Overseas	0.18	.	0	0	1
Within North America ¹⁰	0.15	.	0	0	1
Non-US within continent	0.10	.	0	0	1
Non-US across continent	0.06	.	0	0	1
Type of departure location					
Airport	0.88	.	0	1	1
Multi-airport City	0.08	.	0	0	1
Both	0.04	.	0	0	0

within-EU, or within-country flights); 51% are US Domestic, 18% are for US-Overseas trips, 15% are between US and Mexico or Canada and Mexico, 10% of the searches are for continental trips that do not include the US, and the remaining 6% are for trips across continents that do not include the US. Finally, with respect to the time between search and flight, purchase journeys occur, on average, 50.73 days prior to the departing date; 88% introduce a departing location code for an airport (e.g., JFK), 10% a departing code of a city (e.g., NYC), and the rest include a departing code that refers to both city and airport (e.g., MIA); 21% of the times customers search during the weekend, and 49% during work hours (defined by local time of departing city).

2.3.2.2 Click occasions

Once the query is clearly defined, we need to build a set of “click occasions” faced by the customer. These click occasions are composed by a set of alternatives to click on, and the outcome of what was actually clicked (or not). There are two types of click occasions: (1) clicking occasion on a outbound results page (where clicking in a product leads to an inbound results page), and (2) clicking occasion on an inbound results page (where clicking in a product leads to flight detail page). We observe and allow in our model the customer to click on multiple flights from each click opportunity.

For outbound and inbound results pages, by default, results are sorted increasingly in price.¹³ Some results may be further filtered using the attribute filters such as by airline or number of stops or sort the product results using a different sorting mechanism — while these actions would be valuable pieces of information for inferring preferences, we do not observe explicitly when these actions are taken because the firm does not record these actions. Accordingly, we treat filtering results as if customers were searching again.¹⁴

Putting all clicks made by the customer, we have a list of pages visited by the customer: (1) pages in which no product was clicked (e.g., an outbound result page in which no offer was good enough for the customer, and the customer decided to do another search); (2) pages in which exactly one offer was clicked; and (3) pages in which the customer clicked multiple times, either in different offers or in the same offer.¹⁵ For the pages in which no

¹³Analogously, for one way results.

¹⁴While, in theory, we could try to infer the filter from the set of results, product attributes are highly correlated (e.g., all direct flights being of the same airline), making it difficult to pin down which specific filter was set.

¹⁵When a customer clicks on an inbound offer, it opens a new tab, and therefore a customer can click on several inbound offers from the same outbound page, and open several flight details pages.

product was clicked, we create a click occasion with outcome of: *keep searching* if it is not the last observed page in the journey, or *leave* if it is the last observed page in the journey (and no purchase was made).¹⁶ For the pages in which offer/s were clicked, we create the click(s) occasions directly. We label all offers that have been clicked before for all subsequent click occasions from the same page to account for the fact that a customer has already clicked on a product, and may be less (or more) likely to click on the same product again.^{17,18} This process resulted in 132,665 click occasions, averaging to 4.7 click occasions per journey.

2.3.2.3 Purchase occasion

For each journey, we observe whether it included a purchase or not. We describe the last steps in the purchase funnel, and how we construct the purchase occasion in our customer journey.

We create one purchase occasion per journey by creating a set of products that were likely to be considered for purchase. There are multiple approaches to construct the consideration set (i.e., alternatives that the customer considers before purchase). One could, as it is commonly done in the literature, define considered products as those product the customer clicked on to observe details. This approach seems most appropriate when product attributes are only revealed once the customer clicks on the details page (e.g., Bronnenberg et al., 2016), and therefore customers open those pages to observe these unknown attributes (e.g., photos of the products or reviews). However, in our context, the customer observes

¹⁶To avoid labeling censored journeys as a no-purchase journey, we remove the last observed click occasion in which purchase was not made but the departure date of the flight is after our last day of our observation period.

¹⁷A similar process is used for one-way offers.

¹⁸We could also treat the confirmation of purchase page, as a click occasion from the flight details page. We decided not to this, as the customer has already finished the journey at that point and hence this last step does not provide additional information.

most relevant attributes when inspecting the list of (inbound) results. Thus, there is little incentive for a customer to click on the flight details page for information gathering purposes. For example, customers on average see 0.5 flight details pages for roundtrips and 0.7 flight details pages for one-ways. Therefore, considering only flight that were clicked on as being part of the consideration set would eliminate many flight that were considered without clicking.

At the other extreme, one could include in the consideration set *all* products displayed over the course of a journey. In theory, it is plausible that a customer is considering all the products that s/he saw prior to purchase as the full information for all flights is revealed in the results page. For example, a customer can click on an outbound flight to see a list of corresponding inbound flights, and without clicking on any of those, the customer has almost full information about the features of those inbound flights. However, given that most customers are exposed to hundreds of products per journey (189 flights in average), this approach would be not only unrealistic—it is unlikely that customers consider all the products displayed but rather a subset of these—but also impractical from a computational perspective. Including all these alternatives will increase significantly the computational burden for estimation.

Therefore, we take an intermediate solution by constructing the set of considered products as a combination of the products that were clicked on and therefore observed in the Flight details page, plus the top 20 results from outbound and inbound results page. In other words, we assume a heuristic rule to determine consideration set formation and we model purchases given that consideration set. Finally, we register the outcome of the

purchase occasion as a purchase for the product that was chosen, if any, or we register it as a non-purchase in case no product was purchased.

2.3.2.4 Product attributes

Customers observe multiple product attributes that they consider when making a click and purchase decision for an offer. For a roundtrip journey, all attributes, except price, are specific to *each leg* of the trip. That is, there is a set of attributes that describe the outbound leg of the trip, and there is the same set of attributes that describe the inbound (returning) leg of the trip.¹⁹ We do highlight that there is an important difference between leg-specific outbound and inbound product attributes. Outbound offers are shown as the first step in the journey, and therefore are a more representative sample of offers available in the market. Inbound offers, on the other hand, are shown only after the customer clicks on the specific outbound offer. Therefore, inbound attributes can have a different distribution than their corresponding outbound attributes, as their appearance in the data depends on the customer clicking on the corresponding outbound flight.

We summarize these attributes in Table 2.2.²⁰ Prices are measured at the whole trip level. The average offer displayed is priced at \$1,547; but offers vary significantly in their price, with a standard deviation of \$3,249. Not only the offers within a journey vary in their prices, but also journeys have a different price level that strongly depends on origin-destination and the dates. This variation in price becomes clearer when analyzing the price of the cheapest offer per journey. The cheapest price displayed per journey has an average of \$698 across all purchase journeys, with a standard deviation of \$1,526. This

¹⁹For one-way journeys, clearly only one set of these attributes is observed.

²⁰One-way offers are summarized within the outbound component of the table.

Table 2.2: Summary statistics of product attributes in page results

Product attribute	Mean	SD	Quantiles		
			5%	50%	95%
Product level attributes					
Price	1,547	3,269	196	751	5,320
Cheapest price per journey	698	1,526	98	401	2,117
Outbound level attributes					
Length of trip (hours)	11.28	8.49	2.05	8.42	28.60
Shortest length of trip per journey (hours)	5.86	5.05	1.25	4.07	17.08
Number of stops: Non stop	0.20	.	0	0	1
Number of stops: One stop	0.59	.	0	1	1
Number of stops: 2+ stops	0.21	.	0	0	1
Alliance: Alaska Airlines	0.04	.	0	0	0
Alliance: Frontier	0.01	.	0	0	0
Alliance: JetBlue	0.03	.	0	0	0
Alliance: Multiple alliances	0.07	.	0	0	1
Alliance: Other – No alliance	0.07	.	0	0	1
Alliance: OneWorld (American)	0.27	.	0	0	1
Alliance: Skyteam (Delta)	0.27	.	0	0	1
Alliance: Spirit	0.02	.	0	0	0
Alliance: Star Alliance (United)	0.23	.	0	0	1
Dep. time: Early morning (0:00am - 4:59am)	0.04	.	0	0	0
Dep. time: Morning (5:00am – 11:59am)	0.47	.	0	0	1
Dep. time: Afternoon (12:00pm - 5:59pm)	0.31	.	0	0	1
Dep. time: Evening (6:00pm - 11:59pm)	0.18	.	0	0	1
Arr. time: Early morning (0:00am - 4:59am)	0.05	.	0	0	0
Arr. time: Morning (5:00am – 11:59am)	0.24	.	0	0	1
Arr. time: Afternoon (12:00pm - 5:59pm)	0.34	.	0	0	1
Arr. time: Evening (6:00pm - 11:59pm)	0.37	.	0	0	1
Inbound level attributes					
Length of trip (hours)	11.08	9.02	1.83	7.92	29.50
Shortest length of trip per journey (hours)	6.17	5.31	1.25	4.27	17.75
Number of stops: Non stop	0.19	.	0	0	1
Number of stops: One stop	0.70	.	0	1	1
Number of stops: 2+ stops	0.11	.	0	0	1
Alliance: Alaska Airlines	0.02	.	0	0	0
Alliance: Frontier	0.02	.	0	0	0
Alliance: JetBlue	0.02	.	0	0	0
Alliance: Multiple alliances	0.02	.	0	0	0
Alliance: Other – No alliance	0.07	.	0	0	1
Alliance: OneWorld (American)	0.51	.	0	1	1
Alliance: Skyteam (Delta)	0.13	.	0	0	1
Alliance: Spirit	0.05	.	0	0	1
Alliance: Star Alliance (United)	0.15	.	0	0	1
Dep. time: Early morning (0:00am - 4:59am)	0.03	.	0	0	0
Dep. time: Morning (5:00am – 11:59am)	0.65	.	0	1	1
Dep. time: Afternoon (12:00pm - 5:59pm)	0.18	.	0	0	1
Dep. time: Evening (6:00pm - 11:59pm)	0.14	.	0	0	1
Arr. time: Early morning (0:00am - 4:59am)	0.04	.	0	0	0
Arr. time: Morning (5:00am – 11:59am)	0.55	.	0	1	1
Arr. time: Afternoon (12:00pm - 5:59pm)	0.19	.	0	0	1
Arr. time: Evening (6:00pm - 11:59pm)	0.23	.	0	0	1

indicates that raw prices may not be a good proxy to capture price sensitivity among our customers, as prices are only compared within a journey. For example, a New York - Chicago roundtrip ticket for \$600 may be considered expensive for this trip, whereas a roundtrip flight from New York to Buenos Aires for \$800 may be considered a good deal.

Offers also differ in terms of how long each leg of the trip is. The average outbound leg of a displayed trip takes 11.28 hours, with a large variation within and across journeys. The shortest flight per journey takes, on average, 5.86 hours for the outbound leg. Most displayed flights are one stop flights (59% for outbound legs, and 70% for inbound legs), whereas nonstop flights account for 20% of outbound offers and 19% of inbound offers.

Airline data is fairly sparse and therefore we aggregate airlines into alliances. Alliances are group of airlines that share benefits and usually run in “shared codes” (e.g., a flight from JFK to Madrid that is operated by Iberia might be sold by American Airlines, British Airways, FinnAir and Iberia, all belonging to the same alliance). The three biggest alliances are: OneWorld (including American Airlines), Skyteam (including Delta Airlines), and Star Alliance (including United Airlines). We kept some individual airlines that are not part of any alliance but represent significant proportion of the displayed offers. Particularly, for the US domestic market we keep Alaska Airlines, Frontier, JetBlue, and Spirit. We group all other smaller airlines that are not part of any alliance in “Other - No alliance” category.²¹ Finally, we label as “Multiple alliances” offers that have connecting flights of different alliances in the same leg of the trip.²² We find that the big three alliances account for 77% of all outbound offers, and 79% of inbound offers.

Finally, offers also vary in terms of their departing and arrival times. For most outbound and inbound legs, the first flight departs in the morning. However, the last connecting flight of the outbound leg tends to arrive either in the afternoon or the evening, whereas that for inbound legs, they tend to arrive in the morning.

²¹There are other smaller regional alliances, but they do not represent a significant portion of the offers in our dataset.

²²This should not be confused with offers that have one alliance for outbound and another for inbound, which have the corresponding alliance for each leg, outbound and inbound.

2.3.3 Inferring preferences from purchase journey data

Table 2.3 shows the the total number of customers, journeys, purchases, click steps and clicked products. We observe a total of 28,025 journeys, for which we aim to estimate individual-level preferences. The data indeed exhibit lack of past purchase history at the individual level— while, on average, each customer undertakes 5.606 purchase journeys, the average number of purchases per customers is 0.81. This piece of evidence highlights the lack of purchase history that challenges preference estimation using traditional models that rely on long individual purchase history. Arguably, the lack of past purchases could be caused by the observation window not be long enough (in our case, 7 months) and therefore a longer time horizon would solve the problem.²³ However, privacy concerns and average lifetime of cookies tracking customer behavior are often not much longer than our time horizon. For many high involvement products, customers do not fly as often as they buy certain consumer package goods such as ketchup, which results in this persistent lack of purchase history. On the other hand, customers click on products along the journey (on average, 6.46 clicked products per customer, and 1.15 clicked products per journey). This piece of information is relevant as we can learn preferences and contexts from customers clicking on products even when those actions may not end up in a purchase.

On average, 14.5% of journeys end with a purchase. This number may seem high for an online retailer but there are two caveats to this quantity. First, our data correspond to a sample of active customers and therefore this figure would be lower for the average customer of the firm. Second, in this paper we adopt a broader definition of a journey, which includes

²³Also note that our data include a relatively active set of customers, implying that the lack of past purchases is even more severe for total user base in this firm.

Table 2.3: Data summaries, per customer and per journey.

Variable	Total	Average per...							
		Customer		Journey		Purchased journey		Non-purchased journey	
		Mean	s.e.	Mean	s.e.	Mean	s.e.	Mean	s.e.
Customers	4,999
Journeys	28,025	5.606	0.072
Roundtrip	18,469	3.695	0.054
One-way	9,556	1.911	0.065
Purchases	4,053	0.811	0.015	0.145	0.002	1.000	.	0.000	.
Click occasions	132,665	26.538	0.351	4.734	0.039	8.632	0.128	4.075	0.038
... in OW search	44,015	8.805	0.246	4.606	0.064	6.094	0.150	4.166	0.070
... in RT outbound	56,355	11.273	0.194	3.051	0.033	5.180	0.132	2.811	0.033
... in RT inbound	16,054	3.211	0.045	0.869	0.012	2.556	0.051	0.679	0.011
Clicked products	32,295	6.460	0.069	1.152	0.013	2.960	0.046	0.847	0.012
... in OW search	6,548	1.310	0.033	0.685	0.012	1.687	0.030	0.389	0.011
... in RT outbound	16,054	3.211	0.045	0.869	0.012	2.556	0.051	0.679	0.011
... in RT inbound	9,693	1.939	0.028	0.525	0.008	1.887	0.035	0.371	0.007

multiple searches for the same customer need, whereas a traditional conversion rate would treat different search queries, with different variations of airports or dates as different and independent purchase funnels.²⁴

In each journey, customers click on 1.15 products on average, with this distribution varying by type of trip. For example, in one-way journeys, customers click on an average of 0.69 one-way itineraries per journey to observe the details page. In roundtrip journeys, customers click, on average, on 0.87 outbound itineraries per journey to observe their corresponding inbound flights; and they click on 0.53 inbound results per journey to observe the full flight details page (1.39 clicked products in total). These figures support our choice for a more flexible definition of considered products, which includes product viewed in the outbound and inbound pages as opposed to only clicked products. If one were to treat the consideration set as the products that the customer saw in the details page, as it is

²⁴We provide an example of why this is the case. If a customer searches for three different sets of (very close) dates but s/he only purchases in the last search query, a traditional conversion metric that treats all these searches as independent would summarize this information as 2 non-purchase sessions and 1 purchase session. Using our broader definition of a customer purchase journey, we would measure 1 single journey with a purchase.

commonly done in the consumer search literature, these findings would imply that customers consider between 0.5 and 0.7 flights per journey, which seems unrealistic.

As expected, there are considerable differences between journeys that end with a purchase and those that do not. Customers in journeys that end up in a purchase are exposed to more than twice the amount of occasions for clicking (8.63 vs. 4.07), and they click on more than three times more products than non-purchase journeys (2.96 vs. 0.85). Interestingly, journeys that do not end in a purchase still contain clicked products; we argue that these clicks should inform preferences as well, as these are choices that customers make about some products but not others.

In summary, while we have very limited information on past purchases, the data at the journey level is quite rich; in terms of queries, click occasions, and click behaviors. Our goal is to integrate those behaviors to infer individual preferences for predicting purchases. We move now to describe our model which integrates these sources of data.

2.4 Model

The main goal of the model is to be able to estimate preferences in contexts in which individual-level purchase history is likely to be very sparse, and heterogeneity exists both across consumers and within a customer across journeys. For the sake of clarity and generalizability, we describe our model in the context of our empirical application — customers searching for flight tickets. However, we want to highlight that the proposed model is applicable to other contexts as well. While some components of the model may be more or less relevant for different applications, each component can be easily

adjusted to overcome the challenges specific to other settings. For example, the usage of clicks to inform preferences is widely applicable to most online contexts as firms perfectly observe customer clicking behavior. The lack of individuals' purchase history is most relevant for other high involvement purchases such as cars or durable goods, while journeys with varying contexts and needs would be most relevant for experiential purchases such as hotel stays, restaurants reservations, food delivery or media consumption.

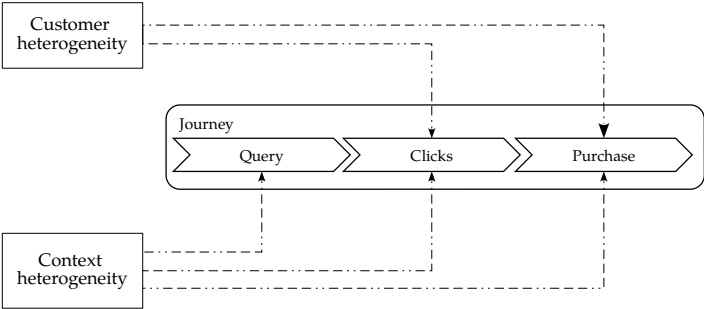
In turn, online travel search is among the most complex settings as it contains all these aforementioned challenges one needs to address. Moreover, in addition to the lack of purchase history and the changing contexts, the set of available products is extremely large, and essentially, unique to each customer. The set of available products depends not only on the origin-destination of the trip and dates, but also on the moment the product is searched, as availability and particularly prices change dramatically over time. This feature of the product space rules out classical approaches to recommendation systems, such as matrix factorization as there is extremely low chance that two customers are exposed to (even less so, buy) the same offer at the same price.

2.4.1 Model intuition

Before describing the model components in detail, we provide some intuition behind the main modelling assumptions. We conceptualize the customer purchase journey as a series of steps that start with a search query, are followed by a series of clicks through different stages of the purchase funnel, and may finalize, eventually, with a purchase (Figure 2.3). These behaviors are realizations of the customer overall preferences and the specific needs that s/he

aims to satisfy. For example, when a customer is searching for a flight, s/he has a trip in mind, and therefore a specific need that these tickets will satisfy. The customer may be looking for flights for a honeymoon, for a summer family vacation, or for a business trip to a nearby city. When inserting the query, the customer would ask for a trip that best matches that kind of trip (e.g., a honeymoon will likely be a trip for two adults, longer than 4 or 5 days, with an exotic destination). Then, in choosing a flight — clicking on it or eventually buying it — the customer will have some stable (journey independent) preferences such as his/her preferences over an airline because s/he is a frequent flyer from that airline. However, different types of trips could affect the customer’s preferences over flights. For example, when the customer is looking for a business trip, s/he may be less price sensitive, or when looking for a summer family vacation s/he may have stronger preferences towards avoiding connections if s/he is flying with kids.

Figure 2.3: The data generating process



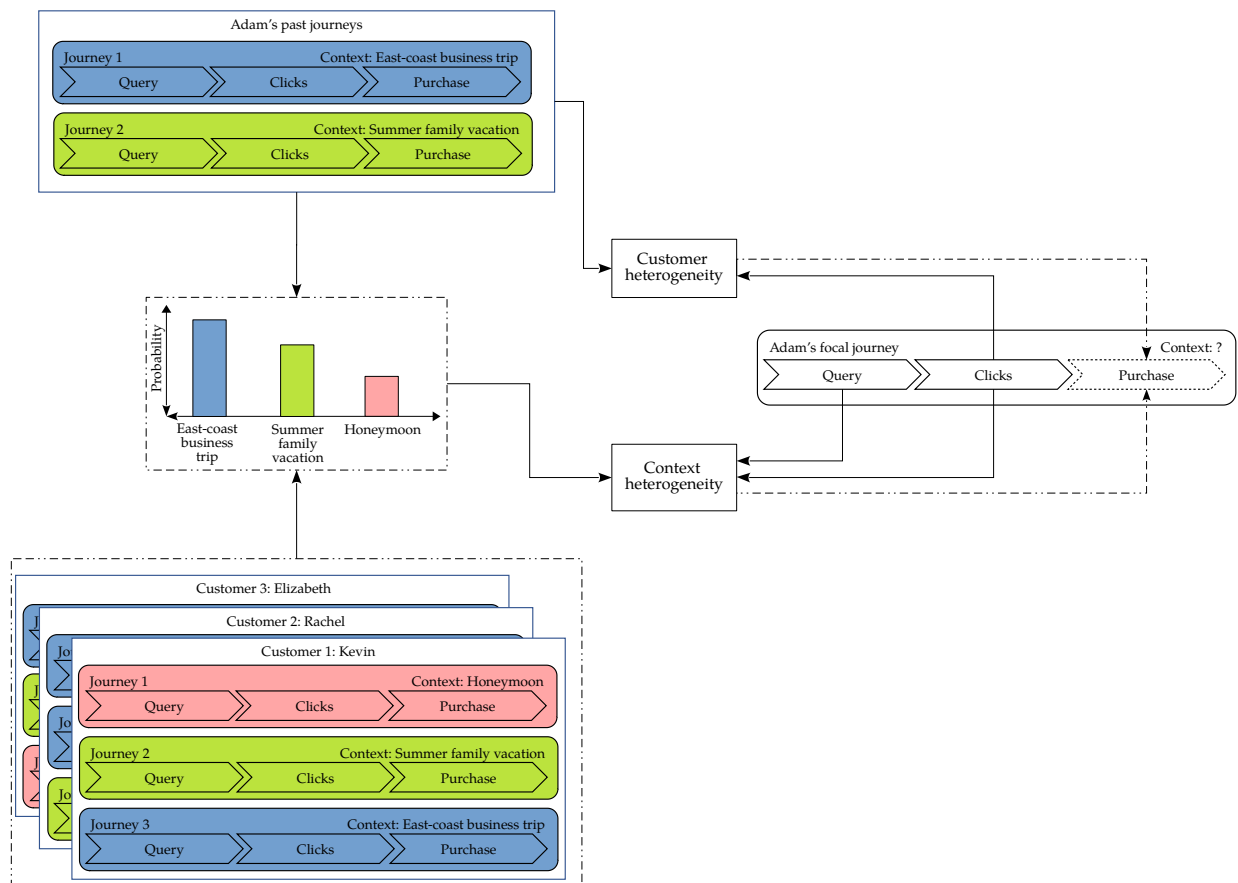
While the trip type or context is unobserved, the firm can infer it from the data, along with the customer preferences that are consistent across trips. It does so by combining information across similar purchase journeys (from other customers), clustering them together in what we define as the *context of a journey*. At an abstract level, the context of a

journey represents the unique needs that the customer seeks to satisfy by purchasing a product in a journey of this type, that are different from the preferences captured by the individual's other journeys.

To better understand how the model combines historical data from the focal customer, with information from other customers with possibly context, Figure 2.4 visually describes how the model learns from the different sources of data. Let us assume that Adam is currently going through a flight search (focal journey in the figure), and the firm wants to predict his actions during the journey, whether he will buy at the end of this journey, and if so, what product he will buy. Adam's behavior in this journey will be determined by both his individual-level stable preferences and the specific needs that he aims to satisfy in this particular trip (the context). The firm has seen Adam in the past (he has two previous journeys) and has also observed Kevin, Rachel, and Elizabeth, going through three purchase journeys each. Taking together Adam's, Kevin's, Rachel's, and Elizabeth's past journeys the firm inferred that the population of travelers and journeys belong to one out of many contexts—in this example we use three contexts, which we identify with three different colors: east-coast business trip (in blue), summer family vacation (in green), and honeymoon (in pink). Moreover, customers are different in nature; i.e., each customer has individual-level preferences that s/he carries for all the journeys that s/he undertakes.

Each purchase journey is composed by: (1) search query, (2) clicks in several steps, and (3) a purchase decision, which includes the no-purchase alternative. Our model treats each of these components as an outcome, which depend on the context and the customer's stable preferences. While the queries are determined only by the context—e.g., Kevin's third and Rachel's first journeys are both of context East-coast business trip (blue)—clicks

Figure 2.4: Model intuition



and purchases are determined jointly by the customer's stable heterogeneous preferences and by the context of the journey. Note that, like customer preferences, contexts are unobserved to the firm and therefore need to be inferred from the behavior. The model infers that Kevin's second journey, Rachel's third journey, and Adam's second journey are all of the same context (summer family vacation, in green), because they have similar queries, but also because they all exhibit similar deviations from the preferences of each customer (e.g., they were more interested in non-stop flights than they would on average). Customer stable preferences are inferred from the customer's consistent behavior across journeys across contexts.

The model would learn from all purchase journeys that have already ended and use these estimates to infer context and preferences for the focal customer and journey.²⁵ Adam started his focal purchase journey by entering a query, and he has possibly also clicked on some products along the way. At this point, the firm is interested in inferring Adam's purchase preferences for this journey in order to predict whether he will buy, and if so, what product he will buy. The context of the focal journey is updated based on the prior distribution of context informed by others customers and Adam's own past purchases *and* the query and clicks of the focal journey. Similarly, Adam's clicks in the focal journey help update Adam's stable heterogeneity preferences, which can be used to predict his choices. It is through these updates of the context and stable heterogeneity distributions from the query and clicks in the focal journey that our model leveraged the within journey information and combines it with historical information from the focal and other customers to predict the choices in the focal journey. For example, if Adam is flying from NYC to Boston, from Monday to Thursday, the model can infer that with higher likelihood this journey is an East-coast business trip. This information may help us determine that Adam's price sensitivity for this trip is lower than usual. In addition, Adam will likely click on flights that are more expensive than what he would usually click on or buy, which will reinforce that he has a lower price sensitivity for this particular trip, giving the model a stronger signal that this trip context is likely to be an East-coast business trip. Finally, combining these different sources of information, the model can now infer Adam's preferences for purchase in this

²⁵The model knows that a journey has ended when the journey ended up with a purchase, or when departure date lies within our observation window. Otherwise, the model accounts for censored data as some journeys may have not ended yet.

journey, and use them to recommend products that are most relevant to Adam at this particular point.

2.4.2 Model development

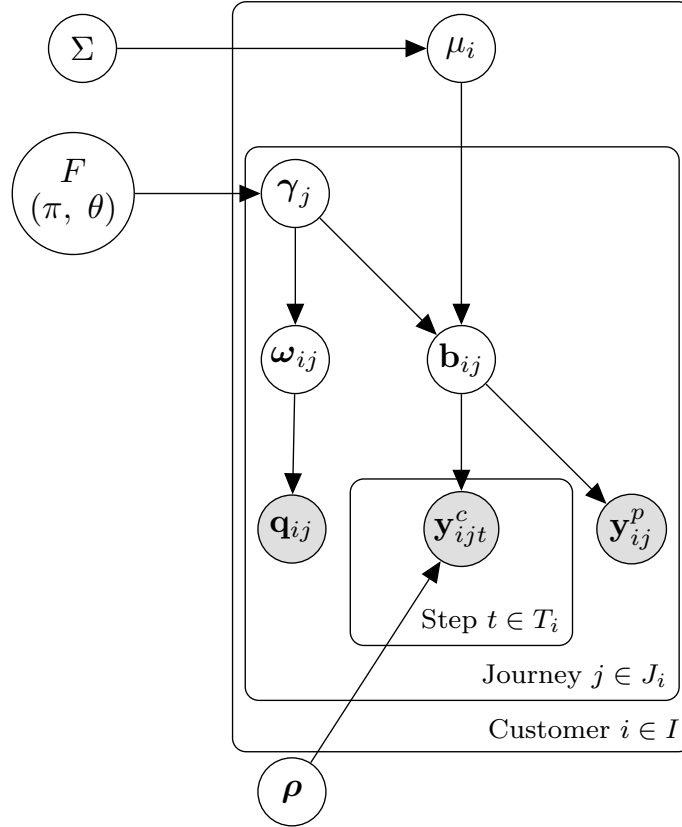
We now present the formal specification of the model. We start by stating the notation and providing an overview on the multiple components of the model and how we bring all information together. Then, we describe each component: query, clicks, and purchase, explain how the model combines these components, and particularly, how we model the journey’s context. We finalize with the details of the model specification tailored to our empirical application and a description of our estimation approach.

2.4.2.1 Model overview

We index customers by $i \in \{1, \dots, I\}$, their journeys by $j \in \{1, \dots, J_i\}$, where J_i is the number of purchase journeys customer i has undertaken, and by $t \in \{1, \dots, T_{ij}\}$ the steps of customer i in journey j . Our model links queries, clicks and purchases over the course of the purchase journey.

At a general level, our model has three major components that we model as outcomes: queries (\mathbf{q}_{ij}), clicks (\mathbf{y}_{ijt}^c), and purchase (\mathbf{y}_{ij}^p). We are modeling each of these components as outcomes from distributions that are parametrized by journey-specific parameters. First, we model the query component by $\mathbf{q}_{ij} \sim f_q(\boldsymbol{\omega}_{ij})$, where $\boldsymbol{\omega}_{ij}$ is a journey-specific vector of parameters for some multivariate distribution f_q , that has both discrete and continuous components. Second, we are modeling each click decision using a choice model, f_c , with

Figure 2.5: DAG of customer journey model



vector of preferences \mathbf{b}_{ij} and each journey's purchase decision using another choice model, f_c , as a function of the same vector of preferences \mathbf{b}_{ij} .

That is, the model combines clicks and purchases from the same journey by specifying both components as a function of preferences \mathbf{b}_{ij} . However, in reality customers may exhibit different preferences when exploring the options versus when they choose a flight to buy. For example, a customer may be more likely to click on expensive outbound options to explore the corresponding inbound offers, but when choosing the actual flight to buy, choose cheaper flights. To account for such a behavior, we introduce the vector $\boldsymbol{\rho}$, which only affects the clicking decision and captures systematic differences in how customers compensate attributes.

We leverage information from the customer’s past journeys, by specifying journey-specific preferences \mathbf{b}_{ij} to be a function of a vector of stable preferences $\boldsymbol{\mu}_i$, which are drawn from the population distribution parametrized by Σ . We further link queries with clicks and purchases from the same journey by assuming that \mathbf{b}_{ij} and $\boldsymbol{\omega}_{ij}$ are a function of a vector $\boldsymbol{\gamma}_j$, that reflects how a specific trip context affects the query, clicks and purchase decisions.

Finally, we leverage the data from other customers with similar journeys by assuming that the context-specific vector $\boldsymbol{\gamma}_j$ is drawn from a common distribution of contexts F shared by all journeys of all customers.²⁶ We can gain intuition by thinking of F as a histogram, or a distribution of segments of journeys, where each segment is described by how frequently it appears in the population (π) and the parameters that describe the context in terms of their meaning for the queries, clicks, and purchase components of the model (θ). We uncover the contexts non-parametrically from the data, as we describe in Section 2.4.2.5.

2.4.2.2 Query model

We index by $m \in \{1, \dots, M\}$ the different types of query variable, where each type variable m relates to one piece of information (e.g., length of the stay, traveling with kids). Because these pieces of information were provided by the customer to obtain a set of products results that match his/her preferences, we treat each query variable as an outcome that depends on some unobserved component that captures the customer’s need in that journey. Moreover, treating query variables as an outcome allows us to easily account for missing query

²⁶Technically, as journey j belongs to customer i , $\boldsymbol{\gamma}_j$ depends on both indexes i and j . However, we drop the explicit dependence on i to remark that the vector is capturing journey-level traits that are not informed by other journeys of customer i

variables, or query variables that are not valid for some journeys. For example, one-way journeys do not have a length of stay as they do not have a returning date.

Each type of query variable m could be of type: (1) binary, (2) categorical, (3) continuous real-valued, or (4) continuous positive-valued. We denote q_{ijm} the realization of query variable m , for customer i and journey j , which we model using a different distribution f_m for each type of variable m by $q_{ijm} \sim f_m(\omega_{mj})$, where ω_{mj} is a journey specific parameter for query variable m . If query variable m is binary (e.g., whether the customer is traveling with kids), we model it by

$$q_{ijm} \sim \text{Bernoulli}(\omega_{ijm}), \quad (2.1)$$

where $\omega_{ijm} \in (0, 1)$ is a scalar parameter. If query variable m is categorical with N possible values (e.g., which market does the trip belong to), we model it by

$$q_{ijm} \sim \text{Categorical}(\omega_{ijm}), \quad (2.2)$$

where ω_{ijm} is the vector of probabilities length N , such that $\omega_{ijmn} \geq 0$ and $\sum_n \omega_{ijmn} = 1$. If query variable m is continuous real-valued (e.g., the log of the distance of the trip), we then model

$$q_{ijm} \sim \mathcal{N}(\omega_{ijm}, \sigma_m^2), \quad (2.3)$$

where ω_{ijm} is a scalar representing the mean and σ_m^2 a positive variance shared across journeys.²⁷ Finally, if m is continuous positive-valued (e.g., length of stay), we then model it by

$$q_{ijm} \sim \exp(\omega_{ijm}), \quad (2.4)$$

where ω_{ijm} is a positive scalar.²⁸

Our model can easily accommodate other distributions if one would aim to capture specific features of the query variables. For example, Poisson or Binomial distributions for count variables, and Student's t-distribution or Cauchy distribution for long-tailed continuous variables. Our choice of distributions is based on the nature of the query variables from our empirical application and computational convenience for drawing efficiently the parameters from the posterior distribution.

We define the vector of query parameters as $\boldsymbol{\omega}_{ij} = \left[\omega_{ij1} \quad \dots \quad \omega_{ijM} \right]'$. We come back to these parameters in Section 2.4.2.5 when we relate the unobserved queries component with the click and purchase preferences in a particular journey.

²⁷We choose to define σ_m fixed across all journeys, to avoid singularity issue. Like Gaussian Mixture Models when variances are cluster specific, this model would behave similarly when we cluster non-parametrically journeys that are of similar characteristics. In these cases, fitting such a model could lead to one cluster fitting a single specific data point with mean equal to the data point value and variance converging to zero, leading to singularity in the Gaussian density.

²⁸We choose to model some query variables using an exponential distribution instead of a log-normal distributions as variables associated with time (e.g., length of stay, time in advance for booking the flights) tend to be distributed closer to an exponential distribution.

2.4.2.3 Click model

Along the journey, the customer clicks through pages of product results in a series of steps from the initial search query, eventually to purchasing or leaving. The customer can navigate back and forth between clicking on flight options and refining his/her searches. In each step of the process, the customer decides among: (1) clicking on one of the products to move further in the journey, (2) continuing to search to receive a new set of results, or (3) leaving and finalizing the journey without a purchase.

The model accounts for the different types of pages in which a customer can click on products. For example, in our empirical context, roundtrip flights have two types of pages where the customer can click on products: “Outbound results page”, where the customer chooses the outbound/departing flight; and “Inbound results page”, where the customer chooses the inbound/returning flight, whereas one-way journeys only show “One-way results page”. We denote $p(t) \in \{1, \dots, P\}$ the type of page of step t .²⁹ These types of pages differ in how products are shown as well as what happens next when the customer clicks on one of the shown products. Accordingly, we allow for different base click rates for outbound, inbound and one-way clicks.

We model the click decision at step t using a discrete choice model, and we index choice alternatives by k . We define \mathcal{K}_{ijt}^c as the set of products displayed to customer i in journey j at step t . The customer faces a decision between: clicking in one a set of products $k \in \mathcal{K}_{ijt}^c$, continue searching ($k = s$), or finish the purchase journey without buying ($k = \ell$). We denote by $y_{ijt}^c = k^*$ the decision made at step t , with $k^* \in \mathcal{K}_{ijt}^c \cup \{s, \ell\}$ being the

²⁹For the case of one-way journeys, only “One-way results page” can occur.

alternative that maximizes utility u_{ijtk}^c , such that

$$\begin{aligned}
u_{ijtk}^c &= \alpha_{1ijp(t)} + \mathbf{x}'_k \cdot \text{diag}(\boldsymbol{\rho}) \cdot \boldsymbol{\beta}_{ij} + \mathbf{Z}'_{ijkt} \cdot \boldsymbol{\eta} + \varepsilon_{ijtk}, & \text{for all } k \in \mathcal{K}_{ijt}^c, \\
u_{ijts}^c &= \alpha_{2ijp(t)} + \varepsilon_{ijtk}, \\
u_{ijtl}^c &= \varepsilon_{ijtk},
\end{aligned} \tag{2.5}$$

where $\varepsilon_{ijtk} \sim \mathcal{N}(0, 1)$ are i.i.d. unobserved (to the researcher) components of utilities. The term $\alpha_{1ijp(t)}$ is the intercept that captures the base rate in page type $p(t)$ for continuing in the journey by either clicking or searching as opposed to leaving the journey without a purchase, and $\alpha_{2ijp(t)}$ is the intercept that captures the base rate in page type $p(t)$ for clicking as opposed to searching again. We define $\boldsymbol{\alpha}_{0ij} = [\alpha_{1ij1}, \dots, \alpha_{1ijP}, \alpha_{2ij1}, \dots, \alpha_{2ijP}]$ the vector of click intercepts for all types of pages. The vector \mathbf{x}_k denotes the observed attributes of product k ,^{30,31} $\boldsymbol{\beta}_{ij}$ is the vector of product attributes preferences, and $\boldsymbol{\rho}$ is a vector of the same length as $\boldsymbol{\beta}_{ij}$ that captures how product attributes preferences manifest differently in clicks relative to purchases.³² Setting $\boldsymbol{\rho}$ to $\mathbf{1}$ corresponds to assuming that the preferences for attributes affect clicks and purchases decisions in the exact same way. We also allow for a set of controls \mathbf{Z}_{ijkt} that may affect click decisions, where $\boldsymbol{\eta}$ captures their impact on clicks.

With reference to the notation introduced in Figure 2.5, $\boldsymbol{\alpha}_{0ij}$ and $\boldsymbol{\beta}_{ij}$ are part of the component \mathbf{b}_{ij} . The vector of product attribute preferences $\boldsymbol{\beta}_{ij}$, is a key parameter in our

³⁰Formally, the vector \mathbf{x}_k also depends on customer i , journey j and step t , but we drop the explicit dependencies of these indexes to keep the notation simpler.

³¹To mimic the customer behavior, we set to zero the attribute levels of \mathbf{x}_k that cannot be observed in step t , this is, inbound attributes for $p(t) = \text{“Outbound page”}$, and outbound attributes for $p(t) = \text{“Inbound page”}$.

³²In practice, we can only identify $\boldsymbol{\rho}$ relative to the differences in scales of the error terms in the click and purchase models. That is, if $\boldsymbol{\rho}$ is a vector where all its entries are equal to the same value, we cannot distinguish $\boldsymbol{\rho}$ from differences in variance of the unobserved component of utility for the click and purchase models.

model since, once inferred, it will allow us to recommend products to customers in ongoing journeys. These preferences are both customer and journey specific. We will describe in Section 2.4.2.5 how these preferences are related across journeys from the same customer as well as across journeys from other customers with the same journey context. These preferences will also play a role in our purchase model, which we describe next.

2.4.2.4 Purchase model

The customer can either buy a product from a subset of displayed products, or not purchase at all from the website.³³ We model the likelihood of purchase as a single decision that happens once per journey, where the customer chooses among a set of considered products.

Formally, consider customer i in journey j . We index by $k \in \mathcal{K}_{ij}^p$ the products considered for purchase by customer i in journey j , and by $k = 0$ the no-purchase in the website outside option. When we consider the journey as a whole, the customer decides between purchasing one of the products $k \in \mathcal{K}_{ij}^p$ and not purchase any product at all ($k = 0$). We denote by $y_{ij}^p = k^*$ the purchase decision of customer i in journey j , where k^* is the alternative in $\mathcal{K}_{ij}^p \cup \{0\}$ that maximizes utility u_{ijk}^p , with

$$\begin{aligned} u_{ijk}^p &= \tau_{0ij} + \mathbf{x}_k' \cdot \boldsymbol{\beta}_{ij} + \epsilon_{ijk}, & \text{for all } k \in \mathcal{K}_{ij}^p, \\ u_{ij0}^p &= \epsilon_{ijk}, & \end{aligned} \tag{2.6}$$

³³Our dataset does not allow us to distinguish between the customer purchasing the product in another website and no purchasing the product all together. However, ultimately the goal for the firm is that the customer buys in their website.

where $\epsilon_{ijk} \sim \mathcal{N}(0, 1)$ are i.i.d. unobserved components of utilities of the products. The element τ_{0ij} is the intercept for purchasing a product, the vector \mathbf{x}_k contains the attributes of product k , and β_{ij} is the vector of product attribute preferences described in the click model.

Note that, unless the customer makes a purchase, it is difficult to disentangle whether s/he has already decided not to purchase from the focal firm, or s/he might do so in the future, making the purchase outcome partially unobserved. In our setting, we can determine that a customer decided not purchase for many of the journeys because for many journeys the start of the trips is included in our data period.

2.4.2.5 Combining different sources of information

One of the key objectives of our modeling effort is to combine different sources of information from the customer journey—i.e., being able to learn from queries, clicks, and purchases—while recognizing that customers might exhibit journey-specific preferences. That is, a customer may exhibit different behavior when looking for a flight domestically versus internationally, or when flying for leisure versus for business. To capture this behavior, we model these journeys as belonging to one of many *journey contexts*. These are unobserved components that capture need-specific preferences that are shared across customers.

Incorporating context-specific preferences presents several methodological challenges. First, the journey contexts are unobserved and therefore need to be inferred from the data. These journey contexts are not individual-specific as different customers are likely to be searching for similar contexts. It is neither the case that customers systematically “transition” from one context to another (like, for example, in hidden Markov models), challenging identification because individual behavior in the previous journey does not necessarily inform

the context in the current journey. Second, we do not know how many contexts are there — and this number is likely to be different across settings — so ideally, we would like to learn the number of contexts from the data, without the need to run the model for each number of contexts. Finally, we want to provide enough flexibility to the model such that it will be able to capture *meaningful* contexts that reflect both queries and behaviors (e.g., a “summer family trip” context that bundles together journeys that are more likely to be international trips, with more than one adults and with children, which may involve strong preferences for non-stop destinations and moderate price sensitivity). The intuition is that, because these journeys share these characteristics, the customers may be interested in covering similar needs, and therefore, their preferences for products in these journeys may also be similar. A model that only groups queries will not necessarily help understanding what the customer is looking for when it comes to product attributes.

To overcome these challenges, we model the journey context as a non-parametric latent segmentation over journeys across customers, using information from the query variables as well as the preferences of these journeys that drive clicks and purchases. In particular, we allow the query parameters ω_{ij} to depend on the journey context by specifying

$$\omega_{ij} = \gamma_j^q,$$

where γ_j^q is the context specific vector of query parameters. Note that the query parameters relate to the individual through the context that individual i uses in trip j . In contrast, click intercepts α_{0ij} , purchase intercept τ_{0ij} , and product attribute preferences β_{ij} are both customer- and journey context-specific by combining the context preferences and the

customer specific stable random-effect parameters in an additive manner. Specifically,

$$\mathbf{b}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\gamma}_j^p, \quad (2.7)$$

where \mathbf{b}_{ij} is the vector of all click and purchase parameters ($\mathbf{b}_{ij} = (\boldsymbol{\alpha}_{0ij}, \tau_{0ij}, \boldsymbol{\beta}_{ij})$), $\boldsymbol{\mu}_i$ is the individual-specific vector of click and purchase parameters, and $\boldsymbol{\gamma}_j^p$ is the context-specific vector of click and purchase parameters. Finally, we define $\boldsymbol{\gamma}_j$ as the vector of all context specific parameters,

$$\boldsymbol{\gamma}_j = \left[\underbrace{\boldsymbol{\gamma}_j^q}_{\text{Query parameters}} \quad \underbrace{\boldsymbol{\gamma}_j^p}_{\text{Click and purchase preferences}} \right]. \quad (2.8)$$

We account for customer heterogeneity by modeling the individual specific vector of parameters

$$\boldsymbol{\mu}_i \sim \mathcal{N}(0, \Sigma), \quad (2.9)$$

where Σ is the covariance matrix. We center the individual specific vector of parameters $\boldsymbol{\mu}_i$ at zero, in order to leave the context-specific vector $\boldsymbol{\gamma}_j^p$ to capture the population mean.

2.4.2.5.1 Modeling contexts

We estimate contexts non-parametrically assuming that $\boldsymbol{\gamma}_j$ are drawn from an unknown discrete distribution F , which we call the context distribution (e.g., a histogram of contexts).

We assume that this histogram F is drawn using a Pitman-Yor Process prior. The Pitman-Yor Process (Pitman and Yor, 1997) is a distribution over infinite almost surely

discrete measures (e.g., infinite histograms) used in non-parametric Bayesian models. Thus, we draw the context specific parameters γ_j from the context distribution F , and we place a Pitman-Yor process prior on the context distribution F , that is,

$$\gamma_j \sim F \tag{2.10}$$

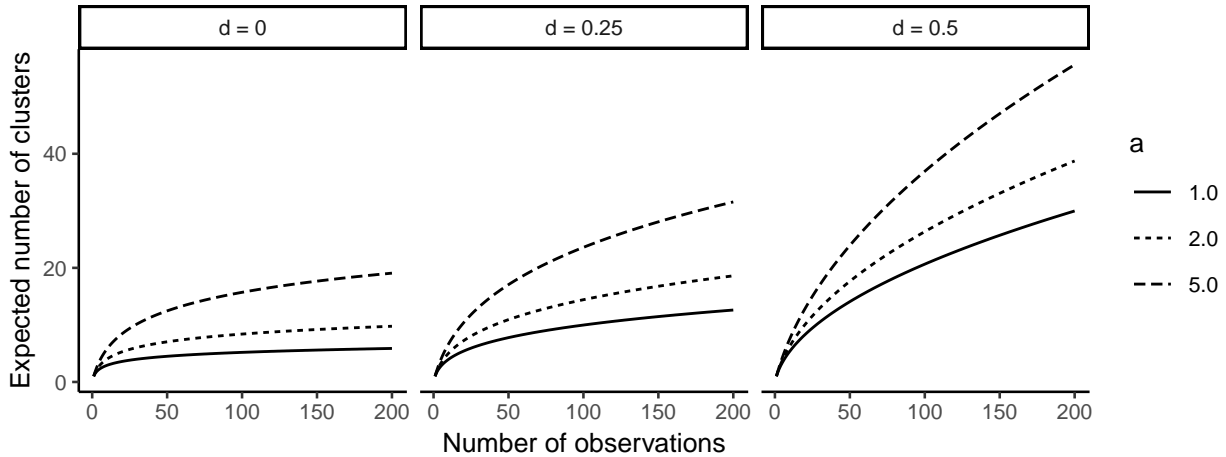
$$F \sim \text{PY}(d, a, F_0), \tag{2.11}$$

where $0 \leq d < 1$ is a discount parameter, $a > -d$ is a strength parameter, and F_0 a base distribution over the same space as γ_j , such that F_0 is the mean distribution of F .

Pitman-Yor processes generalize Dirichlet processes; in particular, when $d = 0$, the Pitman-Yor process reduces to a Dirichlet process with concentration parameter a and base distribution F_0 (i.e., $\text{PY}(0, a, F_0) = \text{DP}(a, F_0)$). This additional parameter allows the drawn distributions from a Pitman-Yor process to exhibit a power-law, long-tail distributions of weights in the histogram, as opposed to histograms with weights decaying exponentially when drawn from Dirichlet processes. This means that the Pitman-Yor process allows for more distinct mass points in the drawn histogram to appear as new observations come in. In particular, this feature of the Pitman-Yor process allows the model to capture new contexts that may not have been observed before, contexts that may happen rather infrequently. In Figure 2.6, we show that as more observations come in, the expected unique number of clusters grows for both models, but it grows rapidly for the Dirichlet process, and then it stops growing significantly ($d = 0$). In contrast, the Pitman-Yor process allow for more flexible patterns of how these unique clusters appear in the data. Moreover, using a

Pitman-Yor process as a prior for our context distribution, similarly to the Dirichlet Process, allows our model to infer the number of contexts directly from the data.

Figure 2.6: Expected number of clusters from a Dirichlet Process ($d = 0$, left) vs. a Pitman-Yor process ($d = 0.25$, middle; and $d = 0.5$, right)



We express the context distribution F in terms of the stick-breaking representation of the Pitman-Yor process (Ishwaran and James, 2001),

$$F = \sum_{c=1}^{\infty} \pi_c \delta_{\theta_c}(\cdot) \quad (2.12)$$

$$\theta_c \sim F_0 \quad (2.13)$$

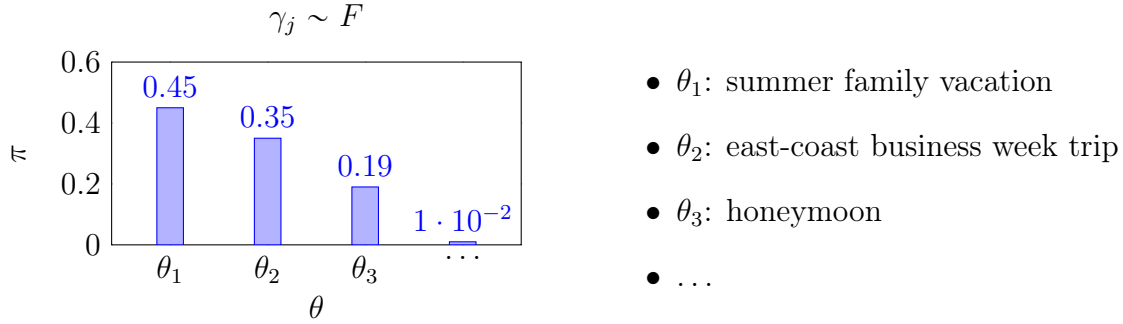
$$\pi_c = V_c \prod_{h=1}^{c-1} (1 - V_h) \quad (2.14)$$

$$V_c \sim \text{Beta}(1 - d, a + c \cdot d). \quad (2.15)$$

This representation allows us to provide some intuition on how this model captures the contexts non-parametrically, which we illustrate using Figure 2.7.

The distribution F acts as a histogram of contexts, where each location $c = 1, 2, \dots$ of the histogram represents a different context (e.g., the summer family vacation, the east-coast

Figure 2.7: Example of a context distribution drawn from a Pitman-Yor process prior



business week trip). For any new journey that a customer undertakes, the model would draw its journey specific parameter γ_j from this histogram of contexts F . The histogram has two main set of parameters, the location θ_c and the context size π_c . The locations θ_c indicate the set of query, click, and purchase preferences that are associated with context c . The context size π_c represent how likely is context c to be drawn. For example, in Figure 2.7, θ_1 represents the set of query parameters, and click and purchase preferences associated with the summer family vacation context. Accordingly, θ_1 would be such that a query is more likely to consist of more than one adult and children, longer than average stays, and farther destinations. At the same time, θ_1 would capture stronger preferences for non-stop flights, and moderate price sensitivity. The value $\pi_1 = 0.45$, represents that when a new journey is drawn, it is of context “summer family vacation” with probability 0.45. Similarly, the “east-coast business week trip” context is drawn with 0.35 chance and so forth.

In summary, our model is able to infer journey-specific preferences using short historical purchase data. It does so by combining queries, clicks, and purchases, both across purchase journeys within a customer and across customers with similar purchase journeys.

In the next section we briefly describe how we specify certain variables in the model tailored to our empirical context and finish the section outlining our estimation procedure.

2.4.3 Specification of query variables, covariates, and sensitivities

With respect to the query component of the model (q_{ijt} in Equations (2.1) through (2.4)), we use all binary and all categorical variables presented in Table 2.1. We apply a log-transformation to the query variable “Trip distance” and model it using a Gaussian distribution according to Equation (2.3). We model “Time in advance to buy”³⁴ using the exponential distribution in Equation (2.4). Finally, we model the query variable “Length of stay” using Equation (2.4), allowing for missing values for one-way journeys.

With respect to the covariates in the click and purchase sub-models (X_{ijt} in Equations (2.5) and (2.6)), we transform price (and length) attributes using $f(x) = \log(1 + x - \min\{x\})$, where $\min\{x\}$ is the minimum price (length) of all displayed offers for that particular click/purchase occasion. We further standardize these variables, by subtracting the mean and dividing by the standard deviation. In addition, we set the following base levels for their corresponding categorical attributes: “One stop” for number of stops, “OneWorld (American)” for alliance, and “Morning (6:00am-11:59am)” for departing and arrival times. We acknowledge that alternatives being ranked at the top may have higher probability of being clicked and bought (Ursu, 2018). However, as the firm ranks products by sorting increasingly by price, we cannot include both price and rank order in the model, as we cannot disentangle the effect of each of these. Therefore, our price coefficient

³⁴We transform the variable using $f(x) = x + 1$, to avoid a zero-valued variable as some journeys search for same-day flights.

captures the effect of both price and ranking. This is a limitation of our dataset, not our model. If products were sorted randomly, we could incorporate both ranking and price simultaneously in our model to control for higher click rates of products displayed at the top of the list. That being said, our price coefficient is the appropriate effect of price in our setting, in which flights are ordered by price.

We further control in our click model (\mathbf{Z}_{ijt} in (2.5)) for products that were previously clicked, in order to capture that in click occasions in which a customer revisits a results page, some products were already clicked before. Theoretically, the sign of this control should be negative, as most theoretical models rule out multiple clicks on the same product, which should translate in lower probabilities of clicking on this products again. Our data shows the opposite pattern — customers are more likely to click on a product they have clicked on before.

Regarding the sensitivity to these covariates, we assume that customers have the same preferences for outbound and inbound legs for length of the leg, number of stops, and alliance and allow them to have different preferences for departing and arrival times, for outbound and inbound legs. Finally, for computational convenience, we set μ_i to zero for all click intercepts and controls (α_{0ij}). One could easily allow for heterogeneity in this components; our results are robust to this change.

2.4.4 Model estimation and prediction for partially observed journeys

We estimate the model parameters using a hierarchical Bayesian framework. We draw from the posterior distribution using Markov Chain Monte Carlo algorithm, specifically, a Gibbs sampler, as we choose priors that allow to us compute full conditional distributions analytically for all parameters, and draw them sequentially. We use the stick-breaking representation of the Pitman-Yor process,³⁵ and we use a blocked Gibbs sampler to be able to implement a fast sampler that can draw the context parameters in parallel for each journey (Ishwaran and James, 2001). We estimate our model using 4,000 iterations of burn-in and 2,000 iterations for drawing the posterior distribution. We assess convergence by observing the mixing distribution of the parameters. Once we have obtain a sample from the posterior distribution, we compute predictions for partially observed journeys. That is, using the queries and clicks of an existing journey, we predict whether the customer will end this journey with a purchase, and if so, which product s/he will buy.

We draw a sample of 4,500 customers for calibrating the model and leave the remaining 499 customers for evaluating the model’s performance on “new customers.” In addition to explore the model performance on hold-out customers, we are predicting ongoing journeys of existing customers. Therefore, for customers with three purchase occasions, we leave the last journey out; and for customers with four or more purchase occasions, we leave the last two journeys out. We will use these held-out journeys to evaluate prediction of purchase incidence and product choice using different depth of information of the journey.

³⁵The results in this essay are estimated using $d = 0$.

2.5 Results

2.5.1 Model estimates

We start by describing mean level preferences for product attributes as well as base intercepts of the click and purchase models. As these parameters vary at the journey level, both across customers as well as across contexts (recall $\mathbf{b}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\gamma}_j^p$ from (2.7)), we compute the population mean estimates of these parameters averaging \mathbf{b}_{ij} across journeys, and $\boldsymbol{\sigma}_b = \sqrt{\text{diag}(\boldsymbol{\Sigma})}$, the standard deviation capturing the level of customer heterogeneity in Equation (2.9).³⁶ Results are presented in Table 2.4. As expected, we find that, on average, customers prefer offers with lower prices and flights (both outbound and inbound) of shorter length. We find that customers do not prefer significantly more non-stops over one stops, but they strongly dislike offers with two or more stops. On average, customers prefer OneWorld alliance over all other alternatives.³⁷ Finally, customers prefer to depart and arrive in the morning for the outbound leg, while they prefer to depart in the morning and arrive either in the morning or in the evening for the return leg.

³⁶Customer heterogeneity is not the only source of heterogeneity across journeys for \mathbf{b}_{ij} . There is another source of heterogeneity, context heterogeneity, which we describe later.

³⁷This is consistent with American Airlines market share. See <https://www.worldatlas.com/articles/the-largest-airlines-in-the-united-states.html>

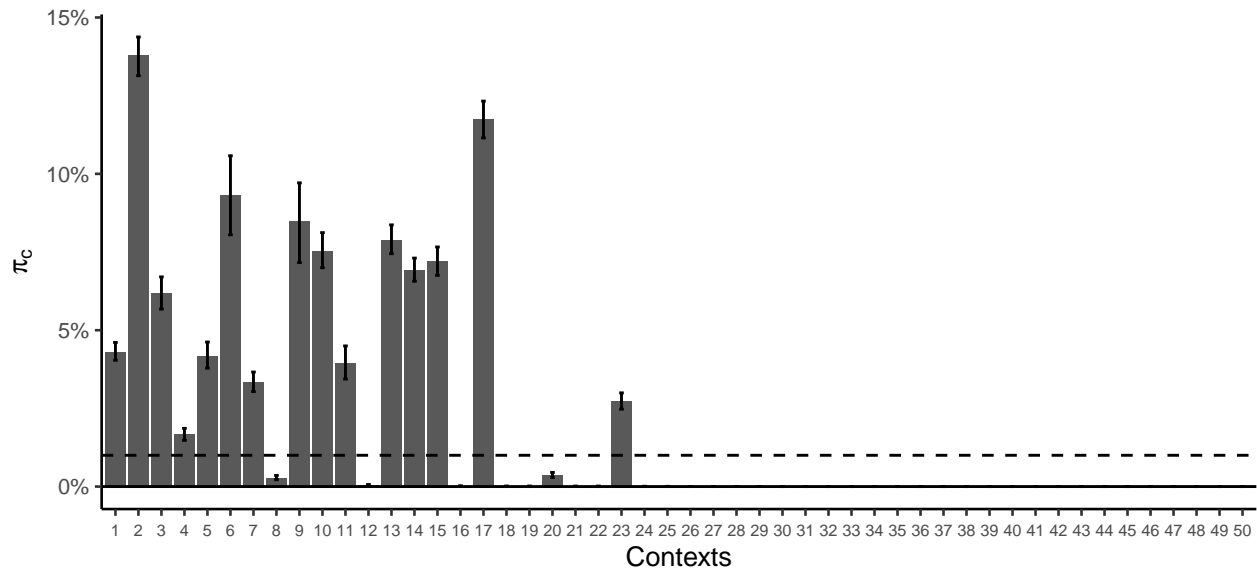
Table 2.4: Parameter estimates of click and purchase models. We show the average across customers and contexts ($\bar{\mathbf{b}}$) and standard deviation across customers (σ_b).

Parameter	Variable	$\bar{\mathbf{b}}$				σ_b			
		Posterior mean	Posterior sd	CPI		Posterior mean			
				2.5%	97.5%				
Click occasions	α_{1ij}	Intercept Click: OW Search	-3.158	0.073	-3.242	-3.001	.		
		Intercept Click: RT Outbound	-3.191	0.053	-3.263	-3.086	.		
		Intercept Click: RT Inbound	-2.563	0.063	-2.671	-2.447	.		
	α_{2ij}	Intercept Search: OW Search	-1.661	0.034	-1.727	-1.597	.		
		Intercept Search: RT Outbound	-2.455	0.027	-2.513	-2.414	.		
		Intercept Search: RT Inbound	-2.651	0.036	-2.726	-2.582	.		
	η	Product was clicked before	1.257	0.050	1.161	1.356	.		
Purchase	τ_{0ij}	Intercept Purchase	-5.550	0.110	-5.669	-5.326	1.054		
		Price	-0.567	0.008	-0.580	-0.551	0.213		
		Length of trip (hours)	-0.737	0.013	-0.766	-0.709	0.339		
		Number of stops: Non stop	0.023	0.018	-0.016	0.052	0.541		
		Number of stops: 2+ stops	-1.621	0.013	-1.652	-1.596	0.580		
		Alliance: Skyteam (Delta)	-0.564	0.049	-0.664	-0.510	0.640		
		Alliance: Star Alliance (United)	-0.367	0.041	-0.452	-0.305	0.897		
		Alliance: Alaska Airlines	-0.497	0.032	-0.561	-0.448	0.641		
		Alliance: Spirit	-0.667	0.045	-0.755	-0.588	0.633		
		Alliance: JetBlue	-0.097	0.042	-0.164	-0.025	0.733		
		Alliance: Frontier	-0.130	0.050	-0.229	-0.039	0.591		
		Alliance: Other – No alliance	-0.228	0.041	-0.329	-0.171	0.941		
		Preferences over attributes	β_{ij}	Alliance: Multiple alliances	-1.542	0.015	-1.571	-1.515	0.590
				Outbound dep. time: Early morning (0:00am - 4:59am)	-0.638	0.040	-0.705	-0.535	0.836
				Outbound dep. time: Afternoon (12:00pm - 5:59pm)	-0.162	0.017	-0.203	-0.140	0.673
				Outbound dep. time: Evening (6:00pm - 11:59pm)	-0.226	0.024	-0.263	-0.179	0.795
				Outbound arr. time: Early morning (0:00am - 4:59am)	-0.835	0.076	-0.943	-0.698	0.823
				Outbound arr. time: Afternoon (12:00pm - 5:59pm)	-0.160	0.025	-0.216	-0.110	0.621
				Outbound arr. time: Evening (6:00pm - 11:59pm)	-0.213	0.017	-0.238	-0.170	0.621
				Inbound dep. time: Early morning (0:00am - 4:59am)	-0.964	0.080	-1.128	-0.848	0.775
Inbound dep. time: Afternoon (12:00pm - 5:59pm)	-0.146			0.019	-0.184	-0.113	0.890		
Inbound dep. time: Evening (6:00pm - 11:59pm)	-0.486			0.049	-0.563	-0.397	0.807		
	Inbound arr. time: Early morning (0:00am - 4:59am)	-0.886	0.155	-1.086	-0.537	0.890			
	Inbound arr. time: Afternoon (12:00pm - 5:59pm)	-0.665	0.047	-0.775	-0.549	0.891			
	Inbound arr. time: Evening (6:00pm - 11:59pm)	-0.078	0.066	-0.240	0.069	0.879			

2.5.2 Contexts in the data

As described in Section 2.4.2.5, the proposed model automatically finds the number of contexts and provides us with a histogram that represents how often each of the contexts appears in the data. We find a very rich set of contexts appearing in the data. Figure 2.8 shows the relative size of each context; of the 19 contexts found in the data (i.e., with at least one journey assigned to them), 15 of these contexts appear in a journey with a probability higher than 1% (dotted line in Figure 2.8).

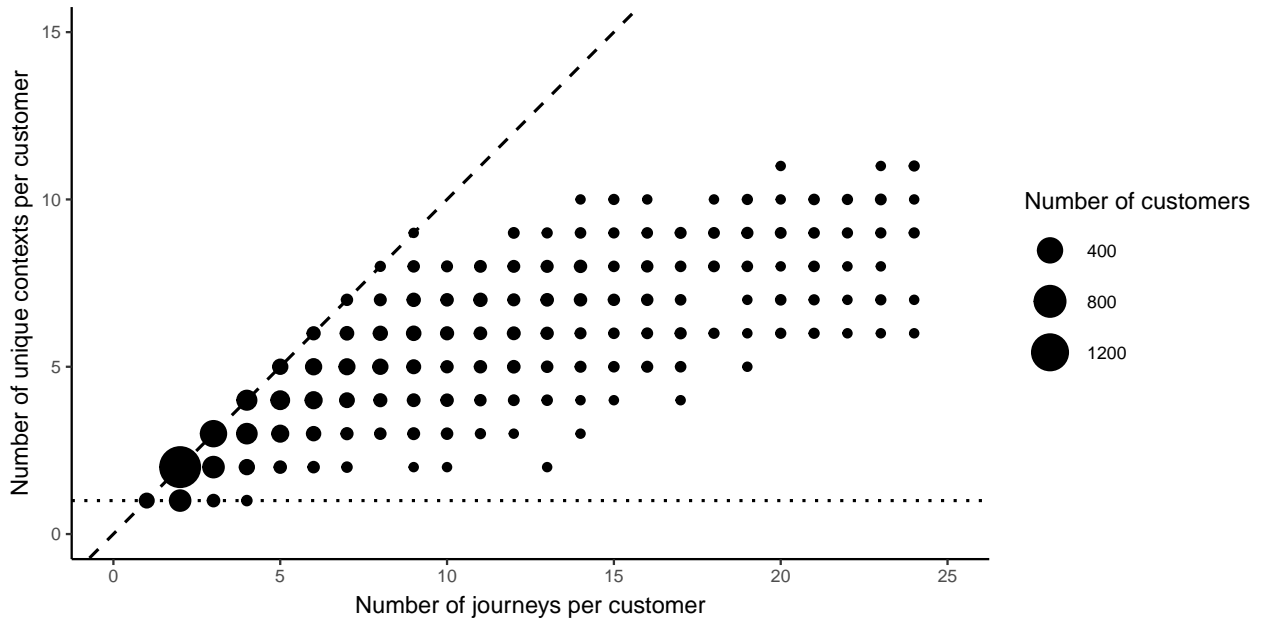
Figure 2.8: Posterior mean and 95% CPI of contexts probabilities, π_c . To be read, an average journey has a 4.5% chance to be of Context 1. (Dotted line marks the 1% chance.)



Arguably, customers could be stable in terms of the contexts their journeys belong to. If that was the case, modeling context heterogeneity would not be necessary as customer heterogeneity could capture those differences. We investigate this issue by computing, per customer, the number of unique contexts his/her journeys belong to. Figure 2.9 shows the distribution of customers (number of customers) by number of unique contexts and number

of journeys per customers.³⁸ The diagonal dashed line represent customers for which all their journeys belong to different contexts. The horizontal dotted line represent customers for which all their journeys belong to a single unique context. We can see that the vast majority of customers have journeys that belong to more than a single context, reinforcing the idea that context heterogeneity and customer heterogeneity capture different sources of variation and that both are present in this setting.

Figure 2.9: Number of contexts per customer. Customers on the horizontal (dotted) line search for flights belonging to only one context, whereas customers above that line search for more than one context. The size of the circle represents the number of customers in that group.

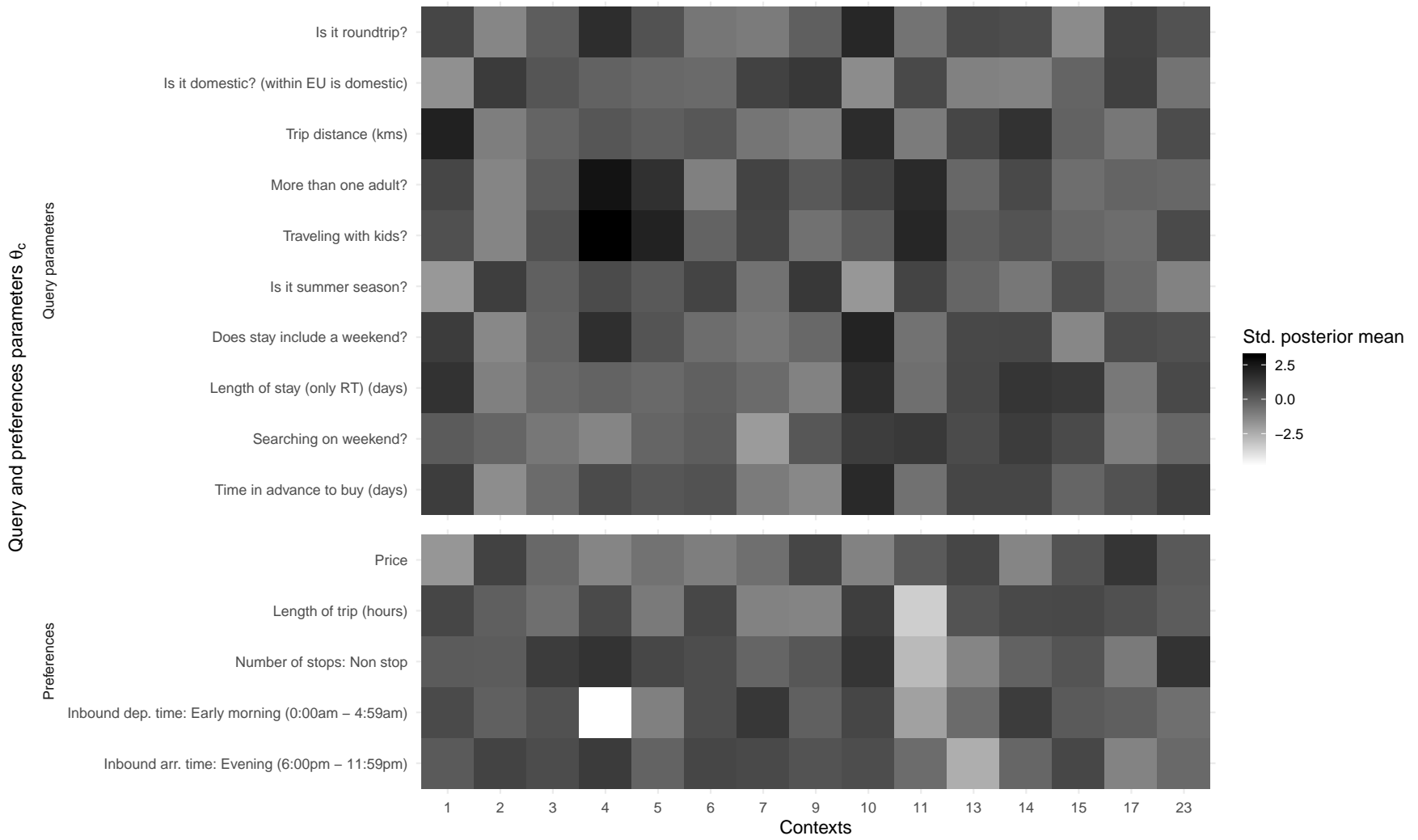


We focus on the 15 most prominent contexts and explore what type of trip each of these contexts is capturing. To do that, we normalize location parameters θ_c (from the stick-breaking representation of the Pitman-Yor process described in (2.13)) which allow us to compare both the query parameters and the flight preferences across the 15 different

³⁸As our sampling algorithm provide us a posterior distribution of the context assignment for each journey, for this exercise, we allocate journeys to the context to its posterior mode.

contexts. First, for each context c , we compute the posterior mean of each location parameter θ_c . Second, we compare these location parameters with the population mean level of those same parameters, similarly to computing $\bar{\mathbf{b}}$, but now we include query parameters as well. We subtract these two, to measure whether contexts are above or below average on each of the query parameters, and click and purchase preferences. Finally, we normalize these differences by dividing by the square-root of the posterior variance across journeys. This variance is composed by two terms (similarly to ANOVA): (1) within-context posterior variance of each θ_c , which measures the posterior uncertainty of each location parameter θ_c ; and (2) the across-context variance of all θ_c with respect to the population mean, which captures how much variance is explained by the differences between contexts. By normalizing the location parameters, we can now compare contexts with respect to whether they score higher or lower than average on each of the query parameters and preferences. We visualize these relative scores in Figure 2.10. A darker color towards black means that the context is higher than average in that dimension; and lighter color towards white means that the context is lower than average in that dimension. Mid-gray color means that contexts does not differ much from the population average. For example, looking at Context 4, this context represents journeys that would likely include more than one adult and kids (compared to the average journey), have a high chance to include a weekend, and where customers show a lower preference of high prices, which means higher price sensitivity.

Figure 2.10: Posterior mean of context location parameters θ_c , relative to the average in the population. The top figure shows how each context deviates from the average with respect to the query variables whereas. The bottom figure shows deviations with respect to the preferences parameters. Darker (lighter) gray means positive (negative) deviation from the average in the population.



Moreover, we identify the top most frequent 50 routes per context (see Figure 2.11) and, combining the insights from these two figures, describe some of these contexts in further detail:

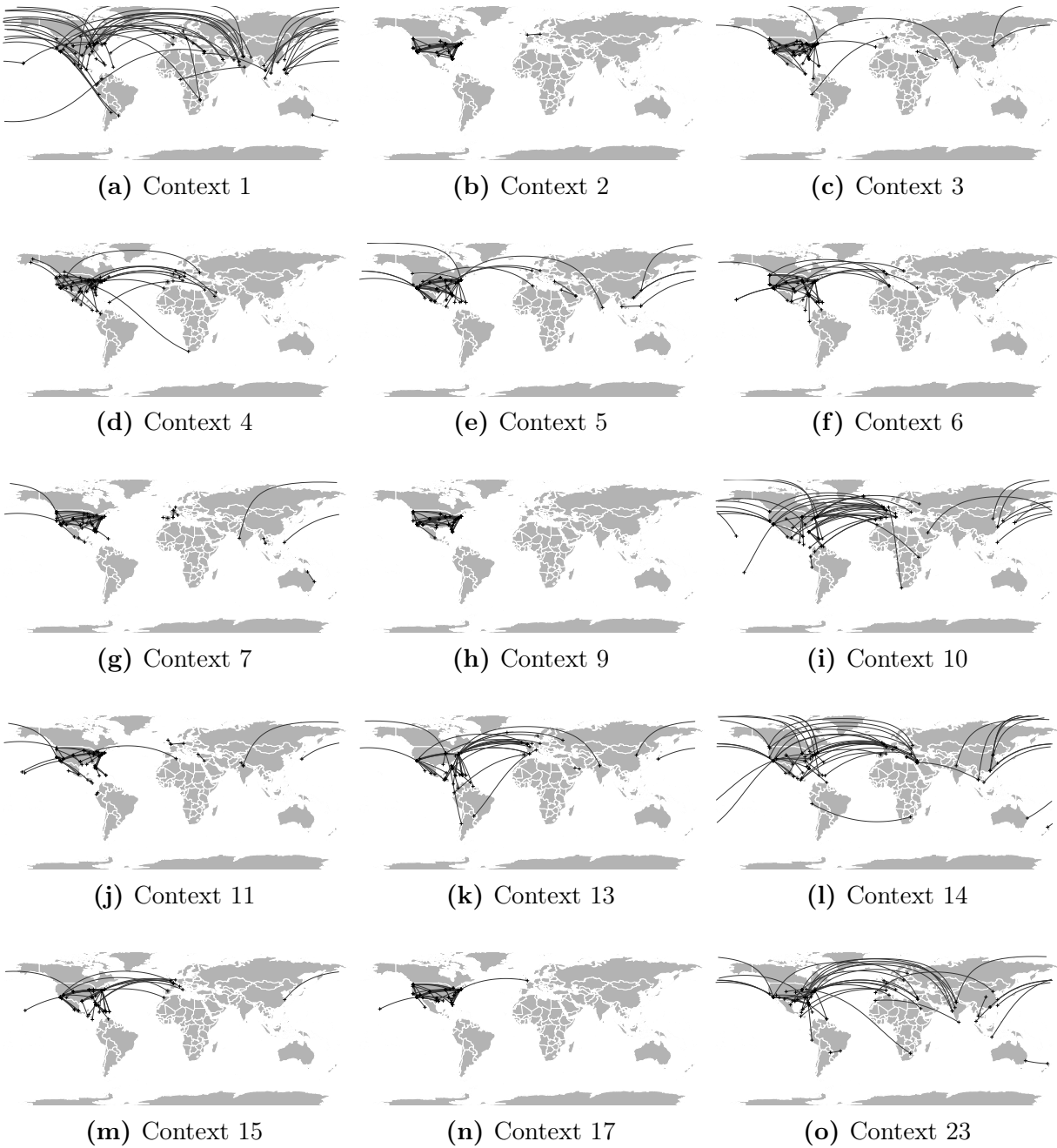
- **Context 4 - No-hassle family vacations (Central America, Europe and Middle East):** This context represents journeys that are more likely to be roundtrip, to include children and other adults, to include weekends. Searches for this type of journeys occur more likely during weekdays. When searching for this type of trips, customers are more price sensitive, probably because they are paying for more tickets, have stronger preferences for non-stop flights, most likely to avoid connections as they travel with kids, and avoid returns that depart in the early morning. Figure 2.11 shows that the top routes of this context include trips between the US Mexico or Central America, as well between the US and Europe, and between the US and the middle east.
- **Context 2 - Short business domestic trip:** This context represents journeys that are more likely to be one-way, domestic and to close destinations in the US. If they are roundtrips, the stays are short and unlikely to include weekends. They neither include kids nor other adults, and searches are made short in advance of the departing date. For these contexts, customers are less price sensitive, and prefer returns that arrive at evening. Figure 2.11 confirms that the top routes of this context are mostly US domestic trips.
- **Context 7 - Close-by family getaway:** This context represents journeys that are more likely to be domestic, and to include another adult and children. Customers search for these trips during the week, and they tend to search with less time in

advance than the average search. Like those in Context 4, customers are more sensitive to price and to longer flights, however they tend to prefer the early morning for the inbound flight. As shown in Figure 2.11, the top routes of this context include mostly trips within the US, and within Europe (with some exceptions to India and Philippines).

- **Context 10 - US overseas winter couple's trip:** This context represents journeys that are more likely to be roundtrip, international, to include another adult, are less likely to occur in the summer-season, and the stays generally include weekends. For these journeys, customers are more price sensitive, less sensitive to longer flights, but with stronger preferences for non-stops. Figure 2.11 confirms that the top routes of this context include trips between the US and Europe, between the US and China or Japan, and between Europe and northern South America.

Not surprisingly, “short business domestic trips” (Context 2) are very prominent in the data (from Figure 2.8, almost 14% of journeys belong to this context) whereas “No-hassle family vacations” and “Close-by family getaways” occur less often (with less than 5% chance each of them). Note that the lower appearance of family vacations in the data does not imply that these contexts are irrelevant, as identifying them earlier on will help the focal firm infer which product offerings will be most preferred in those cases.

Figure 2.11: Top 50 routes per context



To sum, the model identifies several distinct contexts that most customers search for when looking for flights. Not only the model provides valuable information about the characteristics of the contexts and how often these contexts occur in flight searches, but also

identifies what attributes matter the most in each specific context. This information is very relevant for the firm as it can better tailor the needs of the customer as s/he is searching along. Moreover, the firm can leverage this information by identifying cross-selling offerings, either post-purchase or during-search, that enhance the overall trip (e.g., the resort hotel and a tour for the family vacation, or the convenient hotel and a car rental for the business traveler).

2.5.3 Prediction of purchase incidence and product choice

In addition to provide valuable insights about contexts and journey-specific preferences, this model can also be used to predict purchase incidence, and more importantly, product choice. This is relevant when the company wants to predict the likelihood of a purchase outcome in a current journey for example for advertising re-targeting purposes or to prominently display specific flights at the top of the search screen. In particular, we show how the model gets updated as the customer provides more information, by allowing the firm to access his/her past journeys (via cookies), by inserting the query for the current journey, and as the customer clicks on some options as s/he progresses with the journey.

2.5.3.1 Estimated models

We estimate our (*Full*) model, as well as four other benchmark models, all of which are modified versions of our model (see Table 2.5). In the *No context* model, preferences do not depend on context. For this reason, this model informs preferences only through clicks and purchases, and therefore, it does not model queries either. The *No cookies* model ignores which journeys belong to each customer, and therefore does not account for customers'

heterogeneity (preferences are only a function of context). This model could be thought of as the privacy model, in which the website does not store information about customer beyond the focal journey. Finally, the *No clicks* model does not leverage clicks, but it does model queries and purchases, and preferences are a function of both customers' stable preferences and the context of the journey.

Table 2.5: Estimated models

Model name	Queries	Clicks	Heterogeneity	
			Customer	Context
Full	✓	✓	✓	✓
<i>Benchmark models</i>				
No context	.	✓	✓	.
No cookies	✓	✓	.	✓
No clicks	✓	.	✓	✓

2.5.3.2 Prediction of purchase outcomes along the journey

The main idea of this exercise is to explore the models' ability to improve the quality of predictions as the customer moves along the journey and therefore more information becomes available. We consider the two main stages of the journey: *prediction after query*, we assume that the firm is interested in predicting product choice for a pre-specified journey for which the customer has just entered a query.³⁹ The information provided by the customer allows the model to partially infer the context based on the characteristics of the query. While not incorporated directly, the updated context would inform about journey-specific preferences because in our model, contexts are also determined by preferences (learned from other customers who had similar journeys). Then, in the second stage, *prediction after query and*

³⁹Predicting before that point would not be practical because without a query, the firm cannot retrieve a set of product results from which to predict purchase.

clicks, we add the sequence of clicks observed in the current journey (ignoring the click to purchase) and use them to further update the inference of context and therefore preferences. In this case, the information on clicked products not only helps updating the inference about the context, but also helps updating the customer heterogeneity (i.e., stable preferences) as additional clicks can supplement the thin historical data of the particular customers.

Table 2.6: Posterior mean and standard deviations of AUC for prediction of purchase incidence for each model using each piece of information from the customer journey. Higher AUC is better. AUC of 1 corresponds to perfect prediction while AUC of 0.5 corresponds to pure chance.

Model	Holdout journeys	
	Prediction (AUC)	
	...after query	...after query and clicks
Full	0.6043 (0.0102)	0.7406 (0.0066)
No context	0.6003 (0.0090)	0.6359 (0.0047)
No cookies	0.5687 (0.0118)	0.6523 (0.0042)
No clicks	0.6187 (0.0097)	0.6187 (0.0097)

Notes: Higher AUC corresponds to better prediction.

We first analyze the model’s ability to predict purchase incidence, merely, whether the customer will make a purchase in this current journey (Table 2.6). Starting from the middle column, predictions after query, most model specifications have similar predicting power (with AUC of about 0.60) with the *No click* model being slightly superior.⁴⁰ Not surprisingly, the models that either ignore the customer history (*No cookies*) or do not incorporate the query information are the specifications that perform the worst at that stage. Now, when the model starts updating the information collected during the journey the gains from the *Full* model are noteworthy, with the AUC increasing from 0.60 to 0.74. It is important to point

⁴⁰A priori we expected the *Full* model to perform exactly as the *No click* specification because, in principle, they both are inferring from the same information. Our explanation for this slight difference is that as the *No clicks* model cannot leverage clicks, it is trained to extract all its predictive ability from queries, whereas the *Full* model is trained to balance queries and clicks. As clicks are more informative than queries, the *Full* model balances the degree to which queries are exploited for prediction.

out that this increase in prediction power is not coming from knowing that the customer keeps clicking during the current journey (i.e., the fact that customers who click more on the website are more likely to buy and therefore easier to predict). Because if that was the case, all other model specifications, with the exception of *No clicks*, should show the same increase in prediction. Rather, what seems to be happening is that, because the *Full* model is the only one able to update both context and customer heterogeneity separately, doing so allows it to capture journeys' preferences much more accurately. Intuitively, as new information becomes available, the model further pins down each customer's heterogeneity, being able to better identify the specific needs of this particular journey.

Second, we evaluate the model's ability at predicting *the type* of product the customer will buy. The rationale for this analysis is as follows, if the model can indeed infer customers preferences, it should also be able to predict which product the customer will choose. Because in this setting customers are presented with dozens, if not hundreds of results, predicting the exact product that a customer will buy is practically impossible. Therefore, instead, we test whether the model is able to predict the attributes of the product that has been purchased. For example, given that the person purchased a flight in a particular journey, how accurate does the model predict the alliance? Can the model predict the price that the customer paid?

Table 2.7 shows the predictions for four main attributes, namely # stops, alliance, price, and length. (For this analysis we only consider journeys that end up in a purchase.) Across the four product attributes, the *No cookies* and *No clicks* models are those with the poorest performance, both when predicting right after inserting the query and after having incorporated click information. Regarding the prediction right after the query, the *No*

Table 2.7: Posterior mean and standard deviations of hitrate and RMSE (root mean squared error) for predictions of product choice per attribute using each piece of information from the customer journey

Product attribute	Model name	Holdout journeys	
		Prediction (hitrate or RMSE)	
		...after query	...after query and clicks
Number of stops (hitrate)	Full	0.8332 (0.0110)	0.8727 (0.0057)
	No context	0.8360 (0.0096)	0.8516 (0.0061)
	No cookies	0.8066 (0.0096)	0.8118 (0.0057)
	No clicks	0.7862 (0.0122)	0.7862 (0.0122)
Alliance (hitrate)	Full	0.5394 (0.0152)	0.6270 (0.0095)
	No context	0.5406 (0.0124)	0.5746 (0.0074)
	No cookies	0.4707 (0.0110)	0.4960 (0.0082)
	No clicks	0.4758 (0.0165)	0.4758 (0.0165)
Price (RMSE)	Full	0.9553 (0.0280)	0.8291 (0.0120)
	No context	0.9556 (0.0318)	0.8980 (0.0104)
	No cookies	1.0423 (0.0137)	1.0178 (0.0062)
	No clicks	1.0246 (0.0285)	1.0246 (0.0285)
Length (RMSE)	Full	1.4656 (0.0229)	1.3811 (0.0099)
	No context	1.4331 (0.0214)	1.4122 (0.0095)
	No cookies	1.5078 (0.0117)	1.5001 (0.0068)
	No clicks	1.4858 (0.0228)	1.4858 (0.0228)

Notes: Higher hitrate and lower RMSE correspond to better predictions

context model provides the most accurate predictions, with the *Full* being a close second. However, the moment we incorporate the clicks, the *Full* model clearly outperforms all other specifications. Our interpretation for this result is that, not surprisingly, the component capturing customer heterogeneity in product attributes is very much predictive of what attributes the customer will buy — both the *Full* and *No context* specifications incorporate such a component. However, when the model incorporates some clicks from the current journey, not only that information informs individual heterogeneity, but it also helps the model infer more accurately what is the context of the current journey.

Taken all together, these two sets of results suggest that customer historical data, while thin, is very relevant at predicting not only whether the customer will buy, but also which kind of products customers will buy. It also confirms that clicking data contains a very

rich source of information which allows the model to separate between customer and context heterogeneity.

2.5.4 Illustration of how the model infers contexts and preferences along the journey

Finally, we illustrate how the model updates preferences using the different pieces of information collected along the customer journey. For this exercise, we select a customer with two holdout journeys and illustrate how the model updates both the context and the price sensitivity as new information is available to the firm. In the first (holdout) journey, the customer searched for one adult and one child roundtrip from Washington DC to Burlington, VT; departing Friday June 30th, 2017, and returning Wednesday July 5th, 2017. This customer is considering a trip with a child for July 4th's weekend to Burlington, which is a small lake city in Vermont. For the second journey, the customer searched for one adult roundtrip from Washington DC to Knoxville, TN; departing Wednesday September 20th, 2017, and returning Thursday September 21st, 2017. (See the actual queries in Figure 2.12).

Figure 2.12: Queries of two holdout journeys from the same customer

Roundtrip One Way Multi-City

Flying from
Washington, DC (WAS-ALL Airports)

Flying to
Burlington, VT (BTV-Burlington Intl)

Departing
06/30/2017

Returning
07/05/2017

Travelers
1 Adult, 1 Child

Search

(a) Journey 1: WAS (2017 – 06 – 30) – BTV (2017 – 07 – 05)

Roundtrip One Way Multi-City

Flying from
Washington, DC (WAS-ALL Airports)

Flying to
Knoxville, TN (TYS-McGhee Tyson)

Departing
09/20/2017

Returning
09/21/2017

Travelers
1 Adult

Search

(b) Journey 2: WAS (2017 – 09 – 20) – TYS (2017 – 09 – 21)

For each journey, we show how the model updates its inferences for context and price sensitivity, at three different stages of the journey: *Homepage* (before using the query), *Query* (after using the query), and *Query and clicks* (after using the queries and clicks).

We first discuss how the model updates its inference about what context each journey belongs to (see Figure 2.13). At *Homepage* the model does not have any information about the journey; hence, the inference of which context each journey belongs to is equal for these two journeys and corresponds to the average propensities across the population — the small differences between first and second journey are simply due to sampling error. Then, the

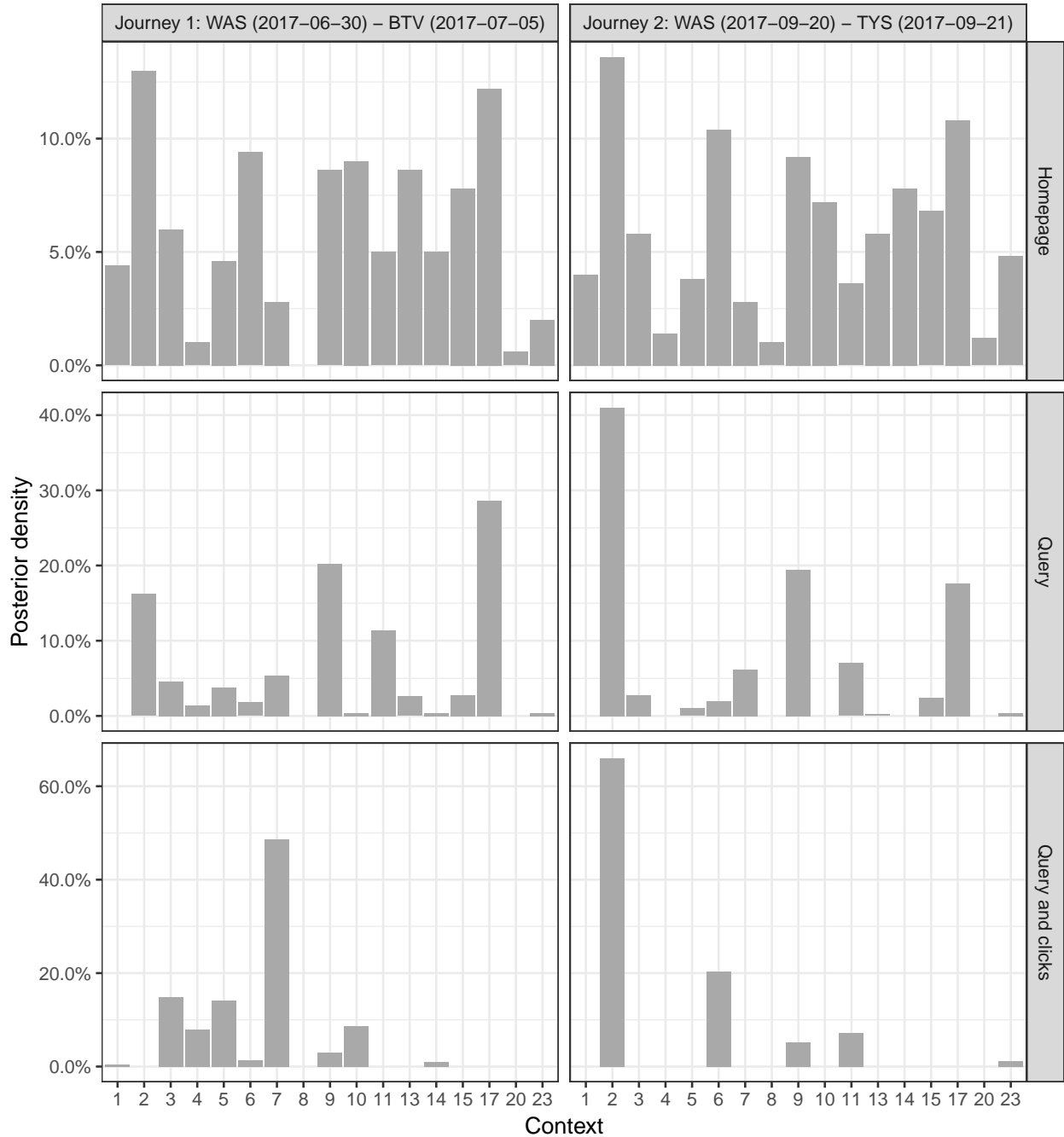


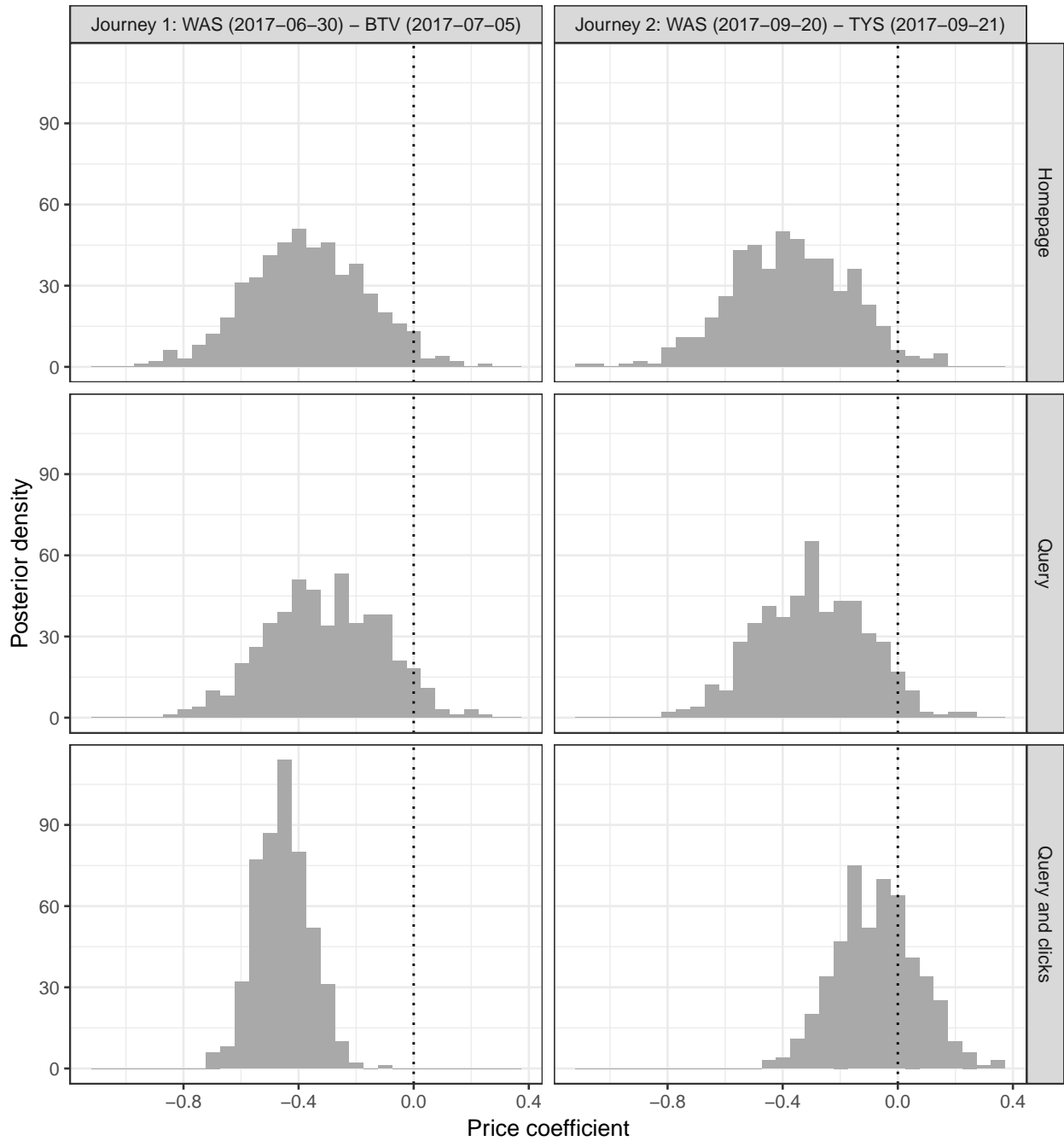
Figure 2.13: Posterior distribution of context for each journey example, using each piece of information from the customer journey.

model incorporates the query information and updates its inference about the context (second row in Figure 2.13). We see that the probability for contexts 2, 9, 11, or 17 increases notably, compared to the *Homepage* step, whereas the probabilities for contexts 1, 10, 13,

and 23 decreases almost to zero. This is not surprising as both trips are within the US and therefore the model infers that these journeys would likely belong to domestic contexts (recall contexts' top destinations from Figure 2.11). Importantly, the inference for these two journeys differ. The second journey, from Washington to Knoxville, is considerably more likely to belong to context 2 (i.e., short business domestic trip), compared to the first journey from Washington to Burlington. In these updates, (similarly to updates of posterior segment membership probabilities in latent class models) the model weights both the information about the query, as well as, how frequent these journeys appear in the population. Finally, once clicks are observed, contexts are updated again (by integrating the information from the clicking component of the model), now showing that the first journey most likely belongs to context 7 (i.e., close-by family getaway), whereas the second journey is more likely to belong to context 2 (i.e., short business domestic trip).

We now show how the model updates its inferences about price sensitivity for these journeys (Figure 2.14). Similarly with context, the inference about price sensitivity is equal for the two journeys at *Homepage*, as they both belong to the same customer and the model is not using any information from the journey. Then, once the model uses each of the queries, inferences are updated, reflecting how the model updates its inference on which contexts these journeys belong to. Because the context inferences for these two journeys are comparable at this stage (middle row in Figure 2.13), both journeys display a similar price coefficient, which is slight larger than that at the *Homepage* stage (both histograms move slightly to the right, particularly the left tail of the distribution). Finally, when the model observes the clicks, the inference about contexts is more certain and different across the two journeys, which results in considerably different inferences on price sensitivity. The first

Figure 2.14: Posterior distribution of price coefficient for each journey example, using each piece of information from the customer journey.



journey, from Washington to Burlington for which the model uncovers a close-by family getaway context, shows a more price sensitive journey (more negative coefficient), compared

to the second journey, from Washington to Knoxville that the model uncovers as a short business domestic trip context.

To sum, this exercise illustrates how the model updates its inferences as new information from an active customer journey is incorporated. It empirically shows that even for different journeys of the same customer, the model can infer different journey-specific preferences, highlighting how relevant is to model context heterogeneity in this setting.

2.6 Conclusion and discussion

We propose a Bayesian non-parametric model for query, click and purchase, to infer customer preferences in settings where historical purchase data is thin. Our model leverages historical data from the previous journeys, data collected during the current journey, and information from other customers with similar journeys. The model accounts for what we call *context heterogeneity*, which are journey-specific preferences that depend on the context in which the journey is undertaken. We model the (unobserved) contexts using a Pitman-Yor process that allow us to uncover non-parametrically the relevant contexts under which customers undertake purchase journeys.

Applying the model to data from one large travel website, we identify 19 different contexts that clearly differ in the specific needs customers are trying to satisfy. For example, one context prevalent in the data is “short business domestic trip.” This context is characterized by trips between close locations, mostly within the US, including one single passenger, lower price sensitivity, and stronger preferences for evening arrivals of the returning trip. In contrast, many other purchase occasions belong to a very different context,

which we call “No-hassle family vacations.” Unlike the business contexts, customers in family trips travel with other adults and kids, look for non-stop flights, have higher price sensitivity, and avoid flights at early morning when returning from their vacations destination.

Interestingly, the same customer searches for different contexts in different points in time.

We find that, among customers who have more than one journey, the average number of contexts they have searched for is 3.3. This figure confirms that context- and customer heterogeneity capture different variation and that both are important drivers of behavior.

Our model and findings are relevant to other industries as well. Experiential purchases such as hotel stays, restaurants reservations, food delivery and media consumption often involve purchase journeys with varying contexts and needs. Firms in these industries collect extensive data from the customer journey, very similarly to the example in our empirical application. Moreover, setting with high involvement products such as cars and durable goods, also exhibit thin individual purchase history. To the extend firms can observe behaviors along the purchase journey, our research suggest that those insights will be very valuable to infer customers’ individual preference.

Our research is not free of limitations. First, our findings regarding context heterogeneity, and particularly the substantive characteristics of such contexts, are based on a subset of highly active customers. Arguably, both the number of contexts as well as their characteristics may vary when using a more representative sample of the customers of the firm, and therefore, these findings should not be taken as representative of the population of customers that search for flight tickets in this market. Second, we do not observe when customers use sorting or filtering tools in the website. If we were to observe these, we could extend our model to further inform both context and customer heterogeneity by modeling

such behaviors as outcomes in our model, similarly as we do with the queries and clicks. Finally, these websites often offer complementary products that are searched for by the customer to fulfill related needs to those in our settings. For example, customers may also search for hotels and car rentals for the same trip they are searching for flights tickets. It may be useful to leverage the information from the context in one category to inform the others. Our model could be extended to share information on related journeys on different categories if such data was available.

References

- Allenby, G. M. and Rossi, P. E. (1998). Marketing models of consumer heterogeneity. *Journal of Econometrics*, 89(1-2):57–78.
- Anderson, E., Chaoqun, C., Israeli, A., and Simester, D. (2020). Do harbinger products signal which new customers will stop purchasing?
- Ansari, A. and Mela, C. F. (2003). E-customization. *Journal of Marketing Research*, 40(2):131–145.
- Artun, O. (2014). What are those new holiday customers worth? [Online; accessed 5-February-2017] <https://www.internetretailer.com/2014/12/19/what-are-those-new-holiday-customers-worth>.
- Ascarza, E. (2018). Retention futility: Targeting high-risk customers might be ineffective. *Journal of Marketing Research*, 55(1):80–98.
- Belk, R. W. (1975). Situational variables and consumer behavior. *Journal of Consumer Research*, 2(3):157.
- Bishop, C. M. (1999). Bayesian PCA. In *Advances in neural information processing systems*, pages 382–388.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- Boughanmi, K., Ansari, A., and Kohli, R. (2019). Dynamics of musical success: A Bayesian nonparametric approach.
- Braun, M. and Schweidel, D. A. (2011). Modeling customer lifetimes with multiple causes of churn. *Marketing Science*, 30(5):881–902.
- Braun, M., Schweidel, D. A., and Stein, E. (2015). Transaction attributes and customer valuation. *Journal of Marketing Research*, 52(6):848–864.
- Bronnenberg, B. J., Kim, J. B., and Mela, C. F. (2016). Zooming in on choice: How do consumers search for cameras online? *Marketing Science*, 35(5):693–712.

- Bruce, N. I. (2019). Bayesian nonparametric dynamic methods: Applications to linear and nonlinear advertising models. *Journal of Marketing Research*, 56(2):211–229.
- Bucklin, R. E. and Lattin, J. M. (1991). A two-state model of purchase incidence and brand choice. *Marketing Science*, 10(1):24–39.
- Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Li, P., and Riddell, A. (2016). Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1):1–29.
- Chan, T. Y., Wu, C., and Xie, Y. (2011). Measuring the lifetime value of customers acquired from Google Search Advertising. *Marketing Science*, 30(5):837–850.
- Chen, F., Liu, X., Proserpio, D., and Troncoso, I. (2020). Product2vec: Understanding product-level competition using representation learning. *Available at SSRN*.
- Chen, Y. and Yao, S. (2017). Sequential search with refinement: Model and application with click-stream data. *Management Science*, 63(12):4345–4365.
- Datta, H., Foubert, B., and Van Heerde, H. J. (2015). The challenge of retaining customers acquired with free trials. *Journal of Marketing Research*, 52(2):217–234.
- De los Santos, B. and Koulayev, S. (2017). Optimizing click-through in online rankings with endogenous search refinement. *Marketing Science*, 36(4):542–564.
- DeSarbo, W. S., Atalay, A. S., LeBaron, D., and Blanchard, S. J. (2008). Estimating multiple consumer segment ideal points from context-dependent survey data. *Journal of Consumer Research*, 35(1):142–153.
- Dew, R. and Ansari, A. (2018). Bayesian nonparametric customer base analysis with model-based visualizations. *Marketing Science*, 37(2):216–235.
- Dew, R., Ansari, A., and Li, Y. (2020). Modeling dynamic heterogeneity using Gaussian processes. *Journal of Marketing Research*, 57(1):55–77.
- Dickson, P. R. (1982). Person-situation: Segmentation’s missing link. *Journal of Marketing*, 46(4):56–64.
- Dong, X., Morozov, I., Seiler, S., and Hou, L. (2019). Estimation of Preference Heterogeneity in Markets with Costly Search.
- Dubé, J.-P. and Misra, S. (2017). Scalable price targeting. Technical report, National Bureau of Economic Research.
- Duvvuri, S. D., Ansari, A., and Gupta, S. (2007). Consumers’ price sensitivities across complementary categories. *Management Science*, 53(12):1933–1945.
- Dzyabura, D. and Hauser, J. R. (2019). Recommending products when consumers learn their preference weights. *Marketing Science*, (June):mksc.2018.1144.

- Erdem, T. and Keane, M. P. (1996). Decision-making under uncertainty: Capturing dynamic brand choice processes in turbulent consumer goods markets. *Marketing Science*, 15(1):1–20.
- Fader, P. S., Hardie, B. G., and Jerath, K. (2007). Estimating CLV using aggregated data: the Tuscan lifestyles case revisited. *Journal of Interactive Marketing*, 21(3):55–71.
- Fader, P. S., Hardie, B. G. S., and Lee, K. L. (2005). “Counting your customers” the easy way: An alternative to the Pareto/NBD model. *Marketing Science*, 24(2):275–284.
- Fader, P. S., Hardie, B. G. S., and Shang, J. (2010). Customer-base analysis in a discrete-time noncontractual setting. *Marketing Science*, 29(6):1086–1108.
- Fiebig, D. G., Keane, M. P., Louviere, J., and Wasi, N. (2010). The generalized multinomial logit model: Accounting for scale and coefficient heterogeneity. *Marketing Science*, 29(3):393–421.
- Forbes (2015). Big data: A game changer in the retail sector. [Online; accessed 23-September-2017] <https://www.forbes.com/sites/bernardmarr/2015/11/10/big-data-a-game-changer-in-the-retail-sector/>.
- Ghose, A., Ipeirotis, P. G., and Li, B. (2012). Designing ranking systems for hotels on travel search engines by mining user-generated and crowdsourced content. *Marketing Science*, 31(3):493–520.
- Ghose, A., Ipeirotis, P. G., and Li, B. (2014). Examining the impact of ranking on consumer behavior and search engine revenue. *Management Science*, 60(7):1632–1654.
- Ghose, A., Ipeirotis, P. G., and Li, B. (2019). Modeling consumer footprints on search engines: An interplay with social media. *Management Science*, 65(3):1363–1385.
- Gopalakrishnan, A., Bradlow, E. T., and Fader, P. S. (2016). A cross-cohort changepoint model for customer-base analysis. *Marketing Science*, (December).
- Hidasi, B., Karatzoglou, A., Baltrunas, L., and Tikk, D. (2016). Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939*.
- Hoffman, M. and Gelman, A. (2014). The No-U-Turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(1):1593–1623.
- Holbrook, M. B. (1984). Situation-specific ideal points and usage of multiple dissimilar brands. *Research in Marketing*, 7(1):93–131.
- Honka, E. (2014). Quantifying search and switching costs in the US auto insurance industry. *RAND Journal of Economics*, 45(4):847–884.
- Honka, E. and Chintagunta, P. (2017). Simultaneous or sequential? Search strategies in the U.S. auto insurance industry. *Marketing Science*, 36(1).

- Ishwaran, H. and James, L. F. (2001). Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 96(453):161–173.
- Iyengar, R., Ansari, A., and Gupta, S. (2003). Leveraging information across categories. *Quantitative Marketing and Economics*, 1(4):425–465.
- Jacobs, B. J., Donkers, B., and Fok, D. (2016). Model-based purchase predictions for large assortments. *Marketing Science*, 35(3):389–404.
- Karaletsos, T. and Rättsch, G. (2015). Automatic relevance determination for deep generative models. *arXiv preprint arXiv:1505.07765*.
- Kim, J. B., Albuquerque, P., and Bronnenberg, B. J. (2010). Online demand under limited consumer search. *Marketing Science*, 29(6):1001–1023.
- Kim, J. B., Albuquerque, P., and Bronnenberg, B. J. (2011). Mapping online consumer search. *Journal of Marketing Research*, 48(1):13–27.
- Kim, J. B., Albuquerque, P., and Bronnenberg, B. J. (2017). The probit choice model under sequential search with an application to online retailing. *Management Science*, 63(11):3911–3929.
- Kim, J. G., Menzefricke, U., and Feinberg, F. M. (2004). Assessing heterogeneity in discrete choice models using a Dirichlet process prior. *Review of Marketing Science*.
- Knox, G. and van Oest, R. (2014). Customer complaints and recovery effectiveness: A customer base approach. *Journal of Marketing*, 78(5):42–57.
- Koulayev, S. (2014). Search for differentiated products: Identification and estimation. *RAND Journal of Economics*, 45(3):553–575.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. (2017). Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474.
- Kumar, M., Eckles, D., and Aral, S. (2020). Scalable bundling via dense product embeddings. *arXiv preprint arXiv:2002.00100*.
- Lee, L., Inman, J. J., Argo, J. J., Böttger, T., Dholakia, U., Gilbride, T., van Ittersum, K., Kahn, B., Kalra, A., Lehmann, D. R., et al. (2018). From browsing to buying and beyond: The needs-adaptive shopper journey model. *Journal of the Association for Consumer Research*, 3(3):277–293.
- Lewis, M. (2006). Customer acquisition promotions and customer asset value. *Journal of Marketing Research*, 43(2):195–203.
- Liu, J. and Toubia, O. (2018). A semantic approach for estimating consumer content preferences from online search queries. *Marketing Science*.

- Liu, L. and Dzyabura, D. (2017). Capturing heterogeneity among consumers with multi-taste preferences.
- Loupos, P., Nathan, A., and Cerf, M. (2019). Starting cold: The power of social networks in predicting non-contractual customer behavior. *Available at SSRN 3001978*.
- MacKay, D. J. (1995). Probable networks and plausible predictions—a review of practical Bayesian methods for supervised neural networks. *Network: computation in neural systems*, 6(3):469–505.
- Manchanda, P., Rossi, P. E., and Chintagunta, P. K. (2004). Response modeling with nonrandom marketing-mix variables. *Journal of Marketing Research*, 41(4):467–478.
- Mcauliffe, J. D. and Blei, D. M. (2008). Supervised topic models. In *Advances in neural information processing systems*, pages 121–128.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mnih, A. and Hinton, G. E. (2009). A scalable hierarchical distributed language model. In *Advances in neural information processing systems*, pages 1081–1088.
- Mohamed, S., Ghahramani, Z., and Heller, K. A. (2008). Bayesian exponential family PCA. *Advances in neural information processing systems*, 21:1089–1096.
- Montgomery, A. L., Li, S., Srinivasan, K., and Liechty, J. C. (2004). Modeling online browsing and path analysis using clickstream data. *Marketing Science*, 23(4):579–595.
- Neal, R. M. (2012). *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media.
- Neslin, S. A., Taylor, G. A., Grantham, K. D., and McNeil, K. R. (2013). Overcoming the “recency trap” in customer relationship management. *Journal of the Academy of Marketing Science*, 41(3):320–337.
- Netzer, O., Feldman, R., Goldenberg, J., and Fresko, M. (2012). Mine your own business: Market-structure surveillance through text mining. *Marketing Science*, 31(3):521–543.
- Netzer, O., Lattin, J. M., and Srinivasan, V. (2008). A hidden Markov model of customer relationship dynamics. *Marketing science*, 27(2):185–204.
- Padilla, N. and Ascarza, E. (2019). The value of first impressions: Leveraging acquisition data for customer management. *Working paper. Columbia Business School Research Paper*, (17-37).
- Park, S. and Gupta, S. (2012). Handling endogenous regressors by joint estimation using copulas. *Marketing Science*, 31(4):567–586.
- Pitman, J. and Yor, M. (1997). The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *The Annals of Probability*, 25(2):855–900.

- Ranganath, R., Perotte, A., Elhadad, N., and Blei, D. (2016). Deep survival analysis. *arXiv preprint arXiv:1608.02158*.
- Ranganath, R., Tang, L., Charlin, L., and Blei, D. (2015). Deep exponential families. In *Artificial Intelligence and Statistics*, pages 762–771.
- Řehůřek, R. and Sojka, P. (2010). Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA. <http://is.muni.cz/publication/884893/en>.
- Research and Markets (2016). Overview & evolution of the global retail industry. [Online; accessed 23-September-2017] <https://www.researchandmarkets.com/research/tqh2xb/>.
- RJMetrics (2016). The ecommerce holiday customer benchmark. [Online; accessed 5-February-2017] <https://rjmetrics.com/resources/reports/the-ecommerce-holiday-customer-benchmark/>.
- Rossi, P. E., McCulloch, R. E., and Allenby, G. M. (1996). The value of purchase history data in target marketing. *Marketing Science*, 15(4):321–340.
- Ruiz, F. J., Athey, S., and Blei, D. M. (2017). Shopper: A probabilistic model of consumer choice with substitutes and complements. *arXiv preprint arXiv:1711.03560*.
- Sarwar, B., Karypis, G., Konstan, J., and Reidl, J. (2001). Item-based collaborative filtering recommendation algorithms. In *Proceedings of the tenth international conference on World Wide Web - WWW '01*, volume 3, pages 285–295, New York, New York, USA. ACM Press.
- Schmitt, P., Skiera, B., and Van den Bulte, C. (2011). Referral programs and customer value. *Journal of Marketing*, 75(1):46–59.
- Schmittlein, D. C., Morrison, D. G., and Colombo, R. (1987). Counting your customers: Who are they and what will they do next? *Management Science*, 33(1):1–24.
- Schweidel, D. A. and Knox, G. (2013). Incorporating direct marketing activity into latent attrition models. *Marketing Science*, 32(3):471–487.
- Schweidel, D. A., Park, Y.-h., and Jamal, Z. (2014). A multiactivity latent attrition model for customer base analysis. *Marketing Science*, 33(2):273–286.
- Seiler, S. (2013). The impact of search costs on consumer behavior: A dynamic approach. *Quantitative Marketing and Economics*, 11(2):155–203.
- Shaffer, G. and Zhang, Z. J. (1995). Competitive coupon targeting. *Marketing Science*, 14(4):395–416.
- Steffes, E. M., Murthi, B. P. S., and Rao, R. C. (2011). Why are some modes of acquisition more profitable? A study of the credit card industry. *Journal of Financial Services Marketing*, 16(2):90–100.

- Thomadsen, R., Rooderkerk, R. P., Amir, O., Arora, N., Bollinger, B., Hansen, K., John, L., Liu, W., Sela, A., Singh, V., Sudhir, K., and Wood, W. (2018). How context affects choice. *Customer Needs and Solutions*, 5(1-2):3–14.
- Uncles, M. D., East, R., and Lomax, W. (2013). Good customers: The value of customers by mode of acquisition. *Australasian Marketing Journal*, 21(2):119–125.
- Ursu, R. M. (2018). The power of rankings: Quantifying the effect of rankings on online consumer search and purchase decisions. *Marketing Science*, 37(4):530–552.
- Verhoef, P. C. and Donkers, B. (2005). The effect of acquisition channels on customer loyalty and cross-buying. *Journal of Interactive Marketing*, 19(2):31–43.
- Villanueva, J., Yoo, S., and Hanssens, D. M. (2008). The impact of marketing-induced versus word-of-mouth customer acquisition on customer equity growth. *Journal of Marketing*, 45(1):48–59.
- Voigt, S. and Hinz, O. (2016). Making digital freemium business models a success: Predicting customers’ lifetime value via initial purchase information. *Business & Information Systems Engineering*, 58(2):107–118.
- Wang, Y. and Blei, D. M. (2019). The blessings of multiple causes. *Journal of the American Statistical Association*, 114(528):1574–1596.
- Watanabe, S. (2010). Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research*, 11:3571–3594.
- Yoganarasimhan, H. (2019). Search Personalization Using Machine Learning. *Management Science*.

Appendix A: Appendix to Essay 1 - Overcoming the Cold Start Problem of CRM using a Probabilistic Machine Learning Approach

A.1 Augmenting the acquisition characteristics via product embeddings

While one could attempt to directly include the product-level purchase incidence as acquisition characteristics, such an approach would suffer from high levels of sparsity (i.e., unique SKUs are purchased rather infrequently over the first transaction of the customers in the calibration data). Instead, we rely on embedding models that have been developed to overcome the challenge that large “vocabularies” have on computing probabilities of multinomial outcomes. (Specifically, how to efficiently compute/approximate the large denominator of the softmax). As described in Section 1.3.2, we use the transactional data from anonymous customers to create product embedding vectors, i.e., vectors representations of all products available, that captures the nature of products, as perceived by the customers. In essence, we leverage the co-occurrences of products in customers’ baskets to infer similarities across products.

A.1.1 Data processing

The anonymous transactions include 304,497 transactions and 4,730 unique product codes (corresponding to unique SKUs specified by the firm). Many of those product codes are very similar in nature, as they only reflect slight modifications of the exact same product,

different sizes, or travel-size packaging. Because those pieces of information are already captured by the acquisition characteristics (`NewProduct`, `Travel`, and `Size`), we aggregate the product code to unique combinations of product sub-category (e.g., liquid soap, bath, beauty oils) and product line (e.g., shea butter, chamomile, fresh-summer). This characterization of product codes results in 515 unique products in the data.

A.1.2 Word2vec algorithm

To capture latent semantic patterns among products in the same transaction, we use Word2vec, a word embedding method in Natural Language Processing (NLP), to map words into numerical vectors. Word2vec is proposed by Mikolov et al. (2013) who develop two architectures to take advantages of word context: continuous bag-of-words (CBOW) and continuous skip-gram (SG). The first model predicts a word based on its neighbor words, and the second model predicts surrounding words based on a given word. We use the SG model to generate a “product vector.”

More specifically, let $T = \{T_1, T_2, \dots, T_H\}$ be the set of transactions, $Q = \{q_1, q_2, \dots, q_M\}$ be the set of unique products, $V = \{V_{q_1}, V_{q_2}, \dots, V_{q_M} | V_{q_i} \in \mathbb{R}^N\}$ be the set of product vectors. Then, the SG model optimizes V by maximizing the loss function:

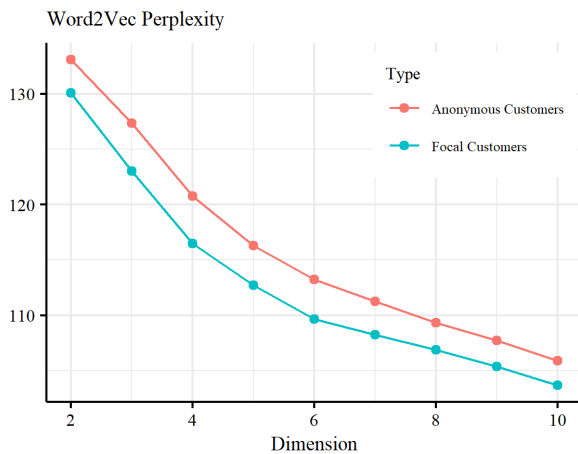
$$L = \sum_{T_i \in T} \sum_{q_i \in Q} \sum_{1 \leq j \leq M, j \neq i} \log P(q_j | q_i), \quad (\text{A.1})$$

where P is the probability of observing product q_j given the occurrence of product q_i in the same transaction. The probability function is defined by the softmax:

$$P(q_j|q_i) = \frac{e^{V_{q_i}^T V_{q_j}}}{\sum_{k=1}^M e^{V_{q_i}^T V_{q_k}}}. \quad (\text{A.2})$$

A straightforward softmax calculation requires an evaluation of all M products in the denominator, so we speed up the computation by using hierarchical softmax (Mnih and Hinton, 2009) to approximate the conditional probability. We implement the model via the Python package Gensim (Řehůřek and Sojka, 2010) and train the model on anonymous customers till the loss L is stable. The hyper parameters in Gensim are: sg=1, negative=0, hs=1, window=10000, min_count=1, random_seed=4. We set a large sliding window size so that all product combinations are selected.

Figure A1: Model selection for Word2vec: Perplexity when varying the number of dimensions from 2 to 10.



We calibrate the Word2vec algorithm using $N = 2, 3, \dots, 10$ dimensions to represent the set of 515 products available in the data and compare the model performance over the

number of dimensions (Figure A1). We select the model with 6 dimensions based on the (lower) rate of decline.¹ As a result, we have a matrix of product embeddings that maps each product to a 6-dimensional vector that represents the position of the product within a multi-dimensional space that captures product similarities.

A.1.3 Interpreting the product dimensions

One could interpret those dimensions by identifying the products that score high in each of the dimensions (Table A1). While not all dimensions are easy to interpret, some clearly capture characteristics defining the nature of the product. For example, looking at the products that score high in the first dimension, we infer that it represents aromas and items for the household. The fifth dimension seems to capture kits and other uncategorized items whereas the sixth dimension represents a specific line of beauty called Fleur Cherie.

In addition to creating the product embeddings that will be used to augment the data, this methodology can also be used to visualize similarities across products. For example, Figure A2 visualizes the 40 most popular products in the anonymous data. Because showing the 6 dimensions would be cumbersome, we apply TSNE (t-distributed stochastic neighbor embedding; algorithm for dimensionality reduction that is well-suited to visualizing high-dimensional data) and visualize the data in a two-dimensional space. It appears to be four clusters representing similarities across these products.

¹A company with a larger product space would calibrate the model with a greater number of dimensions and pick the dimensionality that is best suited for their application.

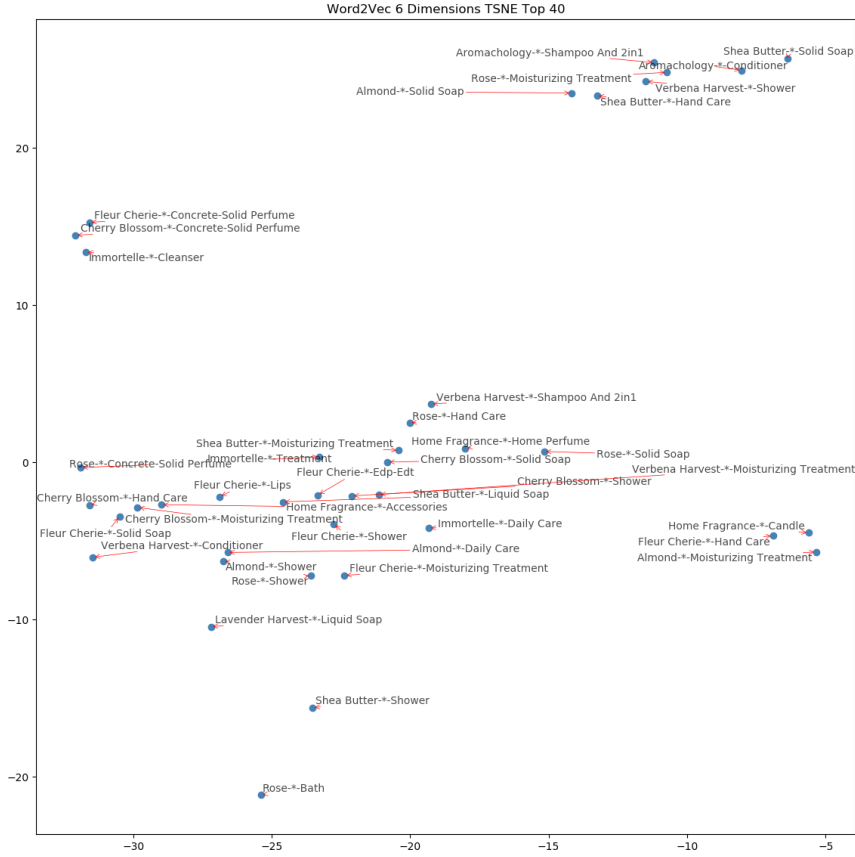
Table A1: Top 5 products per dimension of the product embeddings.

Dimension 1	Dimension 2
Furniture-*-Others	Immortelle-*-Accessories
Aromachology-*-Accessories	Collection De Grasse-*-Accessories
Aromachology-*-Beauty Oils	027-*.Others
Home Fragrance-*-Accessories	Collection De Grasse-*-Shampoo And 2in1
Relaxing Recipe-*-Home Perfume	Verbena Harvest-*.Conditioner
Dimension 3	Dimension 4
Furniture-*-Others	Grape-*.Shower
Orange Harvest-*.Lips	Fleur Cherie-*.Concrete-Solid Perfume
Bonne Mere-*.Others	Olive Harvest-*.Conditioner
Homme-*.Edp-Edt	Shea Butter-*.Body Sun Care
Relaxing Recipe-*.Kits	Grape-*.Body Scrub
Dimension 5	Dimension 6
027-*.Others	Fleur Cherie-*.Solid Soap
Almond-*.Kits	Fleur Cherie-*.Shower
Bonne Mere-*.Kits	Bonne Mere-*.Others
Others-*.Lips	Fleur Cherie-*.Edp-Edt
Immortelle-*.Moisturizing Treatment	Fleur Cherie-*.Moisturizing Treatment

A.1.4 Product mapping for first purchase data

Finally, once the product embeddings are created, we characterize the first purchase from our focal customers by taking the average of the embeddings of each product in the basket (`BasketNature`) and by computing the standard deviation of all products in the basket (`BasketDispersion`), which has missing value if the first purchase only included one product. Note that four products from the first purchase data were not present in the data from anonymous customers and therefore have missing values in the `ProductNature` variable as well.

Figure A2: Visual representation of the product embeddings



A.2 Brief description of DEFs

DEFs are deep generative probabilistic models that describe a set of observations \mathbf{D}_i with latent variables layered following a structure similar to deep neural networks. The lowest layer describes the distribution of the observations, $p(\mathbf{D}_i | \mathbf{z}_i^1, \mathbf{W}^0) = f(\mathbf{D}_i | \mathbf{W}^{0'} \mathbf{z}_i^1)$ and the top layers describe the distribution of the layer just below them. As in deep neural networks, DEFs have two sets of variables: layer variables (\mathbf{z}_i^ℓ) and weights matrices (\mathbf{W}^ℓ) for the ℓ 'th layer. Each layer variable \mathbf{z}_i^ℓ is distributed according to a distribution in the exponential family with parameters equal to the inner product of the previous layer parameters $\mathbf{z}_i^{\ell+1}$ and

the weights \mathbf{W}^ℓ , by

$$p(z_{i,k}^\ell | \mathbf{z}_i^{\ell+1}, \mathbf{w}^\ell) = EXPFAM_\ell \left(z_{i,k}^\ell | g_\ell \left(\mathbf{w}_k^{\ell'} \cdot \mathbf{z}_i^{\ell+1} \right) \right) \quad \ell \in \{1, \dots, L-1\},$$

where $z_{i,k}^\ell$ is the k 'th component of vector \mathbf{z}_i^ℓ , \mathbf{w}_k^ℓ is the k 'th column of weight matrix \mathbf{W}^ℓ , $EXPFAM_\ell(\cdot)$ is a distribution that belongs to the exponential family and governs the ℓ 'th layer, and $g_\ell(\cdot)$ is a link function that maps the inner product to the natural parameter of the distribution, allowing for non-linear relationships between layers. The top layer is purely governed by a hyperparameter η , that is, $p(z_{i,k}^L) = EXPFAM_L(z_{i,k}^L | \eta)$.

Similar to deep unsupervised generative models, DEF models are suitable to find interesting exploratory structure in large data sets. For example, DEFs have been applied to textual data (newspaper articles), binary outcomes (clicks) and counts (movie ratings), being found to give better predictive performance than state-of-the-art models (Ranganath et al., 2015).

A.3 Model priors and automatic relevance determination component

We detail the specification of the automatic relevance determination component that creates sparsity in the weights \mathbf{W}^y , \mathbf{W}^a , and \mathbf{W}^1 and the prior distribution.

A.3.1 Automatic relevance determination

Following Bishop (1999) we define $\boldsymbol{\alpha}$ as a positive vector of length N_1 (number of traits in the lower layer z_i^1), to control the activation of each trait. Note that \mathbf{W}^y is matrix of size $D_y \times N_1$, where D_y is the length of the demand parameters β_i^y ; and \mathbf{W}^a is matrix of size $P \times N_1$, where P is the length of the acquisition parameters β_i^a .

We assume that the component associated with the n 'th row (demand parameter) and k 'th column (trait) of \mathbf{W}^y is modeled by:

$$p(\mathbf{w}_{nk}^y) = \mathcal{N}(\mathbf{w}_{nk}^y | 0, \sigma^y \cdot \alpha_k) \quad (\text{A.3})$$

where σ^y is the parameter that captures the variance of the demand model outcome (e.g., the variance of the error term in a linear regression). For identification purposes, we assume $\sigma^y = 1$ for logistic regressions. For other demand models, σ^y controls the scale of \mathbf{W}^y , and therefore should be defined accordingly. Note that if the vector of covariates \mathbf{x}_{it}^y is not standardized, then this distribution should also consider the scale of the covariates.

Similarly, we model \mathbf{w}_{pk}^a , the component associated with the p 'th row (acquisition behavior) and k 'th column (trait) \mathbf{W}^a , by:

$$p(\mathbf{w}_{pk}^a) = \begin{cases} \mathcal{N}(\mathbf{w}_{pk}^a | 0, \alpha_k) & \text{if } p \text{ is discrete} \\ \mathcal{N}(\mathbf{w}_{pk}^a | 0, \sigma_p^a \cdot \alpha_k) & \text{if } p \text{ is continuous} \end{cases}, \quad (\text{A.4})$$

where σ_p^a is the variance of the error term in the acquisition model for variable p . This variable again corrects for the scale of \mathbf{w}_{pk}^a so it matches the scale of acquisition behavior p .

Finally, note that matrix \mathbf{W}^1 is of size $N_1 \times N_2$. We model \mathbf{w}_{km}^1 , the component associated with the k 'th row (lower layer) and m 'th column (higher layer) of \mathbf{W}^1 , using a sparse gamma distribution:

$$p(\mathbf{w}_{km}^1) = \text{Gamma}(\mathbf{w}_{km}^1 | 0.1, 0.3) \quad (\text{A.5})$$

A.3.2 Model priors

We model the prior distribution of the set of parameters using

$$\begin{aligned} p(\mathbf{W}^y, \mathbf{W}^a, \boldsymbol{\alpha}, \mathbf{W}^1, \boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}^a) &= p(\mathbf{W}^y, \mathbf{W}^a, \boldsymbol{\alpha}, \mathbf{W}^1, \boldsymbol{\mu}^y, \boldsymbol{\mu}^a, \boldsymbol{\sigma}^y, \boldsymbol{\sigma}^a, \mathbf{b}^a) \\ &= p(\mathbf{W}^y | \boldsymbol{\alpha}, \boldsymbol{\sigma}^y) \cdot p(\mathbf{W}^a | \boldsymbol{\alpha}, \boldsymbol{\sigma}^a) \cdot p(\mathbf{W}^1) \cdot p(\boldsymbol{\alpha}) \\ &\quad \cdot p(\boldsymbol{\mu}^y) \cdot p(\boldsymbol{\mu}^a) \cdot p(\boldsymbol{\sigma}^y) \cdot p(\boldsymbol{\sigma}^a) \cdot p(\mathbf{b}^a) \end{aligned}$$

In our estimated models, $\boldsymbol{\sigma}^y$ is a positive scalar σ^y when the demand model is a regression and it does not exist when the demand model is a logistic regression; and $\boldsymbol{\sigma}_p^a$ is a positive

scalar σ_p^a if the p 'th acquisition behavior is continuous, and it does not exist if it is discrete.

We use the automatic relevance determination component, described in Appendix A.3.1, for the terms $p(\mathbf{W}^y|\boldsymbol{\alpha}, \boldsymbol{\sigma}^y)$, $p(\mathbf{W}^a|\boldsymbol{\alpha}, \boldsymbol{\sigma}^a)$, and $p(\mathbf{W}^1)$. Denoting N_{ac} the number of firm-level controls for the acquisition model (i.e., dimension of $\mathbf{x}_{m\tau}^a$), and P_c the number of discrete acquisition variables, we model the remaining terms by:

$$\begin{aligned}
p(\boldsymbol{\alpha}) &= \prod_{k=1}^{N_1} \text{InverseGamma}(\alpha_k|1, 1), \\
p(\boldsymbol{\mu}^y) &= \prod_{k=1}^{D_y} \mathcal{N}(\mu_k^y|0, 5), \\
p(\boldsymbol{\mu}^a) &= \prod_{p=1}^P \mathcal{N}(\mu_p^a|0, 5), \\
p(\mathbf{b}^a) &= \prod_{n=1}^{N_{ac}} \prod_{p=1}^P \mathcal{N}(b_{np}^a|0, 5), \\
p(\boldsymbol{\sigma}^y) &= \log \mathcal{N}(\boldsymbol{\sigma}^y|0, 1), && \text{(if demand model is a regression),} \\
p(\boldsymbol{\sigma}^a) &= \prod_{p=1}^{P_c} \log \mathcal{N}(\sigma_p^a|0, 1) && \text{(A.6)}
\end{aligned}$$

A.4 Rotation of traits

In order to obtain insights about the traits, we post process the posterior sample by carefully rotating the lower weights parameters across draws to define a consistent sign and label of those traits.

First, we define the vectors $\beta_i^{ya} = \begin{pmatrix} \beta_i^y \\ \beta_i^a \end{pmatrix}$, and $\mu^{ya} = \begin{pmatrix} \mu^y \\ \mu^a \end{pmatrix}$ of length $(D_y + P)$, and the matrix $\mathbf{W}^{ya} = \begin{bmatrix} \mathbf{W}^y \\ \mathbf{W}^a \end{bmatrix}$ of size $(D_y + P) \times N_1$. Second, we rewrite (1.5) and (1.6) as:

$$\beta_i^{ya} = \mu^{ya} + \mathbf{W}^{ya} \cdot z_i^1. \quad (\text{A.7})$$

Let D the number of posterior draws obtained using HMC, and $d = 1, \dots, D$ one draw from the posterior distribution. For a sample $\{\mathbf{W}_d^{ya}, \{z_i^1\}_i\}_{d=1}^D$, where traits may switch signs and labels, we are interested in constructing $\{\widetilde{\mathbf{W}}_d^{ya}, \{\widetilde{z}_{id}^1\}_i\}_{d=1}^D$ with “consistent labels and signs”, such that:

$$\mathbf{W}_d^{ya} \cdot z_{id}^1 = \widetilde{\mathbf{W}}_d^{ya} \cdot \widetilde{z}_{id}^1 \quad \forall i, d$$

Intuitively, we are interested in finding the major traits that explain heterogeneity.

In order to build this sample, we use two steps:

1. **Fix labels:**

We obtain the singular value decomposition (SVD) of $\mathbf{W}_d^{ya} = \mathbf{U}_d \cdot \mathbf{D}_d \cdot \mathbf{V}_d'$, where \mathbf{U}_d is an orthogonal matrix of size, $(D_y + P) \times N_1$, \mathbf{D}_d is a diagonal matrix of size $N_1 \times N_1$ with non-negative diagonal values sorted in decreasing order, and \mathbf{V}_d is a orthogonal matrix of size $N_1 \times N_1$. We define $\widehat{\mathbf{W}}_d^{ya} = \mathbf{U}_d \cdot \mathbf{D}_d$, and $\widehat{z}_{id}^1 = \mathbf{V}_d' \cdot z_{id}^1$. Note that we have $\mathbf{W}_d^{ya} \cdot z_{id}^1 = \mathbf{U}_d \cdot \mathbf{D}_d \cdot \mathbf{V}_d' \cdot z_{id}^1 = \widehat{\mathbf{W}}_d^{ya} \cdot \widehat{z}_{id}^1$.

This construction allow us to choose the labels of the traits that explain the most variance in decreasing order, similarly as in Bayesian PCA (Bishop, 2006), which are unlikely to switch across posterior samples for well behaved samples of the product $\mathbf{W}_d^{ya} \cdot z_{id}^1$, which is identified in our model. However, the sign of the traits are not uniquely determined by the SVD. Note that if we multiply by -1 a column of \mathbf{U}_d , and we also multiply by -1 the same corresponding row of \mathbf{V}_d' , then we would also obtain a valid SVD.²

2. Fix signs:

We are interested in fixing a sign for each traits across draws of the posterior distribution, however some trait weights may change sign across the posterior. In order words, the posterior distribution may have its mode close to the origin, and therefore the weights may take values both positive and negative. Therefore, we choose the sign of each trait by observing the behavior it impacts the most (demand or acquisition), and we choose the sign such that the weight of this trait on that behavior does not change sign across draws of the posterior sample.

²Let \widetilde{I} a diagonal matrix of size $N_1 \times N_1$ where each of its diagonal values are either 1 or -1, then we have that $(\mathbf{U}_d \cdot \widetilde{I}) \cdot \mathbf{D}_d \cdot (\mathbf{V}_d \cdot \widetilde{I})' = \mathbf{U}_d \cdot \widetilde{I} \cdot \mathbf{D}_d \cdot \widetilde{I} \cdot \mathbf{V}_d' = \mathbf{U}_d \cdot \mathbf{D}_d \cdot \mathbf{V}_d'$.

More formally, let $k = 1, \dots, N_1$ a trait (a column of \mathbf{W}_d^{ya}), and $n(k)$ the behavior (a row of \mathbf{W}_d^{ya}) that is most impacted by trait k , which we operationalize by computing the posterior mean of the absolute value of \hat{w}_{nk}^{ya} , the weight of trait k on behavior n (i.e., the nk 'th component of matrix $\widehat{\mathbf{W}}^{ya}$), and choosing the maximum:

$$n(k) = \arg \max_{n=1, \dots, (D_y+P)} \left\{ \frac{1}{D} \sum_{d=1}^D \text{abs}(\hat{w}_{nk,d}^{ya}) \right\} \quad (\text{A.8})$$

Then, we change the sign of the trait so $\mathbf{w}_{n(k)k,d}^{ya}$ is always positive, by defining \tilde{I}_d a diagonal matrix of size $N_1 \times N_1$, where its k diagonal value is:

$$(\tilde{I}_d)_{kk} = \text{sign}(\hat{w}_{n(k)k,d}^{ya})$$

Finally, we construct our sample by:

$$\begin{aligned} \widetilde{\mathbf{W}}_d^{ya} &= \widehat{\mathbf{W}}_d^{ya} \cdot \tilde{I}_d && \forall d \\ \tilde{z}_{id}^1 &= \tilde{I}_d \cdot \hat{z}_{id}^1 && \forall i, d \end{aligned}$$

A.5 Algorithm for newly-acquired customers

With reference to (1.10), once we have estimated the full model using the calibration data, we can form first impressions of newly acquired customers using the following procedure:

Algorithm 1 Forming first impressions

Input A sample of the population parameters drawn from the posterior $\{\Theta_m\}_{m=1}^M$
 Acquisition characteristics A_j of focal customer j .
Output A sample of β_j^y drawn from $p(\beta_j^y|A_j, \mathcal{D})$
for all $d \leftarrow 1 : S$ **do**
 Draw $\Theta_d \sim p(\Theta|\mathcal{D})$ from sample $\{\Theta_m\}_{m=1}^M$
 Draw $\mathbf{Z}_{jd} \sim p(\mathbf{Z}_j|\Theta_d, A_j)$ ▷ Using MCMC, HMC or VI
 Compute $\beta_{jd}^y \leftarrow \boldsymbol{\mu}_d^y + \mathbf{W}_d^y \cdot \mathbf{z}_{jd}^1$
end for
Return $\{\beta_{jd}^y\}_{d=1}^S$

Note that the step “Draw $\mathbf{Z}_{jd} \sim p(\mathbf{Z}_j|\Theta_d, A_j)$ ” involves sampling from a posterior distribution for which we do not have access to a closed form distribution. Instead, using the approximation described in (1.10), we use HMC to approximately sample from this posterior for each draw $\Theta_d \sim p(\Theta|\mathcal{D})$. Note that as in this sub-model, only \mathbf{Z}_j of the focal customer j is unknown, an HMC algorithm that samples from this posterior is computationally fast even if this algorithm has to be run inside the loop for each value of d .

A.6 Further details about the simulation analyses

In this appendix we provide further details about the simulation exercise described in Section 1.4.4

A.6.1 Simulation design

We simulate demand and acquisition behavior for 2,200 customers. We first simulate acquisition and demand parameters (β_i^a and β_i^y respectively), and then use those to simulate the observed behaviors (A_i and $y_{i1:T}$ respectively). The data from 2,000 customers will be used to calibrate the models while the remaining 200 individuals will be used to evaluate the performance of each of the estimated models. For those (hold out) customers, we will assume that only the acquisition characteristics are observed, we will use each estimated model to infer customers' demand parameters and then will compare those inferences with the true parameters.

For our simulation study, we assume that acquisition and demand parameters are correlated, that is, observing acquisition behavior can partially inform demand parameters. For this purpose, we generate the individual demand parameters as a function of the acquisition parameters. To cover a variety of relationships among variables we use a linear, quadratic/interactions, and a positive-part (i.e., max) function, therefore exploring linear as well as non-linear relationships. Furthermore, to test whether the model can account for redundancy and irrelevance of variables in the acquisition characteristics collected by the firm, we assume that some acquisition variables are correlated among them and that other acquisition variables are totally independent of future demand. For clarity of exposition and

brevity's sake, we first assume a small number of acquisition variables. Because many empirical contexts will likely have a large number of acquisition variables, we then extend the analysis to incorporate dozens of variables and show how the model performs at a larger scale.

A.6.2 Data generation process

Generate individual-level parameters

First, we generate seven acquisition parameters for seven corresponding acquisition characteristics. In order to resemble what real data would look like, and to test whether our model can account for redundancy in the acquisition data (e.g., the number of items purchased and total amount spent at acquisition being highly correlated), we make some of these acquisition parameters highly correlated among themselves. We operationalize such a relationship by assuming that six of the seven parameters are driven by two main factors $\mathbf{f}_i = \begin{pmatrix} f_{i1} \\ f_{i2} \end{pmatrix}$, where $\mathbf{f}_i \sim N(0, I_2)$. Furthermore, we set the seventh acquisition parameter to be independent of other acquisition parameters as well as independent to future demand parameters. The rationale behind this structure is to resemble the situation in which the acquisition data includes irrelevant data and therefore test whether the model is robust to random noise. More specifically,

$$\begin{aligned}
 \beta_{ip}^a &\sim N(\mu_p^a + B_{1p} \cdot f_{i1}, \sigma_p^{ba}), & p = 1, \dots, 3 \\
 \beta_{ip}^a &\sim N(\mu_p^a + B_{2p} \cdot f_{i2}, \sigma_p^{ba}), & p = 4, \dots, 6 \\
 \beta_{i7}^a &\sim N(\mu_7^a, \sigma_p^{ba}), &
 \end{aligned} \tag{A.9}$$

where β_{ip}^a is the p^{th} component of acquisition vector β_i^a , μ_p^a is the mean of the p^{th} acquisition parameter; B_{1p} and B_{2p} represent the impact of factors 1 and 2 respectively on the p^{th} acquisition parameter; and $\sigma_{ba} = 0.1$ the standard deviation of the uncorrelated variation of the p^{th} acquisition parameter. The values used to generate factors f_{i1} and f_{i2} are presented in Table A2.

Table A2: True values for factors f_{i1} and f_{i2} impact on acquisition parameters (B_{1p} and B_{2p}).

Acquisition parameter	Weight factors	
	B_{1p}	B_{2p}
Factor 1, f_{i1}		
Acq. variable 1	3.0	0.0
Acq. variable 2	2.0	0.0
Acq. variable 3	-2.5	0.0
Factor 2, f_{i2}		
Acq. variable 4	0.0	3.5
Acq. variable 5	0.0	-2.0
Acq. variable 6	0.0	-3.0
Independent		
Acq. variable 7	0.0	0.0

Second, we generate the individual customer parameters for demand; these are the values that the firm is interested in inferring (β_i^y). We generate three parameters governing the demand model: an intercept and two covariate effects. We generate these individual demand parameters β_{ik}^y as a function of the acquisition parameters β_i^a , following a general form

$$\beta_{ik}^y \sim N\left(\mu_k^y + g_k(\beta_i^a | \Omega_k), \sigma_k^{by}\right), \quad k = 1, \dots, 3, \quad (\text{A.10})$$

where $g_k(\beta_i^a | \Omega_k)$ is the function that represents the relationship between acquisition and demand parameters. Because our goal is to investigate the accuracy of the model (compared

to several benchmarks) in contexts in which the relationship between acquisition and demand parameters could take different forms, we vary g_k to capture a variety of scenarios:

- Scenario 1: Linear

$$g_k(\beta_i^a | \Omega_k) = \omega_k^{1'} \cdot \beta_i^a \tag{A.11}$$

This relationship would exist when, for example, customers with a strong preference for discounted products at the moment of acquisition are also more likely to be price sensitive in future purchases.

Table A3: Simulated values for ω_k^1 in the Linear scenario

Variable	Demand variables		
	Intercept	Covariate 1	Covariate 2
ω_{k1}^1	0.30	-0.69	-0.03
ω_{k2}^1	0.86	-0.61	-1.37
ω_{k3}^1	-1.44	-0.35	-0.03
ω_{k4}^1	-0.05	-0.10	0.12
ω_{k5}^1	1.16	-0.06	0.71
ω_{k6}^1	-0.12	0.10	0.93

- Scenario 2: Quadratic/interactions

$$g_k(\beta_i^a | \Omega_k) = \omega_k^{1'} \cdot \beta_i^a + \beta_i^{a'} \cdot \Omega_k^2 \cdot \beta_i^a \tag{A.12}$$

This pattern captures situations in which the relationship between an acquisition variable and future demand depends on other acquisition-related parameters, or when

such a relationship is quadratic. For example, it is possible that a strong preference for discounted products at the acquisition moment relates to price sensitivity in future demand *only* if the customer was purchasing for herself/himself, or outside the holiday period. In that case, the relationship between demand parameters and acquisition variables will be best represented by an interaction term.

Table A4: Simulated values for ω_k^1 and Ω_k^2 in the Quadratic/Interaction scenario

Variable	Demand variables		
	Intercept	Covariate 1	Covariate 2
ω_{k1}^1	0.30	-0.69	-0.05
ω_{k2}^1	0.86	-0.61	-1.04
ω_{k5}^1	1.16	-0.06	0.36
ω_{k3}^1	-1.44	-0.35	-0.27
ω_{k4}^1	-0.05	-0.10	0.10
ω_{k6}^1	-0.12	0.10	-1.11
Ω_{k11}^2	-0.01	0.06	0.00
Ω_{k22}^2	0.41	0.34	0.00
Ω_{k33}^2	-0.01	0.05	0.00
Ω_{k44}^2	0.01	-0.04	0.00
Ω_{k55}^2	0.17	-0.24	0.00
Ω_{k66}^2	-0.21	-0.11	0.00
Ω_{k12}^2	-0.36	-0.27	0.00
Ω_{k13}^2	-0.01	0.12	0.00
Ω_{k14}^2	-0.05	-0.01	0.00
Ω_{k15}^2	0.11	-0.08	0.00
Ω_{k16}^2	0.08	-0.16	0.00
Ω_{k23}^2	-0.01	-0.18	0.00
Ω_{k24}^2	0.24	0.10	0.00
Ω_{k25}^2	-0.24	-0.29	0.00
Ω_{k26}^2	-0.06	0.04	0.00
Ω_{k34}^2	0.17	0.07	0.00
Ω_{k35}^2	0.14	-0.14	0.00
Ω_{k36}^2	0.36	-0.10	0.00
Ω_{k45}^2	0.08	0.04	0.00
Ω_{k46}^2	-0.17	-0.15	0.00
Ω_{k56}^2	0.29	-0.17	0.00

- Scenario 3: Positive part

$$g_k(\beta_i^a | \Omega_k) = \omega_k^{1'} \cdot \begin{pmatrix} \max\{\beta_{i1}^a, 0\} \\ \vdots \\ \max\{\beta_{iP}^a, 0\} \end{pmatrix} \quad (\text{A.13})$$

This pattern captures situations in which an acquisition variable relates to future demand parameters, but only if the former passes a certain threshold. For example, the number of items purchased at the moment of acquisition might relate to the likelihood of purchasing again in the category, but only above a certain threshold that reflects strong parameters for such a category.

Table A5: Simulated values for ω_k^1 in the Positive part scenario

Variable	Demand variables		
	Intercept	Covariate 1	Covariate 2
ω_{k1}^1	0.34	0.00	0.30
ω_{k2}^1	0.00	0.00	0.86
ω_{k3}^1	0.00	0.00	-1.44
ω_{k4}^1	0.00	0.28	-0.05
ω_{k5}^1	0.00	0.00	1.16
ω_{k6}^1	0.00	0.00	-0.12

For each scenario, we generate the intercept (β_{i1}^y) and the effect of the first covariate (β_{i2}^y) according to the functions $g_1(\cdot)$ and $g_2(\cdot)$ as described in equations (A.11)–(A.13), while maintaining the effect of the second covariate (β_{i3}^y) to be a linear function of the acquisition variables. Furthermore, to compare parameters in the same scale across scenarios, we scale

demand parameters such that the standard deviation across individuals is equal across all scenarios.

Simulate individual-level behaviors

Once the individual-level parameters are generated, we simulate behaviors using the generated acquisition and demand parameters for each scenario, a set of market-level covariates $\mathbf{x}_{m(i)}^a$ for the acquisition model, and individual and time-variant covariates \mathbf{x}_{it}^y for the demand model. We assume a Gaussian distribution for all behaviors,

$$A_{ip} \sim N(\beta_{ip}^a + \mathbf{x}_{m(i)}^a \cdot \mathbf{b}_p^a, \sigma_p^a), \quad p = 1, \dots, 7 \quad (\text{A.14})$$

$$y_{it} \sim N(\mathbf{x}_{it}^{y'} \cdot \boldsymbol{\beta}_i^y, \sigma^y), \quad t = 1, \dots, 20. \quad (\text{A.15})$$

with $\sigma^a = 0.5$, $\mathbf{x}_{m(i)}^a \sim \mathcal{N}(0, 1)$, $\mathbf{b}^a \sim \mathcal{N}(0, 2)$, $\sigma^y = 0.5$, and $\mathbf{x}_{it}^y \sim \text{Bernoulli}(0.5)$.

A.6.3 Estimated models

Given the observed behaviors (A_{ip} and y_{it}) and the covariates ($\mathbf{x}_{m(i)}^a$ and \mathbf{x}_{it}^y), we estimate the model parameters. In addition to our proposed FIM, we use four benchmark models to infer β_j^y : (1) a hierarchical Bayesian demand-only model in which acquisition variables are not incorporated, (2) a linear model, where individual demand parameters are a linear function of the acquisition characteristics, (3) a full hierarchical model, where individual demand and acquisition parameters are jointly distributed according to a multivariate Gaussian distribution with a flexible covariance matrix, and (4) a Bayesian PCA model, identical to our proposed model, without the higher layer. For all models we assume the same linear demand model as in the data generation process, equation (A.15). We describe these models in more detail.

Hierarchical Bayesian (HB) demand-only model

This first benchmark is a *HB demand-only* model that does not incorporate acquisition variables. That is,

$$\beta_i^y | \mu^y, \Sigma^y \sim \mathcal{N}(\mu^y, \Sigma^y),$$

where μ^y , and Σ^y are the population mean vector and covariance matrix respectively.

We acknowledge that such a model would fail to provide individual-level demand parameter estimates for customers that are not in the calibration sample. In other words, the best this model can provide is to draw the estimates from the population distribution. We include this benchmark to illustrate the problem of estimating parameters when only one observation per customer is observed and most importantly, to have a reference of how much error we should obtain if the model only captured random noise.

Linear HB model

The second benchmark is the *linear HB model*, which is an extension of the previous model with the mean demand parameters being a linear function of the acquisition characteristics and market level covariates. That is,

$$\beta_i^y = \mu^y + \Gamma \cdot A_i + \Delta \cdot \mathbf{x}_{m(i)}^a + \mathbf{u}_i^y, \quad \mathbf{u}_i^y \sim \mathcal{N}(0, \Sigma^y),$$

where Γ capture the linear explanatory power of acquisition characteristics A_i , and Δ allows to control for market-level covariates $\mathbf{x}_{m(i)}^a$.

In this model, we incorporate both acquisition characteristics as well as market-level covariates to control for firm's actions that may be correlated with acquisition characteristics

(e.g. average price paid and promotional activity). Note that this model resembles the first simulated scenario in which the relationship between acquisition and demand parameters was assumed to be linear. As such, this model should be able to predict demand parameters in the first scenario most accurately.

Full hierarchical model

For the third benchmark, we endogenize the acquisition characteristics by modeling them as an outcome. Similar to our proposed FIM (described in Section 1.4.1), the full hierarchical model estimates acquisition and demand parameters jointly, with the difference that these two sets of parameters are modeled using a standard hierarchical model, rather than connected via DEF models. That is, the full hierarchical model assumes that

$$\beta_i = \begin{pmatrix} \beta_i^y \\ \beta_i^a \end{pmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma),$$

where $\boldsymbol{\mu}$ is the population mean vector of all individual parameters (demand and acquisition), and Σ is the population covariance matrix of these parameters, capturing correlations within demand and acquisition parameters as well as across those types of parameters.

Because of the Gaussian specification for β_i , this model imposes a linear relationship between β_i^y and β_i^a ; this is, the conditional expectation of β_i^y given β_i^a , is linear in β_i^a . As such, this model is mathematically equivalent to the linear HB model. However, the full hierarchical model differs from the linear model if acquisition behavior A_i is not linear in β_i^a (e.g. logit or log-normal. Moreover, if the number of acquisition characteristics increases, the full hierarchical model becomes more difficult to estimate due to the dimensionality of the

covariance matrix. In this simulation exercise we assume a linear (Gaussian) acquisition model and therefore the linear and full hierarchical models should provide equivalent results. Nevertheless, this is not the case in the empirical application as we incorporate binary acquisition characteristics modeled using a logit specification.

Bayesian PCA

The fourth benchmark is the closest to our proposed model, with the omission of the higher layer of traits (\mathbf{z}_i^2). Analogously as in our model, we model individual demand and acquisition parameters as a linear function of a set of traits,

$$\beta_i^y = \boldsymbol{\mu}^y + \mathbf{W}^y \cdot \mathbf{z}_i^1 \quad (\text{A.16})$$

$$\beta_i^a = \boldsymbol{\mu}^a + \mathbf{W}^a \cdot \mathbf{z}_i^1. \quad (\text{A.17})$$

In this Bayesian PCA model, we model the first layer \mathbf{z}_i^1 as a vector of independent standard Gaussian variables,

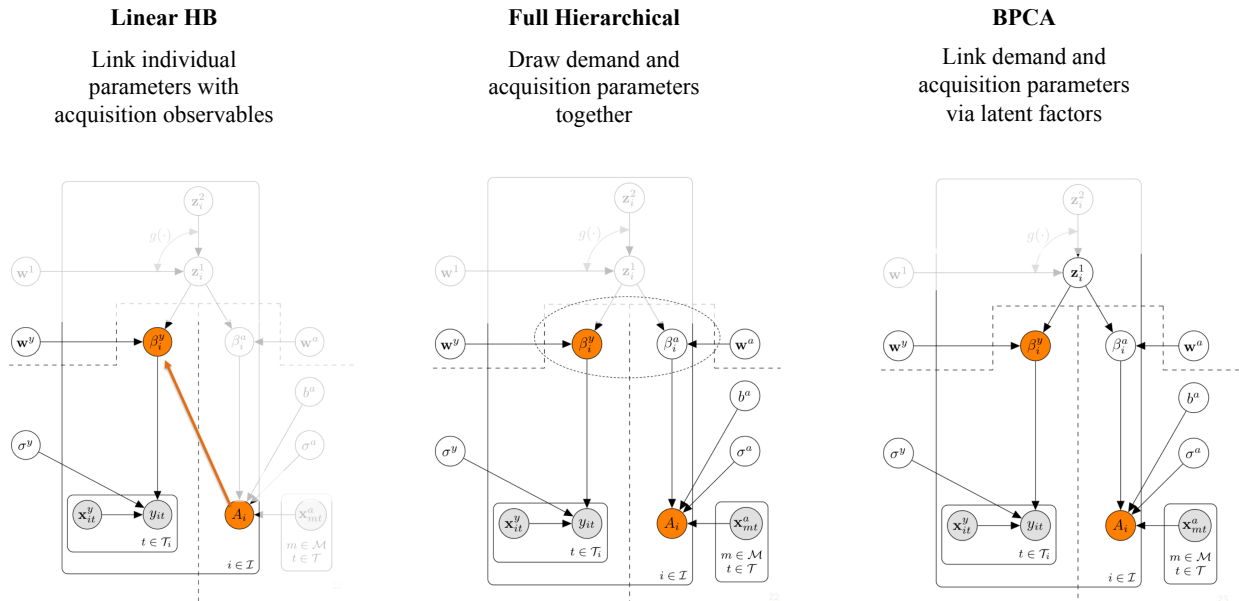
$$\mathbf{z}_{ik}^1 \sim \mathcal{N}(0, 1).$$

Note that like the linear HB and full hierarchical specifications, the PCA also imposes a linear relationship between β_i^y and β_i^a . However this approach is different from those because it allows for data dimensionality reduction via the latent factors. Similarly, as in our proposed model, we use sparse Gaussian priors on \mathbf{W}^y and \mathbf{W}^a , using an automatic relevance determination model to automatically select the number of traits.

As discussed in Section 1.4.1.5, the Bayesian PCA model is a nested specification of the proposed FIM (in which the second layer does not exist) whereas the full hierarchical

model and HB-linear specifications reflect alternative (simpler) ways in which past research has modeled these types of data. Figure A3 visually shows how each of these approaches compares with our proposed modeling framework.

Figure A3: Visualization of the benchmark models

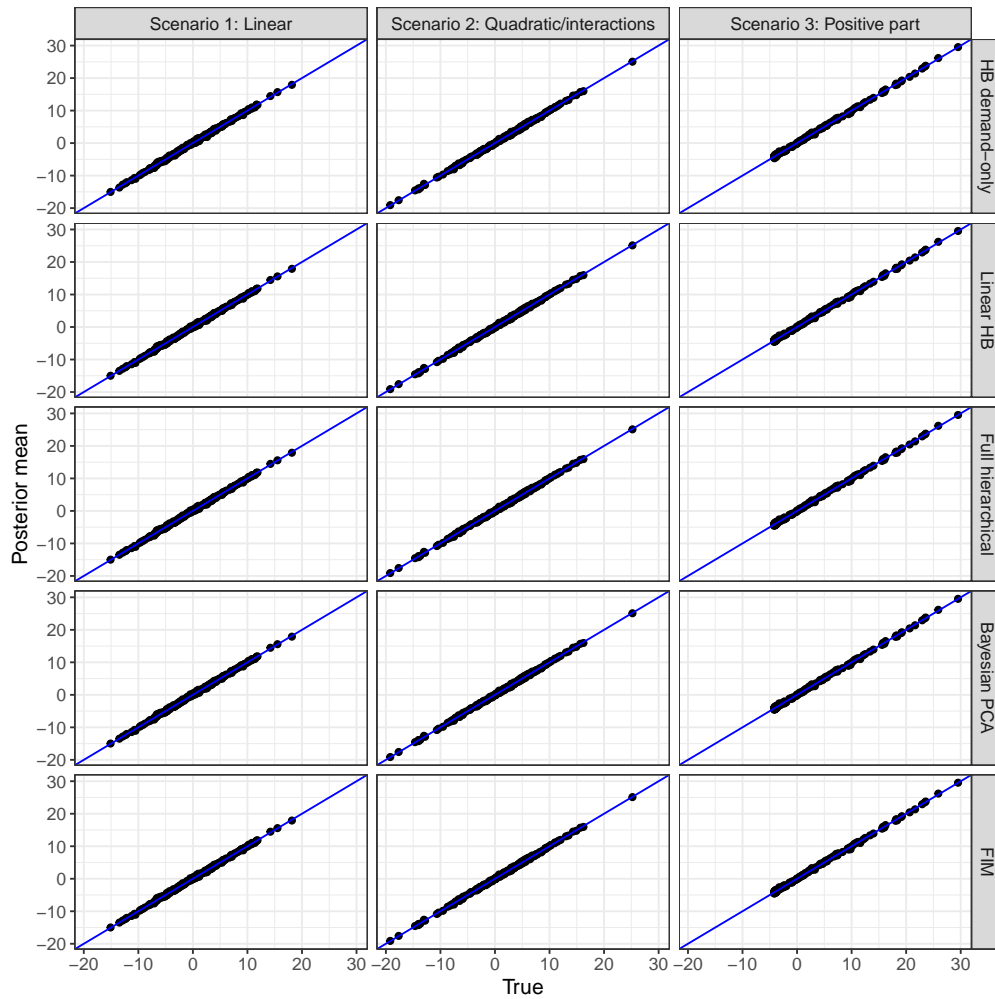


A.6.4 Assessing model performance

We calibrate each model using acquisition and demand data for 2,000 customers. This step resembles the firm calibrating each of the models (our proposed model as well as the benchmark models) with the historical data. First, we corroborate that all models are equally capable of recovering the individual-level parameters for customers in the calibration sample. In particular, we confirm that the in-sample predictions for β_i^y are almost perfect for all model specifications and for all scenarios (see Figure A4 for the in-sample predictions). In

other words, all models are equally capable of accurately estimating individual-level demand parameters for in-sample customers.

Figure A4: In sample individual posterior mean vs. true intercepts of the demand model. Each dot represents a customer from the calibration set. In blue, the 45 degree line represents perfect predictive power.



Then, we evaluate the ability of each model to form first impressions of newly-acquired customers. Under each scenario, we use the estimates of each model to predict the individual-level demand parameters for the remaining 200 customers, using only their acquisition data, and compare those predictions with the true values. As explained in the previous section, this task requires the computation of the individual posterior mean for

each individual ($\hat{\beta}_j^y = E(\beta_j^y|A_j, \mathcal{D})$) by integrating over the estimated density $p(\beta_j^y|A_j, \mathcal{D})$,

$$\hat{\beta}_j^y = \int \beta_j^y \cdot p(\beta_j^y|A_j, \mathcal{D})d\beta_j^y.$$

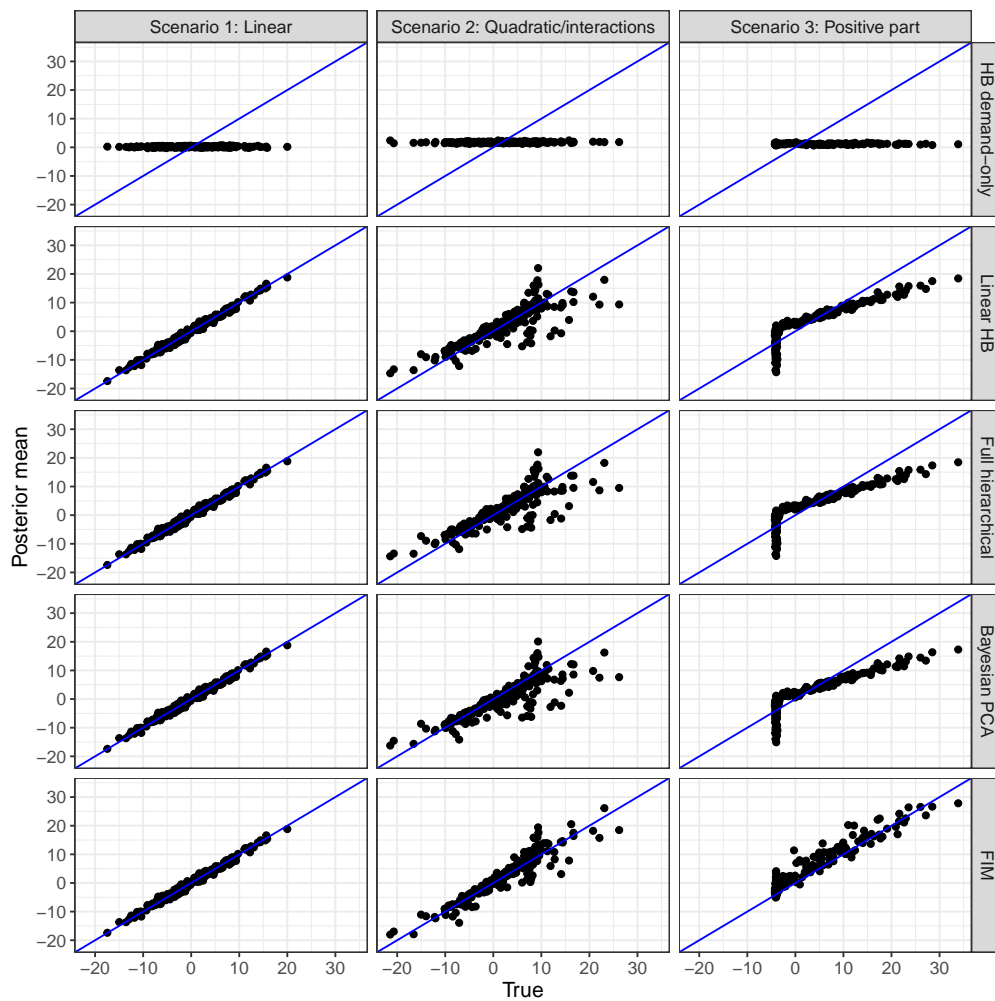
While the procedure described in Section 1.4.3 is valid for all models, the expectation $E(\beta_j^y|A_j, \mathcal{D})$ can be computed directly for some of the benchmark models, which we do for simplicity. For example, for the HB demand-only model, this procedure reduces to compute the expectation of individual draws of β_j^y from the population mean, which converges to the posterior mean of the population mean μ^y . For the linear HB model, it reduces to use the linear formulation and the posterior mean estimates of μ^y , Γ , and Δ . For the full hierarchical model, the Bayesian PCA model, and our proposed FIM, where acquisition is modeled as an outcome, we compute the posterior of β_j^y given A_j using HMC as described in Section 1.4.3.

Figure A5 shows the scatter plot of the predicted ($\hat{\beta}_{j1}^y$) versus actual (β_{j1}^y) individual demand intercepts from each model, for each scenario.³ Not surprisingly, the HB demand-only model that does not incorporate acquisition behavior in the model (top row of Figure A5) cannot distinguish (hold out) individuals from their population mean. Turning our attention to the other model specifications, we start analyzing the scenario in which the relationship between acquisition and demand parameters is linear (left-most column of Figure A5). Under this scenario, all models are equally capable of predicting demand estimates for (hold out) customers using only their acquisition data. This result is not surprising for the benchmark models as their mathematical specification resembles that of the simulated data. However, when the relationship between the acquisition and demand

³For brevity's sake, we present the results for one parameter of the demand model (the intercept), but the results hold for all other parameters as well.

parameters is not perfectly linear (as it is the case in scenarios 2 and 3), all benchmark models struggle to predict these individual-level estimates accurately. On the contrary, the proposed FIM is flexible enough to recover these parameters rather accurately. Note that the flexibility of the FIM comes at no overfitting cost; that is, even when the relationship is a simple linear relationship, our model recovers the parameters as well as the benchmark models, which assume a linear relationship by construction.

Figure A5: Out of sample individual posterior mean vs. true intercepts of the demand model. Each dot represents a customer from the hold out set; i.e., only their acquisition characteristics are used to form first impressions about their individual-level parameters. In blue, the 45 degree line represents perfect predictive power.



To explore the differences in accuracy more systematically, we compute two different measures of fit: (1) the (squared) correlation between true β_j^y and predicted $\hat{\beta}_j^y$ (i.e., R-squared)—measuring the model’s accuracy in sorting customers (e.g., differentiating customers with high vs. low value, more vs. less sensitivity to marketing actions)—and the root mean square error (RMSE)—measuring the accuracy on predicting the value/magnitude of the parameter itself.

The results are presented in Table 1.1 of the main manuscript, confirming the results from Figure A5. Under a true linear relationship (Scenario 1), the FIM predicts the individual parameters as good as the benchmark models. The RMSE of the FIM is comparable to the benchmark models, and the R-squared is equal to the benchmark models. However, when the relationship among the model parameters is not perfectly linear (Scenarios 2 and 3), the FIM significantly outperforms the benchmark models in all dimensions. In particular, the R-squared of the FIM is higher than that of the benchmarks, demonstrating that the model is superior at sorting customers based on their demand parameters. Moreover, the RMSE for the FIM is substantially lower than that of the benchmarks, indicating that the proposed model predicts the exact magnitude of customer parameters (e.g., purchase probability, sensitivity to marketing actions) more accurately than any of the benchmarks.

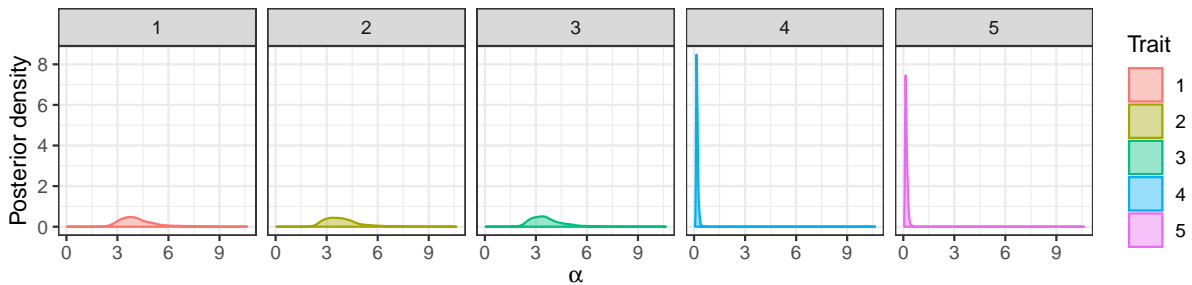
A.6.5 Interpreting the model parameters and results

To get a better sense of what the model is doing and what its parameters capture, we explore in detail the model estimates and compare those with the parameters used to simulate the data. We do so for the linear case, as it is the easiest to interpret the

relationships among variables. For this particular exercise, we select the FIM with 5 dimensions in the lower layer and 3 in the top layer.⁴ We start by evaluating the number of traits captured by the FIM; this is an insight that can be obtained in two ways. First, looking at the posterior estimates for α , parameters that determined the weights of the lower layer to check how many dimensions of the lower layer are activated in the model. Second, by looking at the specific weights, \mathbf{W}^y and \mathbf{W}^a , between the lower layer and the model parameters and interpret their meaning based on their magnitude.

We know from the simulations (Section A.6.1) that the data was generated from three factors: two factors generating 6 acquisition characteristics that relate to demand parameters, and another independent factor that generated one acquisition variable that was irrelevant for the demand model. Figure A6 shows the posterior distribution for α . While the model was specified to have 5 dimensions in the lower layer, it is obvious that the model only “needs” three, one of which is irrelevant in the demand specification.

Figure A6: Posterior distribution of α



We show in Table A6 the posterior mean of the rotated weight traits on demand parameters and acquisition parameters. The first two traits capture most of the variance across individuals for demand and acquisition parameters, while the other traits capture residual variance. First, trait 1 captures the correlation among acquisition variables 1

⁴Results are equivalent for other specifications of the model.

through 3, whereas trait 2 captures the correlation of acquisition variables 4 through 6. Second, both traits capture relationships with demand: trait 1 is negatively correlated with intercept and positively correlated with both covariates, whereas trait 2 is negatively correlated with intercept and covariate 2 (effect on covariate 1 is not significantly different from zero).

Table A6: Posterior mean of lower layer weights (\mathbf{W}^y and \mathbf{W}^a) for FIM.

Variable	Trait 1	Trait 2	Trait 3	Trait 4	Trait 5
Intercept	-5.55	-2.14	0.04	0.00	-0.00
Covariate 1	2.28	-0.53	0.10	-0.00	0.00
Covariate 2	2.91	-3.63	-0.04	-0.00	0.00
Acq. variable 1	-2.78	0.07	-0.04	-0.01	0.00
Acq. variable 2	-1.84	-0.03	0.02	0.00	0.00
Acq. variable 3	2.30	0.05	-0.02	0.00	-0.00
Acq. variable 4	-0.31	3.40	0.02	-0.01	0.01
Acq. variable 5	0.18	-1.95	0.00	0.01	0.02
Acq. variable 6	0.26	-2.91	-0.05	0.01	0.01
Acq. variable 7	-0.01	0.02	-0.03	-0.02	0.01

Note: In bold parameters such that corresponding CPI do not contain zero

Now, we are interested in comparing these insights with the true values used for the simulation, specifically how these estimated traits relate to the true factors in the data generation process. In the data generation process, demand parameters are generated from acquisition parameters. Instead, the FIM gives us the overall correlation of the traits with demand parameters, and not the one-to-one relationships between acquisition variables and

demand parameters. Therefore, in order to assess whether our model can capture the essence of the insights the “true” effect of factors f_{i1} and f_{i2} on acquisition parameters and demand parameters in Table A7. For the acquisition parameters, these true effects are B_{1p} and B_{2p} from (A.9) (whose values are shown in Table A2). For the demand parameters, these effects can be obtained by replacing (A.9) in (A.10), which reduces to $\omega_k^{1'}B_1$ and $\omega_k^{1'}B_2$ for the effects of factors 1 and 2, respectively.

Table A7: True associated effects of factors on demand and acquisition variables.

Demand/acquisition parameter	Variable	Factors	
		1	2
Intercept	$\omega_1^{1'}B_f$	6.20	-2.10
Covariate 1	$\omega_2^{1'}B_f$	-2.40	-0.57
Covariate 2	$\omega_3^{1'}B_f$	-2.77	-3.76
Acq. variable 1	B_{f1}	3.00	0.10
Acq. variable 2	B_{f2}	2.00	0.00
Acq. variable 3	B_{f3}	-2.50	0.00
Acq. variable 4	B_{f4}	0.00	3.50
Acq. variable 5	B_{f5}	0.00	-2.00
Acq. variable 6	B_{f6}	0.00	-3.00
Acq. variable 7	B_{f7}	0.00	0.00

By comparing Tables A6 and A7 we observe that: (1) trait 1 captures the reverse of factor 1 ($\hat{z}_{i1}^1 \approx -f_{i1}$); and (2) trait 2 captures factor 2 ($\hat{z}_{i2}^1 \approx f_{i2}$). This result implies that our model is able to capture and deliver meaningful insights that relate to the true data generation process.

A.6.6 Why is the model giving superior performance?

A natural question to ask is, why is the proposed model outperforming the benchmark models? As described in Section 1.4.1, the DEF component of the proposed model is very flexible at capturing underlying relationships between the model parameters. Such a property enables the model to capture non-linear relationships between acquisition characteristics and the parameters that drive customer demand. This is unlike the benchmarks whose specification imposes a linear relationship among the variables. As such, even though the in-sample predictions of all the models are very accurate (Figure A4), when any of the benchmark models are used to make (out-of-sample) predictions for newly-acquired customers, the predicted values differ dramatically from the actual values (Figure A5).

To better corroborate that it is the DEF component that brings the non-linearities, we compare in greater detail the predictions of the BPCA model with those of the FIM. We pick the BPCA (among the other benchmarks) because that is the only model that is mathematically nested to our proposed model. In turn, the BPCA is the closest to the FIM, with the difference that it does not have an upper layer (and its corresponding non-linear link function). Table A8 shows the squared correlation (true vs. predicted) for Covariate 1 of the second scenario (Quadratic/Interaction), for the BPCA and the FIM models, as we vary the number of dimensions. The first column corresponds to the fit of the BPCA model, as we increase the number of dimensions. We see an improvement in fit as we increase the number of dimensions from 1 to 2, and to 3; and no improvement after that, with the best fit obtained being around 0.25. However, the jump in fit is tremendous when we allow the

Table A8: Squared correlation (true vs predicted) for Covariate 1; Quadratic/Interaction Scenario.

Dim.	Lower layer	Dim. Upper layer		
		Bayesian PCA	FIM	
		0	1	2
	1	0.209	0.207	0.209
	2	0.237	0.304	0.306
	3	0.257	0.402	0.404
	4	0.250	0.539	0.425
	5	0.252	0.538	0.641
	6	0.250	0.509	0.612
	7	0.250	0.451	0.627
	8	0.243	0.525	0.571

model to have an upper layer (even if it only includes 1 dimension).⁵ Such an upper layer is the model component that allows for flexible relationships relationships. The same results hold when looking at the third scenario (Positive-part).

To conclude, the upper layer of the DEF — the component that allows the model to capture non-linear relationships among variables — is responsible for the great improvement in the model’s ability to predict (out-of-sample) individual-level parameters when the underlying relationship between acquisition characteristic and the demand parameter is not linear.

A.6.7 Exploring the number of dimensions per layer

As described in Section 1.4.1.3, we take a hybrid approach to model selection in which we make sure that the number of pre-specified dimensions is large enough — a phenomenon that can be validated from the model parameters — while we rely on the priors of the model to ensure regularization. In this appendix we leverage the simulation results to provide further

⁵We discuss the importance of the dimensionality of the upper layer in Section A.6.7.

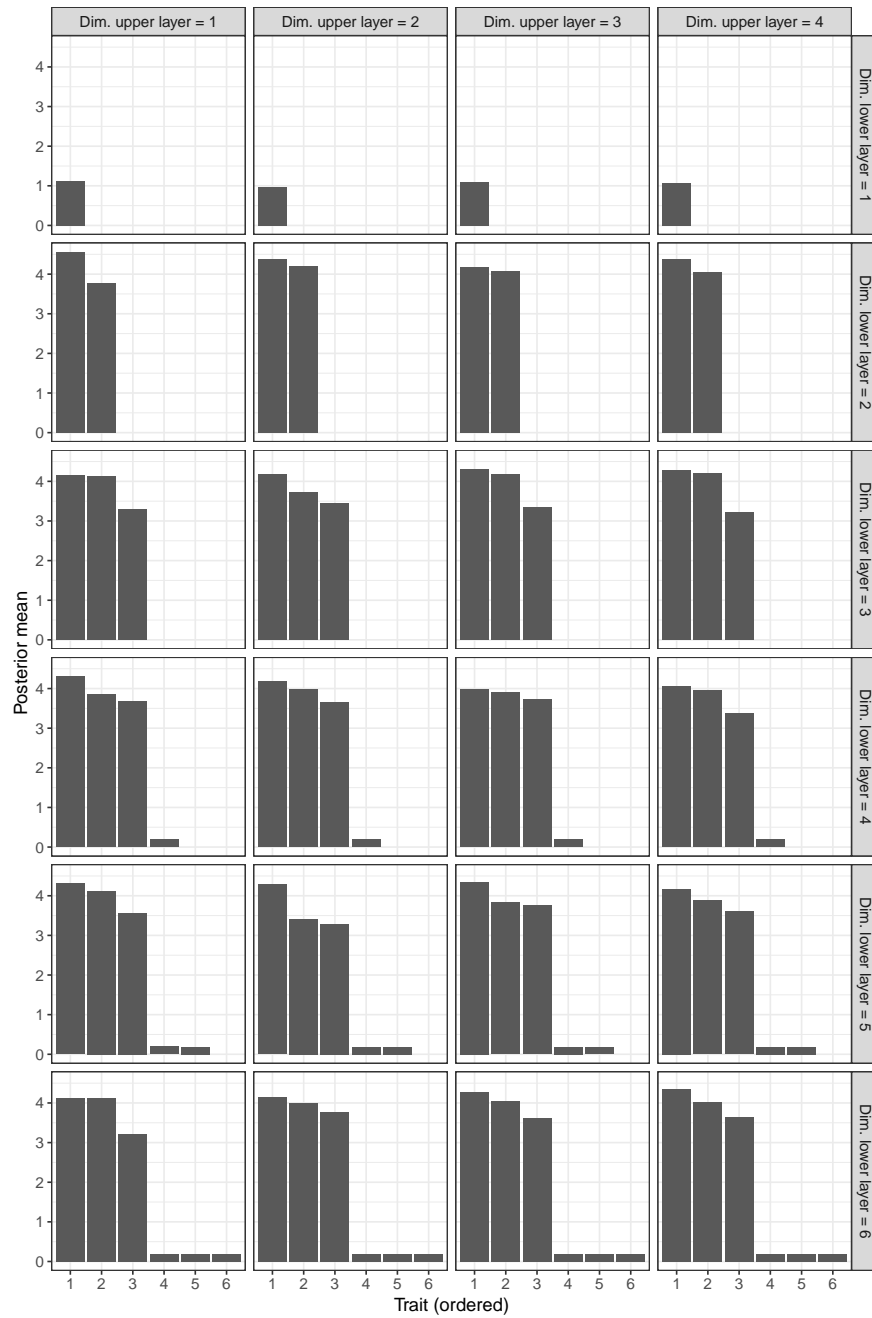
details about the model selection procedure and to corroborate the two premises that drive our model selection approach. Specifically, we present empirical evidence that (a) one can ensure that the model has a “large enough” number of dimensions by examining the posteriors of the Gaussian ARD priors, and (b) as long as the layers have enough dimensions to capture meaningful correlations and priors induce sparsity on the weight traits, increasing the number of dimensions on each layer would only lead to higher computational cost, without the corresponding loss in out of sample performance.

To illustrate how one can use the posterior of the Gaussian ARD priors to ensure that the number of dimensions is “large enough,” we revisit the model examined in Section A.6.5 in which the simulated behavior was generated by three factors, one of which had no impact on the demand parameters, and the FIM specification included 5 traits in the lower layer (e.g., $N_1 = 5$). As seen in that section, the FIM results not only recover that data generation process (Tables A6 and A6), but also informs of the number of dimensions in the lower layer (Figure A6). In this appendix we expand the results presented in Figure A6 by showing the posterior estimates for α for FIM specifications with different values for N_1 and N_2 (Figure A7).

As it is evident from the figure, the model detects that the data was generated from three latent traits (as long as the FIM is specified with $N_1 \geq 3$) and in cases where the FIM allows for larger dimensionality, the model “shuts down” the rest of the traits. In other words, regardless of the dimensionality of the top layer (N_2), when the number of traits in the lower layer is not enough, the model does not “shut down” any component. However, once N_1 is large enough (in this case $N_1 = 3$, as it was used to generate the data), the posterior mean of α_4 , α_5 and so on, are all close to zero. These results corroborate that the

posterior distribution of the Gaussian ARD variances can be used to show when the model has a “large enough” number of dimensions.

Figure A7: Posterior mean of α as a function of number of dimensions in lower layer (N_1) and upper layer (N_2). Components are sorted in decreasing order per model.

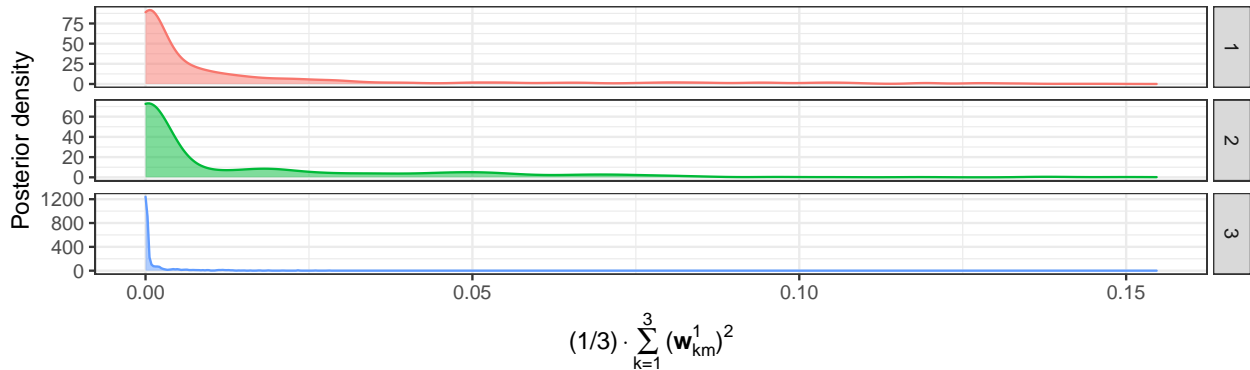


In contrast to the usefulness of α to detect relevant lower traits, our model does not have an analogous parameter to explore how many upper level traits are enough to capture the relevant non-linear correlations. Instead, each component of the upper weight \mathbf{W}^1 , \mathbf{w}_{km}^1 (for lower trait k and upper trait m), has i.i.d. sparse gamma priors, which by themselves induce regularization. In order to summarize each upper level trait in a way that can help us determine whether they make an impact on those 6 relevant lower layer traits, we compute a pseudo- α_m^1 for each upper trait m using the weight matrix \mathbf{W}^1 . Similarly to how the lower level weights \mathbf{W}^y and \mathbf{W}^a are related to α (i.e., variance of zero-centered Gaussian distributions), we compute these pseudo- α^1 's by averaging the square of all weights associated with a fixed upper level trait and those 6 relevant lower level traits, as described by

$$\text{pseudo-}\alpha_m^1 = \frac{1}{6} \sum_{k=1}^6 (\mathbf{w}_{km}^1)^2.$$

We show the posterior this quantity in Figure A8. Not surprisingly, these posterior distributions are concentrated close to the origin, which suggests that no upper trait is relevant for this scenario (as the data was generated linearly).

Figure A8: Posterior distribution of pseudo- α^1 (Linear scenario).

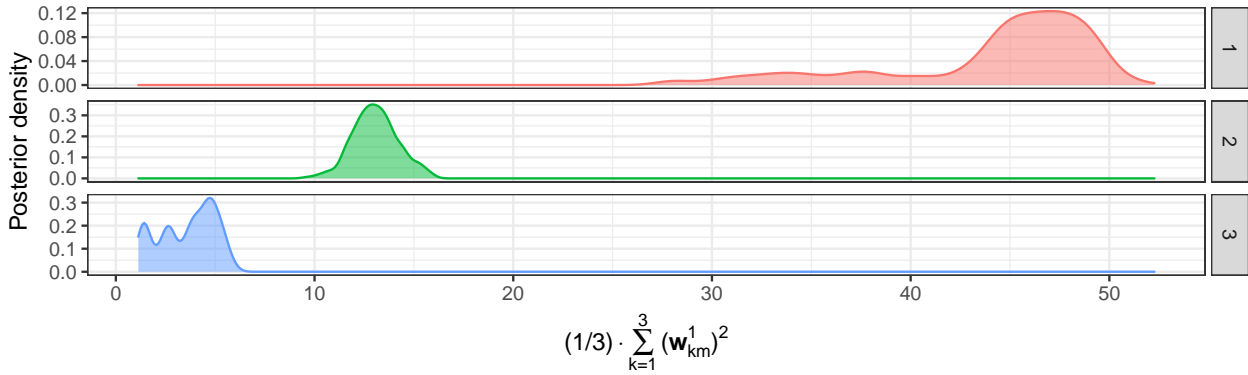


More interestingly, we further explore this quantity using a scenario in which the model requires to capture non-linear relationships, such as the one with Interactions. Figure A9 shows the posterior of pseudo- α^1 for two FIM specifications with different values of N_2 . First, Figure A9a clearly shows that the FIM with $N_1 = 5$ and N_2 estimated for the Interactions scenario, unlike the FIM estimated using the linearly simulated data, has all three upper traits being relevant in the model. Second, if we estimate a FIM with more upper traits ($N_1 = 5, N_2 = 5$) the model starts to “shut down” the less relevant traits, indicating that such a model is enough to recover the non-linear relationships present in those data.

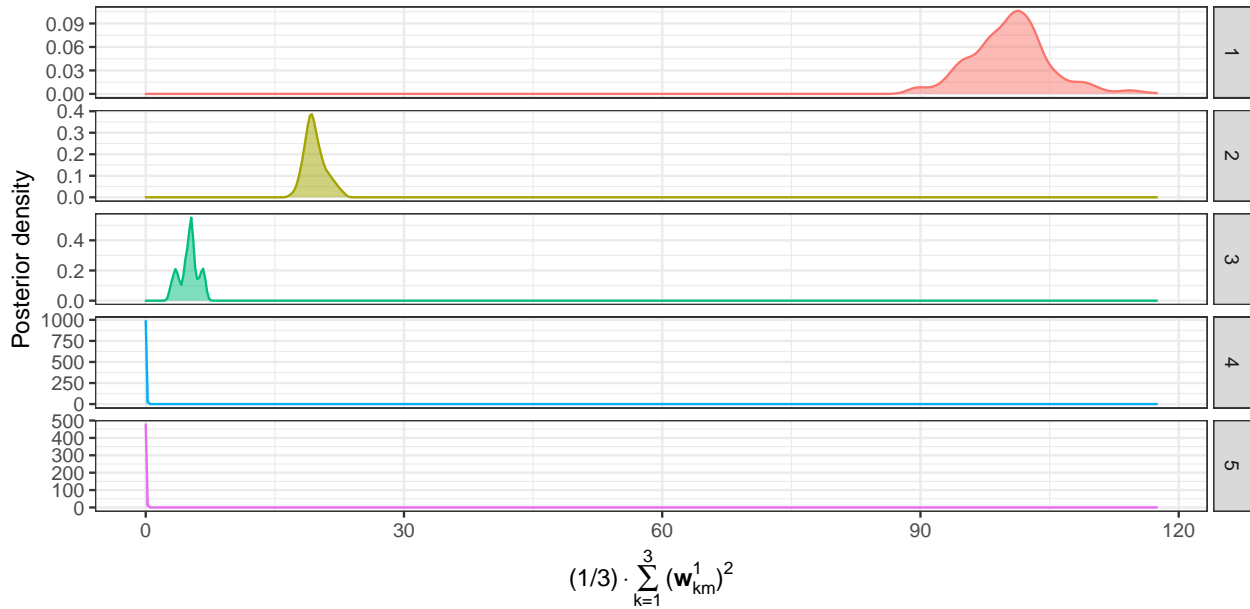
Finally, we leverage the results of multiple estimated FIM specifications over the Interactions scenario and show that once the FIM specification contains the dimensions “needed” by the data, the performance of the model remains the same even if we add dimensions to the DEF component. To illustrate this phenomenon, we focus on the performance of the FIM at predicting the parameter for the sensitivity to the first covariate (bottom half of middle columns in Table 1.1). Figure A10 shows the squared correlation between simulated and predicted values of the parameter of interest (higher numbers imply better model performance). The figure shows a notable improvement in performance as we increase the dimensionality of the lower layer from 1 to 2, 3, and 4. However, once $N_1 > 3$, the model performs very similarly as more layers are added to DEF. Similarly, we observe a radical increase in performance as one increases the dimensionality of the upper layer (from 0 to 1, 2 and 3); reaching a point in which more dimensions do not alter the performance of the model. In other words, the performance in out-of-sample recovery of demand parameters flattens, once the model has a “large enough” number of dimensions.

Figure A9: Posterior distribution of pseudo- α^1 (Interactions scenario).

(a) FIM ($N_1 = 5, N_2 = 3$).



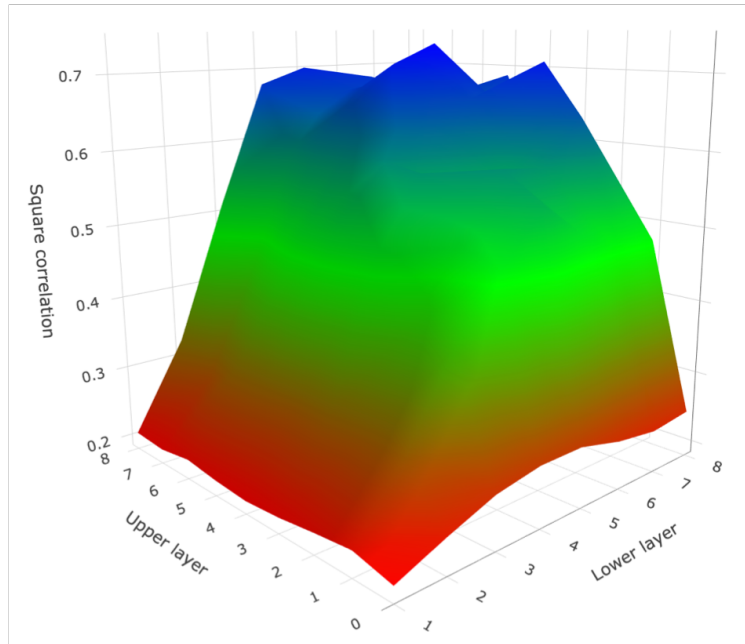
(b) FIM ($N_1 = 5, N_2 = 5$).



A.6.8 Model performance “at scale”

While the analysis thus far assumed a handful number of acquisition variables, many firms collect a larger quantity of behaviors when a customer makes their first transaction. These firms do not necessarily know a priori which variables can be most predictive of demand parameters, and if so, what the underlying relationship between these variables would be. In this section we show that models that incorporate all interactions fail to recover demand

Figure A10: Square correlation between simulated and predicted β for Covariate 1 in Scenario 2: Interaction



parameters when the number of acquisition variables is large, whereas the FIM can accurately infer these non-linear relationships. We maintain a similar simulation structure, where acquisition parameters are driven by factors, but instead we now have 5 factors and 60 acquisition behaviors, where acquisition behavior is driven by one and only one factor, and each factor generates 12 acquisition parameters. We start by describing the simulation details and their differences to the main analysis in Section A.6.1. Then, we describe the additional estimated models, specifically those that include interactions. Finally, similarly as in Section A.6.4, we show the models' ability to infer demand parameters for out of sample customers.

Simulation details

We assume there are 3 demand parameters (intercept and two covariates) and 60 acquisition parameters, for 60 acquisition characteristics. We generate these acquisition parameters as being highly correlated among each other by assuming these parameters are driven by one of five factors f_{i1}, \dots, f_{i5} . Similarly as in Equation (A.9), we generate acquisition parameters by:

$$\begin{aligned}
 \beta_{ip}^a &\sim N(\mu_p^a + B_{1p} \cdot f_{i1}, \sigma_p^b), & p = 1, \dots, 12 \\
 \beta_{ip}^a &\sim N(\mu_p^a + B_{2p} \cdot f_{i2}, \sigma_p^b), & p = 13, \dots, 24 \\
 \beta_{ip}^a &\sim N(\mu_p^a + B_{3p} \cdot f_{i3}, \sigma_p^b), & p = 25, \dots, 36 \\
 \beta_{ip}^a &\sim N(\mu_p^a + B_{4p} \cdot f_{i4}, \sigma_p^b), & p = 37, \dots, 48 \\
 \beta_{ip}^a &\sim N(\mu_p^a + B_{5p} \cdot f_{i5}, \sigma_p^b), & p = 49, \dots, 60,
 \end{aligned} \tag{A.18}$$

where μ_p^a is the mean of the p^{th} acquisition parameter; $B_{\ell p}$ represent the impact of factor ℓ respectively on the p^{th} acquisition parameter; and σ_p the standard deviation of the uncorrelated variation of the p^{th} acquisition parameter.

The rest of the simulation design is identical as the simulation in Section 1.4.3, with a different set of parameters Ω . In order to incorporate noise and to allow for different acquisition parameters to inform demand parameters, we relate demand parameters only to a subset of acquisition parameters. Specifically, we choose Ω such that demand parameters are only affected by acquisition parameters from three out of the five factors. We achieve this by setting to zero Ω values for the remaining acquisition parameters. The intercept is a

function of the acquisition parameters from factors 1, 2 and 3 (i.e.,

$\Omega_{1p} = 0, \forall p = 37, \dots, 60$). Covariate 1 is a function of the acquisition parameters from

factors 1, 2 and 4 (i.e., $\Omega_{2p} = 0, \forall p = 25, \dots, 36, 49, \dots, 60$). Covariate 2 is a function of the

acquisition parameters from factors 2, 3 and 4 (i.e., $\Omega_{3p} = 0, \forall p = 1, \dots, 12, 49, \dots, 60$).

Similarly as in the main simulation analysis, Covariate 2 is always a linear function of

acquisition parameters for all scenarios. The values we use for Ω are specific to each scenario:

- **Linear:** Following (A.11), we define $\omega_{kp}^1 \sim \mathcal{N}(0, 2)$ for all non-zero ω_{kp}^1 .
- **Quadratic/Interaction:** Following (A.12), we define $\omega_{kp}^1 \sim \mathcal{N}(0, 2)$ for all non-zero ω_{kp}^1 ; and $\Omega_{kpp'}^2 \sim \mathcal{N}(0, 1)$ for all non-zero $\Omega_{kpp'}^2$.
- **Positive part:** To avoid attenuating the effect of the non-linear function by combining a large number of non-linear functions of correlated acquisition parameters, we fix the effect to the intercept and the first covariate as a function of only one acquisition parameter from each of the three factors that determine that demand parameter. Following (A.13), we define $\omega_{3p}^1 \sim \mathcal{N}(0, 2)$ for all non-zero ω_{3p}^1 , and:

$$\begin{array}{lll} \omega_{1,1}^1 = 12.5 & \omega_{1,13}^1 = -8 & \omega_{1,25}^1 = 4 \\ \omega_{2,1}^1 = -7.5 & \omega_{2,13}^1 = -4 & \omega_{2,37}^1 = 8. \end{array}$$

Finally, to compare parameters in the same scale across scenarios, we standardize demand parameters such that the population standard deviation is 2.

Estimated models

In addition to all models described in Section A.6.3, we estimate a Linear HB model where we include all interactions of acquisition parameters,

$$\beta_i^y = \mu^y + \Gamma \cdot \tilde{A}_i + \Delta \cdot \mathbf{x}_{m(i)}^a + \mathbf{u}_i^y, \quad \mathbf{u}_i^y \sim \mathcal{N}(0, \Sigma^y),$$

where \tilde{A}_i includes all acquisition characteristics, their squares, and all two-way interactions among them.

We also estimate a Lasso model with all interactions, which is identical to the Linear HB model with interactions, but we exchange the Gaussian prior for a Laplace prior to enforce regularization using a different functional form.

Results

We estimate all models except the full hierarchical model, which is computationally unstable given that now there are 60 acquisition variables, and therefore we need a 63×63 covariance matrix. Note that in theory, and in practice as we showed in Section A.6.4, the full hierarchical model is equivalent to a Linear HB model. Therefore, removing this model from the analysis does not bias our benchmark.

We show in Table A9 the out of sample prediction of intercept, and the two covariates under all three scenarios for all models. We replicate the main results from Section A.6.4. Both the Linear HB and Bayesian PCA models perform well in the Linear scenario. The FIM performs as good as these models in the Linear scenario, and outperforms these linear

models in the Quadratic/Interaction and the Positive part scenarios. More importantly, both models that include all interactions, Linear and Lasso, do not perform well in any scenario.

Table A9: Model at scale results

Model	Intercept		Covariate 1		Covariate 2	
	R-squared	RMSE	R-squared	RMSE	R-squared	RMSE
Linear						
HB demand-only	0.000	2.018	0.000	2.038	0.000	2.003
Linear HB	0.990	0.198	0.987	0.231	0.983	0.264
Linear with interactions	0.202	4.267	0.166	4.825	0.121	5.265
Lasso with interactions	0.161	5.916	0.115	6.129	0.108	5.561
Bayesian PCA	0.990	0.197	0.988	0.229	0.983	0.265
FIM	0.990	0.206	0.987	0.230	0.983	0.262
Quadratic/Interaction						
HB demand-only	0.004	2.060	0.000	2.133	0.007	2.084
Linear HB	0.231	1.808	0.398	1.663	0.994	0.167
Linear with interactions	0.147	4.064	0.201	4.331	0.246	4.125
Lasso with interactions	0.147	4.212	0.211	4.871	0.236	4.181
Bayesian PCA	0.243	1.790	0.408	1.646	0.994	0.167
FIM	0.598	1.456	0.681	1.432	0.994	0.165
Positive part						
HB demand-only	0.003	2.010	0.005	2.030	0.017	1.965
Linear HB	0.723	1.059	0.746	1.019	0.990	0.201
Linear with interactions	0.232	3.990	0.165	4.916	0.122	4.414
Lasso with interactions	0.161	4.493	0.088	5.336	0.186	5.032
Bayesian PCA	0.728	1.052	0.747	1.017	0.991	0.196
FIM	0.884	0.699	0.853	0.825	0.991	0.192

A.7 Empirical application: Additional results

A.7.1 Possible sources of endogeneity in the model components

Like most demand models including firm's marketing actions, we face the risk of introducing endogenous variables in our model, potentially preventing us from obtaining unbiased estimates of the customers' parameters. If that were the case, the relationships between acquisition characteristics and demand parameters captured by the model would likely reflect the firms strategies, and not the true underlying correlations that the FIM intends to capture.

Given the intended applications for this modeling framework, there are three mechanisms by which endogeneity concerns would arise: (unobserved) *temporal shifts* that systematically affect both the time-varying covariate and the overall demand, *static targeting rules*, whereby some customer characteristics (unobserved to the researcher) makes a customer more/less prone to receive marketing actions, while such a characteristic is also correlated to other components of the model, and *dynamic targeting rules*, whereby the presence/absence of the marketing action is driven by an unobserved customer state, which is also correlated with the individual propensity to transact with the firm. The former case is likely to be present if, for example, the firm introduced products or ran specific campaigns only when periods of lower/higher level of demands were expected. The second case corresponds to situations in which marketing actions such as e-mails are prioritized to customers of certain characteristics, for example, those who usually transact online, which is likely to be correlated with one of the acquisition characteristics. The third case is that in which the firm targets only customers who exhibit a behavior that is correlated with demand,

for example, send an email to customers who have visited the online store in the last week, or those who abandoned a basket before purchase, etc.

First, we explore the extent to which these phenomena might present in our application. According to the managers of the focal firm, marketing actions are decided in two steps. First, the firm chooses periods in which it will engage in promotional activity (i.e., run a marketing campaign). This decision is made from the headquarters, runs several times through the year (with special campaigns run during the holidays) and affects all markets simultaneously. Second, managers in each focal market choose the set of customers who will receive each campaign, with the proportion of customers not being determined consistently. The only variable that some markets include in their targeting rules is recency (i.e., time since last purchase). The introduction of new products follows a similar process—i.e., the decision being made globally, the implementation affected also by local factors such as distribution shocks in each of the markets—with the main difference being that the second step does not vary across customers of the same market.

Therefore, regarding potential (omitted) temporal shifts, the only variable that could systematically affect the presence/absence of promotional activity in all markets is the holiday season, which is not omitted as it is included in the model. Regarding (static) targeting rules, we confirm with the firm and verify with the data that these were not present in our application. Nonetheless, it is worth noting that when such targeting rules are present (e.g., the firm contacts customers based on demographic information), because the model includes unobserved heterogeneity on purchase frequency (first element of β_i^y), the identification of the individual-level sensitivity to promotional activity comes mainly from individual differences across periods, for which we have rich variation during the four years of

available data. Finally, regarding dynamic targeting rules, it is indeed the case that some customers (in the most sophisticated markets) are more/less prone to receive emails and DM based on their purchase activity. However, our model not only includes unobserved heterogeneity on purchase frequency — capturing the customers’ base level of activity — but also includes the recency of purchase, alleviating the endogeneity concerns arising from potential correlation between the firm’s targeting policies and customers’ propensity to transact in a particular period.

To conclude, given the business nature of our application, the rich variation in our data (Section 1.5.1.2), and our model specification, we argue that the potential endogenous nature of the marketing actions is not a main concern in this research. Nevertheless, in situations where these conditions do not hold (due to different strategic behavior by the firm or for data limitations), the demand model should be adjusted to account for the firm’s targeting decisions. Given the flexibility of our modeling framework, those adjustments would merely involve extending the demand model to capture unobserved shocks between firm’s actions and individual-level responsiveness (Manchanda et al., 2004) or adding correlations between firm decisions and unobserved demand shocks through copulas (Park and Gupta, 2012), depending on how these actions are determined by the firm. Those changes would only affect the demand (sub)model and not the overall specification of the FIM.

A.7.2 Latent attrition benchmarks models

We estimate three additional non-nested benchmark models (borrowed from the CRM literature) that do account for latent attrition: (1) Linear model with marketing actions + logistic attrition process (without acquisition covariates), (2) Linear model (without

marketing actions) + logistic attrition with acquisition covariates, and (3) Linear model with marketing actions + logistic attrition with acquisition covariates.

For all the aforementioned models we define purchase incidence (y_{it}) given attrition, which we denote as h_{it} , and we have that $p(y_{it} = 1|h_{it} = 1) = 0$, $p(h_{it} = 0|h_{it-1} = 1) = 0$, and

$$p(h_{it} = 1|h_{it-1} = 0) = \text{logit}^{-1} [\beta_i^h],$$

where β_i^h is a (scalar) parameter that captures the individual log-odds of attrition. In all specifications, we model the purchase incidence parameters β_i^y as a linear function of acquisition characteristics as described in Appendix A.6.3.

The models differ in the inclusion of marketing actions into the demand given attrition component and modeling of the attrition parameter β_i^h as displayed in Table A10.

Table A10: Latent attrition benchmarks models.

	Demand $p(y_{it} = 1 h_{it} = 0)$	Attrition parameter β_i^h
Latent Attrition		
w/ Acq.	$\text{logit}^{-1} [\beta_{i1}^y + \alpha_m]$	$\beta_i^h = \mu^h + \Gamma^h \cdot A_i + \Delta^h \cdot \mathbf{x}_{m(i)}^a + u_i^h$
w/ Mktg. Actions	$\text{logit}^{-1} [\mathbf{x}_{it}^{y'} \cdot \beta_i^y + \alpha_m]$	$\beta_i^h = \mu^h + u_i^h$
w/ Acq.+Mktg. Actions	$\text{logit}^{-1} [\mathbf{x}_{it}^{y'} \cdot \beta_i^y + \alpha_m]$	$\beta_i^h = \mu^h + \Gamma^h \cdot A_i + \Delta^h \cdot \mathbf{x}_{m(i)}^a + u_i^h$

Note that in all specifications we model jointly the unobserved individual components of purchase incidence and attrition parameters by $[\mathbf{u}_i^y, u_i^h] \sim \mathcal{N}(0, \Sigma^{yh})$.

A.7.3 Interpreting the latent traits

Finally, we further explore the posterior estimates of the (lower layer) hidden traits and their relationship with the demand and acquisition parameters to provide additional insights into customer traits and behaviors. We begin by analyzing which latent traits capture the most salient relationships in the data. We do so by exploring the posterior estimates of the

parameters governing the ARD component of the model and find that six traits carry most of the “weight” at connecting acquisition and demand parameters. (Please see Appendix A.7.6 for details.) Then, we investigate the correlations among these traits (Table A11), exploring whether customers that score high in a particular trait also score high (or low) in another trait. Note that these traits *do not capture segments* in the population (e.g., groups of customers of similar characteristics) but rather traits that capture the multiple dimensions of customer behavior. In other words, every customer has a score for each of the traits, being not only possible but very likely that customers score high in more than one trait. In our data, customers who score high in Trait 4 also tend to score high in Trait 6 (correlation= 0.553). On the contrary, those same customers have the tendency to score low in Trait 5 (correlation= -0.268).

Table A11: Posterior mean of correlations across customers of individual lower level traits \mathbf{z}_i^1 .

	Trait 1	Trait 2	Trait 3	Trait 4	Trait 5
Trait 1	1.000				
Trait 2	-0.144				
Trait 3	0.101	-0.113			
Trait 4	0.130	0.185	0.170		
Trait 5	-0.026	-0.141	-0.057	-0.268	
Trait 6	0.129	0.242	0.258	0.553	-0.361

An obvious question to ask is: What do these traits represent? To answer that question we compute the posterior mean of the weights of each of the rotated trait on each of the acquisition and demand parameters (Table A12). Looking at the weights to the demand parameters, we learn that the first trait is the most relevant in explaining heterogeneity in the base propensity to buy. Scoring high on this “high-frequency” trait also relates to a positive response to product introductions in future demand. This first trait is negatively correlated with whether the first purchase was made online and whether that purchase

contained a product in the Home category; but positively correlated with whether the customer purchased a product in the Hair Care category. Interestingly this trait is also positively correlated with first transaction baskets containing products that score high on dimension 4 of the Basket Nature product embeddings. Moreover, customers that score high on this trait are more likely to buy at their first purchase smaller sized products and travel sized products.

Table A12: Rotated traits weights' on acquisition and demand variables

Parameter	Trait					
	1	2	3	4	5	6
Demand (W^y)						
Intercept	0.133	0.129	-0.106	-0.072	-0.002	0.024
Email	-0.018	-0.016	0.046	0.027	-0.015	-0.004
DM	0.010	0.038	-0.003	-0.001	0.013	-0.004
Product introductions	0.044	0.085	0.001	-0.029	-0.026	0.009
Season	-0.025	0.058	0.027	0.085	0.004	0.005
Acquisition (W^a)						
Avg. price (log)	-0.109	0.022	-0.644	-0.370	0.039	0.313
Amount (log)	-0.021	0.076	-0.541	0.305	0.209	0.425
Quantity (log-log)	0.074	0.066	0.050	0.647	0.174	0.130
Package size (log)	-0.143	0.052	-0.087	-0.205	0.016	0.217
Holiday	0.029	-0.110	0.053	0.159	0.085	0.170
Discount	0.298	-0.073	0.280	0.414	0.133	0.029
Online	-0.382	1.368	0.581	6.830	0.019	0.146
New product	0.007	0.216	-0.283	0.544	0.354	0.234
Travel	0.470	-0.928	0.440	0.724	0.413	0.037
Category: Body Care	0.248	-4.922	-0.112	2.916	-0.072	-0.016
Category: Body Perfume	-0.025	0.436	-1.152	0.554	0.462	0.079
Category: Face Care	0.352	0.610	0.051	0.745	0.234	0.718
Category: Hair Care	1.267	1.178	-0.514	1.930	-0.631	-0.595
Category: Home	-1.097	-0.051	-0.336	1.836	1.073	-0.417
Category: Kits	0.285	0.227	-0.469	0.803	-0.100	0.225
Category: Make Up	0.377	0.528	0.334	1.149	-0.137	0.001
Category: Others	-0.134	0.230	0.623	1.845	0.387	0.029
Category: Services	-0.006	0.110	-0.501	5.762	-0.545	0.102
Category: Toiletries	0.239	0.733	0.200	1.190	0.607	-0.268
BasketNature dimension 1	-0.104	-0.022	-0.071	0.083	0.078	-0.112
BasketNature dimension 2	0.042	0.012	-0.011	-0.003	0.110	-0.035
BasketNature dimension 3	0.193	0.082	0.034	-0.040	-0.180	0.153
BasketNature dimension 4	0.200	0.105	-0.021	0.136	-0.167	0.005
BasketNature dimension 5	-0.035	0.003	0.001	0.025	0.009	0.154
BasketNature dimension 6	0.120	-0.017	0.141	-0.102	0.012	0.010
BasketDispersion dimension 1	-0.150	0.012	-0.166	0.256	0.237	-0.238
BasketDispersion dimension 2	-0.033	0.026	-0.105	0.196	0.114	-0.151
BasketDispersion dimension 3	-0.045	-0.094	-0.155	0.379	0.039	-0.120
BasketDispersion dimension 4	0.113	0.086	-0.216	0.406	-0.087	-0.082
BasketDispersion dimension 5	-0.137	0.123	-0.154	0.360	0.155	-0.195
BasketDispersion dimension 6	-0.033	-0.020	-0.159	0.462	0.078	-0.160

Another interesting trait is number four, which is associated with lower propensities to buy (intercept) and higher activity during the holiday season (Season variable). This “holiday-customer” trait is positively correlated with whether customers have been acquired online and during the Holiday season. This trait is positively associated with less expensive products and more units on the first transaction. With respect to the type of products associated with the first purchase, customers that score high on this trait are more likely to buy in the Body Care, Hair Care and Home categories. (Note that this trait is capturing some of the correlations among acquisition variables reported in Table 1.4—e.g., [Online-FaceCare]= 0.48—allowing the model to clean redundancies in the acquisition characteristics and tie the main trait to demand variables.) Finally, this “holiday-customer” trait is related with very diverse baskets (with respect to the type of products purchased in the first transaction), as indicated by its positive weights on Basket dispersion in all six dimensions.

A.7.4 FIM predictive accuracy using in-sample customers

Table A13 shows the performance of all models on the *Training* sample. The first two columns show the in-sample fit for each of the models, for which we compute log-likelihood and Watanabe-Akaike Information Criterion (WAIC) (Watanabe, 2010). Columns 3 through 6 show different measures of out-of-sample prediction accuracy, computed for customers in the training sample, but using the time periods that were not included in the estimation (i.e., periods after April 2014). We compute log-likelihood as well as the root mean square error (RMSE) for behavioral predictions. In particular, we compare the predicted and actual number of transactions at the observation level (i.e., at the customer/period level), at the

customer level, calculating the total number of transactions per customer (in “future” periods), and at the period level, computing the total number of transactions per period. While the HB benchmark model fit the in-sample data better than our proposed model, the FIM outperforms all benchmarks in the out-of-sample predictions. In other words, whereas the hierarchical models are very flexible at capturing heterogeneity in the training data, such a model is likely overfitting the data, as reflected in the out-of-sample predictions. On the other hand, the FIM forecasts the out-of-sample behavior of existing customers with greater accuracy.

Table A13: Model fit and prediction accuracy for the *Training* sample

Model	In-sample		Out-of-sample (future periods)			
	Log-Like	WAIC	Log-Like	RMSE		
				Observation	Customer	Period
HB - Linear	-7843.0	17807.8	-5511.1	0.202	0.723	62.841
Latent Attrition w/ Acq	-7880.1	17507.7	-6126.5	0.201	0.750	78.810
Latent Attrition w/ Mktg. Actions	-7781.1	17715.5	-5786.0	0.206	0.767	74.525
Latent Attrition w/ Acq+Mktg. Actions	-7612.8	17438.2	-6476.8	0.209	0.812	81.143
Bayesian PPCA	-8482.4	18361.4	-5137.2	0.191	0.573	35.696
Feed-Forward DNN	--	--	--	0.189	0.556	53.410
Random Forest	--	--	--	0.193	0.616	133.598
FIM ($N_1 = 13, N_2 = 5$)	-9135.4	18885.7	-5096.4	0.190	0.533	32.313
Other FIM specifications						
FIM ($N_1 = 12, N_2 = 2$)	-8654.0	18555.7	-5097.2	0.191	0.558	32.612
FIM ($N_1 = 12, N_2 = 5$)	-8952.1	18927.6	-5116.7	0.190	0.541	32.762
FIM ($N_1 = 13, N_2 = 2$)	-8587.6	18399.0	-5140.1	0.192	0.578	35.454
FIM ($N_1 = 14, N_2 = 2$)	-8683.6	18531.9	-5131.8	0.191	0.561	33.824
FIM ($N_1 = 14, N_2 = 5$)	-8613.9	18465.3	-5147.6	0.191	0.571	34.423

A.7.5 Population distribution and individual-level posterior distributions

Figure A11 summarizes the inferred individual posterior distributions of the demand parameters of *Test* customers using their acquisition characteristics. The top row of Figure A11 shows the degree of heterogeneity that the FIM infers. How uncertain are those

inferences at the individual level? In order to answer that question, for each demand parameter, we sort customers based on their posterior means, and compute their 95% CPI. The second row of Figure A11 shows the uncertainty at the individual level that the model can infer these parameters: each customer is represented horizontally, where the shaded area shows their 95% CPI and the white line, their posterior mean. Using this figure we can show that for the case of the intercept of the demand model, can clearly separate some customers based on their acquisition characteristics: the bottom customers in the figure (i.e., those with individual posterior means between -2.5 and -2) have clearly higher intercept than the top customers (i.e., those with individual posterior means around -4) as the 95% CPI of the latter group does not overlap with the posterior means of the former.

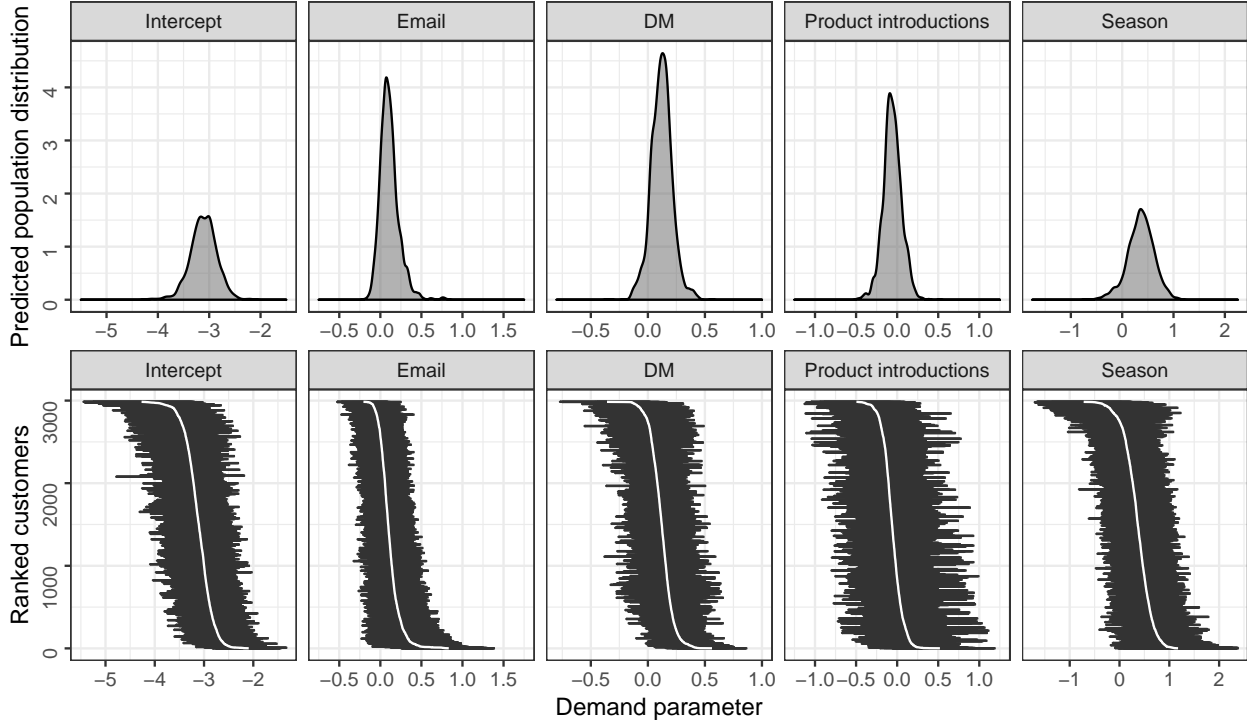
A.7.6 Exploring the latent factors

Figure A12 shows the posterior distribution of weight variances α for each one of the 13 traits. As described in Appendices A.3.1 and A.6.7, each trait parameter α_k controls whether traits are activated by regularizing the weights (\mathbf{W}^y and \mathbf{W}^a) related to the k 'th trait.

We conclude that the first 6 traits carry most of the weight at “connecting” acquisition and demand variables. (Note that the convergence of these parameters, in Figure A13, shows no evidence of label switching or rotation of these traits.) This is not to say that the other traits irrelevant. In turn, those other traits add to the prediction accuracy of the model. However, for deriving insights from the model parameters, we choose to explore the handful of traits that carry most of the information.

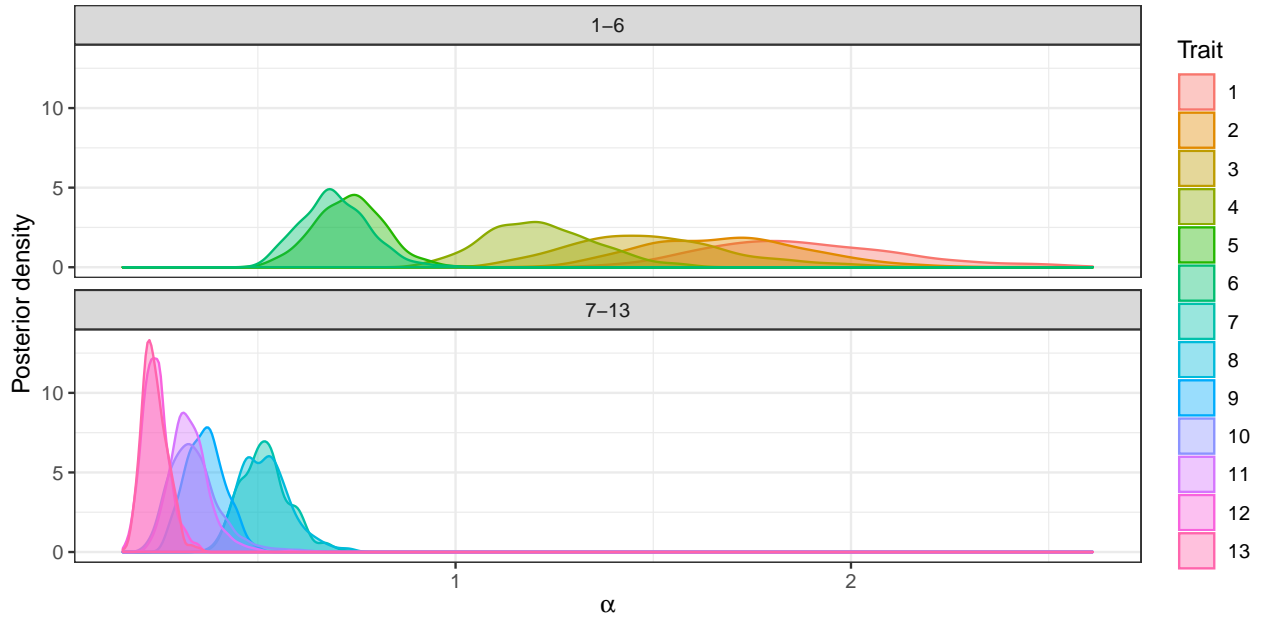
Following the discussion in Appendix A.6.7, we plot in Figure A14a the posterior density of the computed pseudo- α for each upper trait for the FIM model used in our

Figure A11: Population distribution and individual-level posterior distribution for customers in the *Test* sample. The top row shows an histogram of individual-level posterior means for each demand parameter. The bottom row shows customers sorted by posterior means, where the shaded area and the white line represent the individual-level 95% CPI and posterior mean, respectively.



empirical application ($N_1 = 13$, $N_2 = 5$). We find that the relevance of the fifth upper traits is significantly lower than the relevance of the first three traits. This result suggests that $N_2 = 5$ is enough to capture the non-linear correlations present in the data. For robustness, we estimate another FIM specification with $N_2 = 2$ instead, and we find that all upper traits are relevant, suggesting that $N_2 = 2$ may not be enough to capture the non-linear relationships present in the data.

Figure A12: Posterior distribution of α .



A.7.7 Details on the (Machine Learning) benchmark models

We estimate the Feed-Forward DNN model (hidden layer with ReLu as activation function, sigmoid output and cross-entropy loss) using package `torch` in R. We select the value of the weight decay based on the loss calculated using hold-out data in the training sample. After evaluating the values = 0.01, 0.005, 0.001, 0.0005, 0.0001, the value that provides better performance is 0.0001, which we use to estimate the model on the full training sample using 10 epochs. We set the number of hidden dimensions to 128 after corroborating that larger dimensionality does not lead to better fit of the model.

We estimate the Random Forest (RF) using the package `ranger` in R. We finetune the number of trees (`num.trees`), number of variables to possibly split at in each node (`mtry`), and fraction to sample (`sample.fraction`) via cross-validation using the training sample. The resulting values, which we use to estimate the model in the full training data are, `num.trees = 1000`, `sample.fraction = 0.3`, `mtry = 6`.

Figure A13: Convergence of α .

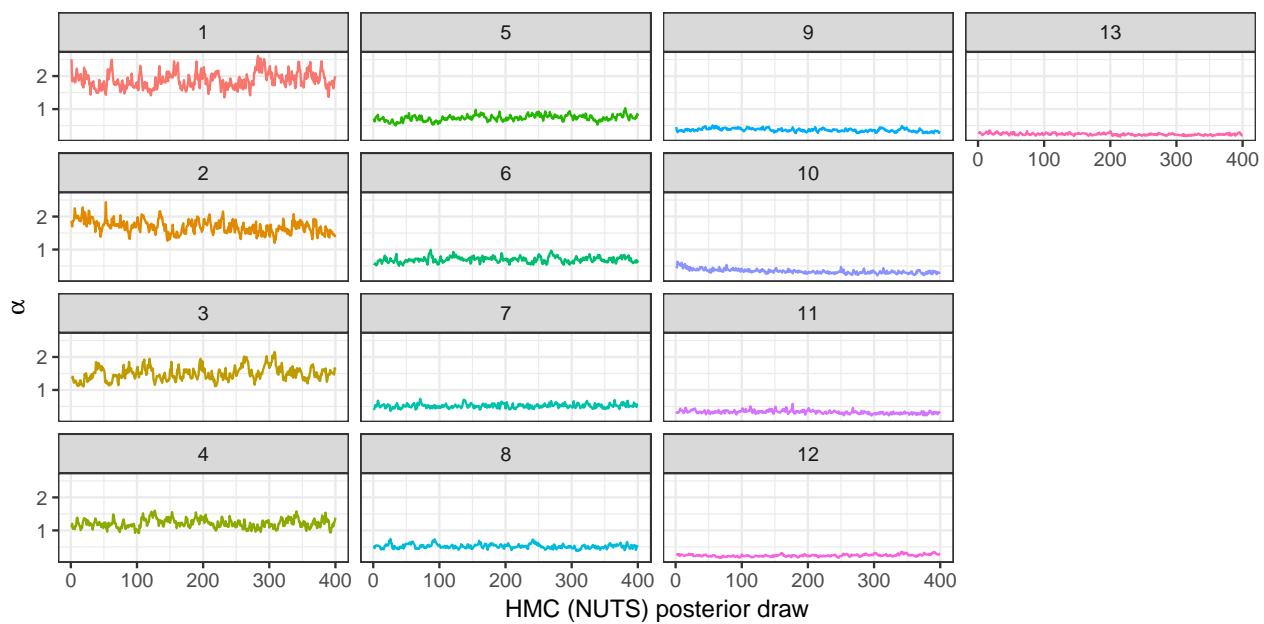
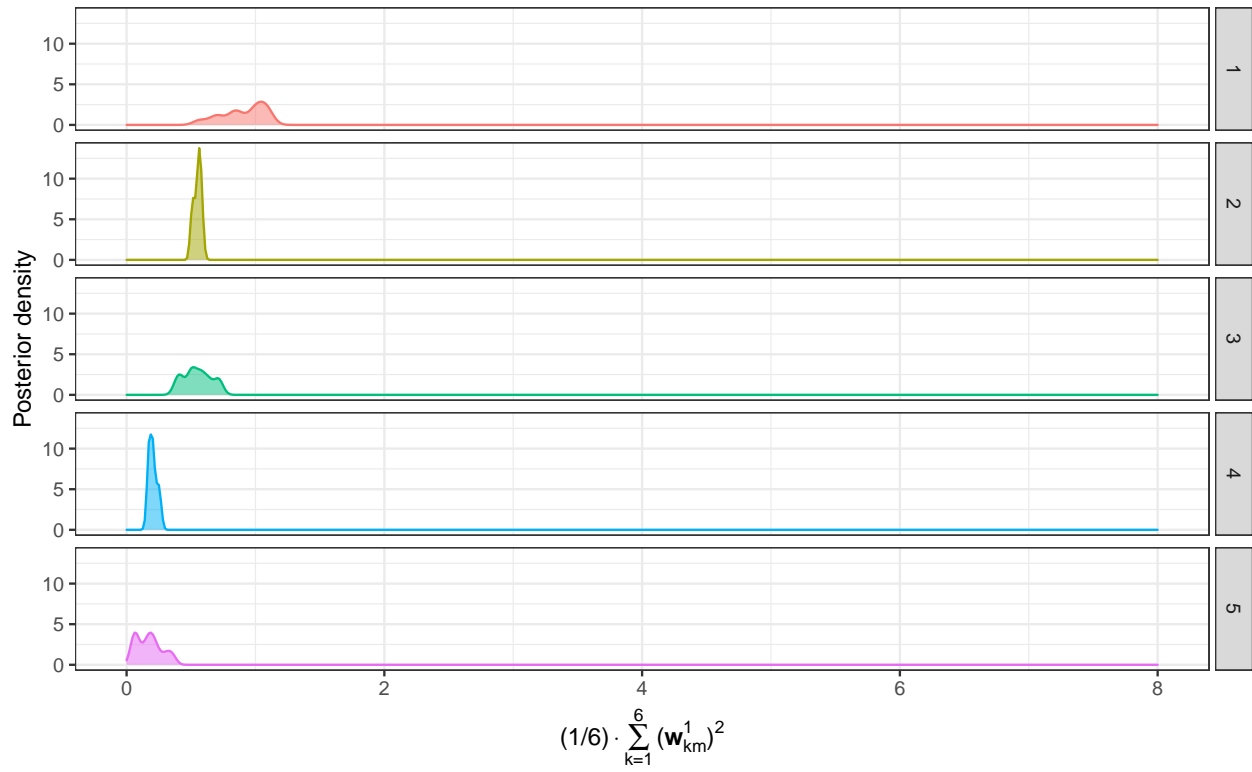
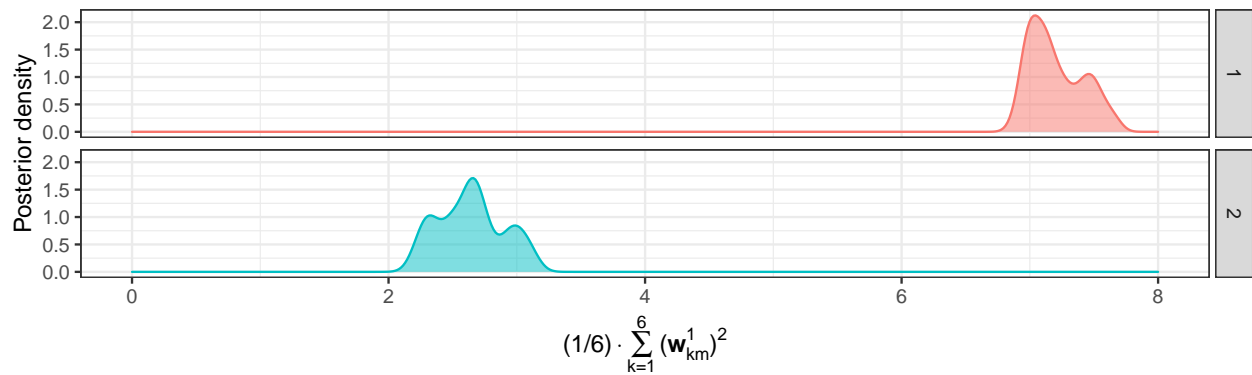


Figure A14: Posterior distribution of pseudo- α^1 .

(a) FIM ($N_1 = 13, N_2 = 5$).



(b) FIM ($N_1 = 13, N_2 = 2$).



Appendix B: Appendix to Essay 2 - The Customer Journey as a Source of Information

B.1 Model priors

We detail the specification of the prior distribution for parameters Σ and $\boldsymbol{\rho}$, and the specification for the base distribution of the Pitman-Yor process F_0 .

We choose the standard Wishart prior for the precision matrix Σ^{-1} ,

$$\Sigma^{-1} \sim \text{Wishart}(r_0, R_0).$$

We put Multivariate Gaussian priors on parameter vector $\boldsymbol{\rho}$,

$$\boldsymbol{\rho} \sim \mathcal{N}(\mathbf{1}, \sigma_\rho^2 \cdot \mathbb{I}),$$

centered at 1.0, which reflects that a priori we do not know whether click decisions are made differently than purchase decisions.

Finally, we put assume a multivariate distribution F_0 , as a product of distributions for each of the components of θ . Following the notation in (2.8), consider θ^q θ^a and θ^p the components of θ that correspond to query parameters and click-purchase parameters. The location parameters are drawn from $\theta_c \sim F_0(\phi_0)$. We assume the multivariate distribution

F_0 to be defined by

$$F_0(\theta|\phi_0) = \left(\prod_{m=1}^M F_{0m}^q(\theta_m^q|\phi_{0m}) \right) \times \mathcal{N}(\theta^p|\mu_0, V_0),$$

where F_{0m}^q is: a Beta distribution if query variable m is a binary variable described in (2.1) by $F_{0m}^q(\phi_{0m}) = \text{Beta}(\phi_{0ma}, \phi_{0mb})$, a Dirichlet distribution if query variable m is a categorical variable described in (2.2) by $F_{0m}^q(\phi_{0m}) = \text{Dirichlet}(\phi_{0m})$, a Gaussian distribution if query variable m is a continuous variable described in (2.3) by $F_{0m}^q(\phi_{0m}) = \mathcal{N}(\phi_{0m\mu}, \phi_{0m\sigma})$, and a Gamma distribution if query variable m is a continuous positive-valued variable described in (2.4) by $F_{0m}^q(\phi_{0m}) = \text{Gamma}(\phi_{0ma}, \phi_{0mb})$.

B.2 Blocked-Gibbs sampler algorithm

Our sampling algorithm is based on Ishwaran and James (2001) approximation using the stick-breaking representation of the Pitman-Yor Process, truncating the infinite mixture by setting $V_C = 1$ for a large enough integer C . This approximation allows us to draw context memberships of different journeys in parallel, which significantly increases the speed of our sampling scheme.

We denote $z_j \in \{1, \dots, C\}$ the context membership latent variable of journey j that captures which context journey j belongs to. Consider a set of drawn values for parameters $\{z_j\}_j, \mathbf{b}_{ij}, \boldsymbol{\mu}_i, \boldsymbol{\mu}_i, \{u_{ijtk}^c\}_{ijtk}, \{u_{ijk}^p\}_{ijk}$. We sequentially update this parameters by,

1. Draw latent click utilities for alternative k by,¹

$$u_{ijtk}^c \sim \begin{cases} \text{Truncated- } \mathcal{N}(\alpha_{1ijp(t)} + \mathbf{x}_k^{p'} \cdot \text{diag}(\boldsymbol{\rho}) \cdot \boldsymbol{\beta}_{ij}, 1, l = -\infty, u = 0) & \text{if } y_{ijt}^c = \ell \\ \text{Truncated- } \mathcal{N}(\alpha_{1ijp(t)} + \mathbf{x}_k^{p'} \cdot \text{diag}(\boldsymbol{\rho}) \cdot \boldsymbol{\beta}_{ij}, 1, l = \max\{u_{ijt-k}^c, 0\}, u = \infty) & \text{if } y_{ij}^p = k \\ \text{Truncated- } \mathcal{N}(\alpha_{1ijp(t)} + \mathbf{x}_k^{p'} \cdot \text{diag}(\boldsymbol{\rho}) \cdot \boldsymbol{\beta}_{ij}, 1, l = -\infty, u = \max\{u_{ijt-k}^c\}) & \text{otherwise.} \end{cases}$$

2. Draw latent purchase utilities by,

$$u_{ijk}^p \sim \begin{cases} \text{Truncated- } \mathcal{N}(\tau_{0ij} + \mathbf{x}_k^{p'} \boldsymbol{\beta}_{ij}, 1, l = -\infty, u = 0) & \text{if } y_{ij}^p = 0 \\ \text{Truncated- } \mathcal{N}(\tau_{0ij} + \mathbf{x}_k^{p'} \boldsymbol{\beta}_{ij}, 1, l = \max\{u_{ij-k}^p, 0\}, u = \infty) & \text{if } y_{ij}^p = k \\ \text{Truncated- } \mathcal{N}(\tau_{0ij} + \mathbf{x}_k^{p'} \boldsymbol{\beta}_{ij}, 1, l = -\infty, u = \max\{u_{ij-k}^p\}) & \text{otherwise.} \end{cases}$$

3. Draw individual level stable preferences $\boldsymbol{\mu}_i$. We define a vector of click and purchase latent utilities for journey j , $\tilde{u}_{ij} = \left[\{u_{ijtk}^c\}_{tk}, \{u_{ijk}^p\}_k \right]'$, and the corresponding matrix of covariates

$$\tilde{\mathbf{X}}_{ij} = \begin{bmatrix} \left[\tilde{\mathbf{X}}_{ij}^c \right]_t \\ \tilde{\mathbf{X}}_{ij}^p \end{bmatrix} = \begin{bmatrix} \left[\begin{array}{cccc} 0 & 1 & 0 & 0 \\ 1 & \mathbf{0} & \mathbf{0} & \mathbf{x}'_{t,1:K} \cdot \text{diag}(\boldsymbol{\rho}) \end{array} \right]_t \\ \left[\begin{array}{ccc} \mathbf{0} & \mathbf{0} & \mathbf{1} \\ & & \mathbf{x}'_{1:K} \end{array} \right] \end{bmatrix},$$

where $\tilde{\mathbf{X}}_{ij}^c$ is the matrix of covariates of click occasion t for customer i and journey j (i.e., intercept dummy variables and product attributes with systematic deviations of preferences $\boldsymbol{\rho}$), and $\tilde{\mathbf{X}}_{ij}^p$ is the matrix of covariates for purchase occasion of customer i

¹For $k = s$ is similar, but using the conditional mean $\alpha_{2ijp(t)}$ instead.

and journey j . The columns of each of these matrices multiply $\alpha_{1ijp(t)}$, $\alpha_{2ijp(t)}$, τ_{0ij} , β_{ij} , respectively; which yields the terms in Equations (2.5) and (2.6).

We also further define $\tilde{\mathbf{X}}_i$ and $\tilde{\mathbf{u}}_i$ as

$$\tilde{\mathbf{X}}_i = \begin{bmatrix} \tilde{\mathbf{X}}_{i1} \\ \vdots \\ \tilde{\mathbf{X}}_{ij} \\ \vdots \\ \tilde{\mathbf{X}}_{iJ_i} \end{bmatrix}, \text{ and} \quad \tilde{\mathbf{u}}_i = \begin{bmatrix} \tilde{u}_{i1} - \tilde{\mathbf{X}}_{i1} \cdot \gamma_{z_1}^p \\ \vdots \\ \tilde{u}_{ij} - \tilde{\mathbf{X}}_{ij} \cdot \gamma_{z_j}^p \\ \vdots \\ \tilde{u}_{iJ_i} - \tilde{\mathbf{X}}_{iJ_i} \cdot \gamma_{z_{J_i}}^p \end{bmatrix}.$$

Finally, we draw

$$\boldsymbol{\mu}_i \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_i, \tilde{S}_i),$$

where

$$\begin{aligned} \tilde{S}_i^{-1} &= \Sigma^{-1} + \tilde{\mathbf{X}}_i' \tilde{\mathbf{X}}_i \\ \tilde{\boldsymbol{\mu}}_i &= \tilde{S}_i \left(\Sigma^{-1} \cdot \mathbf{0} + \tilde{\mathbf{X}}_i' \tilde{\mathbf{u}}_i \right). \end{aligned}$$

4. Draw the vector $\boldsymbol{\rho}$. We define \mathbf{X}_{ij}^ρ by

$$\mathbf{X}_{ij}^\rho = \left[\left[\mathbf{1} \quad \mathbf{0} \quad \mathbf{0} \quad \mathbf{x}'_{t,1:K} \cdot \text{diag}(\mathbf{b}_{ij}) \right]_{tk} \right]$$

and $\mathbf{u}_{ij}^\rho = \left[\{u_{ijtk}^c\}_{tk} \right]'$. We further define $\mathbf{X}^\rho = \left[\{\mathbf{X}_{ij}^\rho\}_{ij} \right]'$, and $\mathbf{u}^\rho = \left[\{\mathbf{u}_{ij}^\rho\}_{ij} \right]'$.

Finally, we draw $\boldsymbol{\rho}$ by,

$$\boldsymbol{\rho} \sim \mathcal{N}(\boldsymbol{\mu}_\rho, S_\rho)$$

where

$$S_\rho^{-1} = \sigma_\rho^{-2} \mathbb{I} + \mathbf{X}^{\rho'} \mathbf{X}^\rho$$

$$\boldsymbol{\mu}_\rho = S_\rho \left(\sigma_\rho^{-2} \cdot \mathbf{1} + \mathbf{X}^{\rho'} \mathbf{u}^\rho \right)$$

5. Draw context membership z_j by

$$p(z_j = c | \cdot) = \frac{\pi_c \mathcal{L}_{jc}}{\sum_{c'=1}^C \pi_{c'} \mathcal{L}_{jc'}}$$

where $\mathcal{L}_{jc} = \left(\prod_{m=1}^M p(q_{ijm} | \theta_{cm}^q) \right) \cdot p\left(\tilde{u}_{ij} - \tilde{\mathbf{X}}_{ij} \boldsymbol{\mu}_i | \tilde{\mathbf{X}}_{ij} \theta_j^p, 1\right)$, $p(q_{ijm} | \theta_{cm}^q)$ is the pdf of query variables as defined in (2.1)-(2.4), and $p\left(\tilde{u}_{ij} - \tilde{\mathbf{X}}_{ij} \cdot \boldsymbol{\mu}_i | \tilde{\mathbf{X}}_{ij} \cdot \theta_j^p, 1\right)$ is the product of elementwise normal pdf evaluated at each components of $\tilde{u}_{ij} - \tilde{\mathbf{X}}_{ij} \cdot \boldsymbol{\mu}_i$ with mean $\tilde{\mathbf{X}}_{ij} \cdot \theta_j^p$ and variance 1.

6. Draw the query components of context location parameters θ_c^q for each context c . We denote $\mathcal{J}(c)$ the set of journeys j such that $z_j = c$, and n_c the number of journeys in that set. For each query variable m , we draw θ_{cm}^q depending on the type of query

variable modeled in (2.1)-(2.4). If m is binary as described in (2.1), we draw

$$\theta_{cm}^q \sim \text{Beta} \left(\phi_{0ma} + \sum_{j \in \mathcal{J}(c)} q_{ijm}, \phi_{0mb} + n_c - \sum_{j \in \mathcal{J}(c)} q_{ijm} \right).$$

If m is categorical as described in (2.2), we draw

$$\theta_{cm}^q \sim \text{Dirichlet}(\phi_{0m} + nq_{cm})$$

where

$$nq_{cm} = \begin{bmatrix} \sum_{j \in \mathcal{J}(c)} \mathbb{1}(q_{ijm} = 1) \\ \vdots \\ \sum_{j \in \mathcal{J}(c)} \mathbb{1}(q_{ijm} = n) \\ \vdots \\ \sum_{j \in \mathcal{J}(c)} \mathbb{1}(q_{ijm} = N_m) \end{bmatrix}_{N_m \times 1},$$

and N_m is the number of possible values of query variable m .

If m is continuous real-valued as described in (2.3), we draw

$$\theta_{cm}^q \sim \mathcal{N}(\tilde{\mu}_{cm}, \tilde{s}_{cm})$$

where $\tilde{s}_{cm}^{-1} = [\phi_{0m\sigma}^{-1} + \sigma_m^{-2}]$ and $\tilde{\mu}_{cm} = \tilde{s}_{cm} \sum_{j \in \mathcal{J}(c)} q_{ijm}$. Finally, if m is positive-valued as described in (2.4), we draw²

$$\theta_{cm}^q \sim \text{Gamma} \left(\phi_{0ma} + n_c, \phi_{0mb} + \sum_{j \in \mathcal{J}(c)} q_{ijm} \right).$$

7. Draw the click-purchase context location parameters θ^p . We denote by $i(j)$ the customer journey j belongs to. We define $\bar{\mathbf{X}}_c$ and $\bar{\mathbf{u}}_c$ as

$$\bar{\mathbf{X}}_c = \left[\left[\tilde{\mathbf{X}}_{i(j)j} \right]_{j \in \mathcal{J}(c)} \right], \text{ and } \bar{\mathbf{u}}_c = \left[\left[\tilde{u}_{i(j)j} - \tilde{\mathbf{X}}_{i(j)j} \cdot \boldsymbol{\mu}_{i(j)} \right]_{j \in \mathcal{J}(c)} \right].$$

We draw θ_c^p by

$$\theta_c^p \sim \mathcal{N}(\bar{\boldsymbol{\mu}}_c, \bar{S}_c),$$

where

$$\begin{aligned} \bar{S}_c^{-1} &= V_0^{-1} + \bar{\mathbf{X}}_c' \bar{\mathbf{X}}_c \\ \bar{\boldsymbol{\mu}}_c &= \bar{S}_c (V_0^{-1} \cdot \boldsymbol{\mu}_0 + \bar{\mathbf{X}}_c' \bar{\mathbf{u}}_c). \end{aligned}$$

8. Draw context probabilities π_c , by drawing the stick parameters V_c from

$$V_c \sim \text{Beta} \left(1 - d + n_c, a + c \cdot d + \sum_{c'=c+1}^C n_{c'} \right)$$

²We use the shape-rate specification of the Gamma distribution (i.e., $\text{Gamma}(\alpha, \beta)$).

9. Draw population covariance matrix Σ , by

$$\Sigma^{-1} \sim \text{Wishart}(r_1, R_1),$$

where

$$r_1 = r_0 + I$$

$$R_1^{-1} = R_0^{-1} + \sum_i \boldsymbol{\mu}_i \cdot \boldsymbol{\mu}_i'$$

B.3 Parameter estimates per context

B.3.1 Query context location parameters

Parameter	Population mean	Context														
		1	2	3	4	5	6	7	9	10	11	13	14	15	17	23
Is it roundtrip?	0.667	0.792	0.455	0.667	0.938	0.728	0.538	0.508	0.656	0.976	0.550	0.769	0.755	0.430	0.812	0.724
Is it domestic? (within EU is domestic)	0.577	0.193	0.854	0.636	0.541	0.493	0.475	0.792	0.876	0.225	0.738	0.304	0.290	0.520	0.808	0.411
Flying from international airport?	0.737	0.772	0.649	0.650	0.631	0.784	0.765	0.653	0.621	0.832	0.761	0.858	0.803	0.822	0.714	0.764
Market:																
Non-US across continent	0.060	0.158	0.013	0.058	0.084	0.052	0.070	0.018	0.008	0.120	0.027	0.084	0.145	0.056	0.020	0.090
Non-US within continent	0.103	0.105	0.122	0.077	0.065	0.111	0.072	0.133	0.125	0.053	0.215	0.076	0.064	0.155	0.103	0.085
US Domestic	0.494	0.146	0.750	0.578	0.456	0.412	0.413	0.665	0.763	0.186	0.553	0.259	0.250	0.407	0.709	0.362
US North America	0.147	0.170	0.073	0.128	0.196	0.240	0.215	0.097	0.063	0.209	0.131	0.239	0.167	0.173	0.091	0.161
US Overseas	0.195	0.421	0.040	0.160	0.199	0.184	0.231	0.087	0.041	0.431	0.075	0.342	0.374	0.209	0.077	0.302
Airport	0.880	0.862	0.881	0.906	0.934	0.903	0.897	0.852	0.909	0.881	0.875	0.837	0.884	0.835	0.885	0.880
Type of location searched:																
Both	0.042	0.078	0.041	0.033	0.021	0.038	0.032	0.053	0.031	0.039	0.035	0.051	0.046	0.060	0.030	0.055
City	0.078	0.060	0.078	0.061	0.045	0.058	0.071	0.095	0.060	0.080	0.090	0.112	0.070	0.105	0.085	0.064
Trip distance (kms)	3820.256	7221.284	2099.494	3452.193	4188.788	3788.162	4161.992	2589.313	2119.567	6558.447	2237.067	5056.982	6198.981	3561.534	2481.910	4715.544
More than one adult?	0.289	0.316	0.243	0.291	0.383	0.344	0.249	0.320	0.293	0.320	0.351	0.277	0.313	0.269	0.280	0.277
Traveling with kids?	0.084	0.094	0.057	0.093	0.152	0.129	0.080	0.102	0.070	0.086	0.126	0.084	0.092	0.077	0.073	0.098
Is it summer season?	0.343	0.255	0.394	0.338	0.372	0.349	0.383	0.310	0.404	0.256	0.383	0.329	0.303	0.365	0.323	0.286
Holiday season?	0.040	0.051	0.028	0.028	0.036	0.055	0.046	0.040	0.024	0.038	0.043	0.046	0.025	0.044	0.053	0.053
Does stay include a weekend?	0.667	0.819	0.478	0.640	0.887	0.708	0.595	0.554	0.617	0.947	0.573	0.766	0.769	0.483	0.745	0.729
Length of stay (only RT) (days)	10.386	16.881	5.446	8.579	9.570	8.656	10.002	8.391	5.112	17.608	7.759	13.687	16.471	15.749	6.588	13.401
Searching on weekend?	0.216	0.217	0.211	0.200	0.192	0.211	0.216	0.179	0.220	0.237	0.240	0.227	0.237	0.228	0.195	0.210
Searching during work hours?	0.487	0.463	0.497	0.542	0.443	0.531	0.480	0.497	0.519	0.494	0.446	0.463	0.478	0.427	0.513	0.459
Time in advance to buy (days)	55.339	67.616	38.060	50.084	62.049	57.785	59.355	44.259	39.881	76.156	47.818	64.167	64.045	52.447	59.501	66.915

Table B1: Posterior mean of query location parameters per context

B.3.2 Click and purchase context location parameters

Parameter	Population mean	Context														
		1	2	3	4	5	6	7	9	10	11	13	14	15	17	23
Click intercepts																
Intercept Search: OW Search	-1.661	-0.934	-1.881	-1.609	3.261	-1.962	-2.091	-1.160	-1.952	-0.986	-2.349	-1.533	-1.401	-2.411	-1.984	-1.645
Intercept Search: RT Outbound	-2.455	-2.317	-2.222	-2.212	-2.126	-2.610	-2.383	-2.055	-2.459	-2.341	-3.331	-2.734	-2.331	-2.529	-2.854	-2.323
Intercept Search: RT Inbound	-2.651	-2.353	-2.405	-2.342	-2.167	-2.751	-2.406	-2.124	-2.577	-2.474	-3.883	-3.031	-2.437	-2.154	-3.551	-2.553
Intercept Click: OW Search	-1.498	-0.597	-2.011	-1.293	-0.771	-1.469	-1.573	-1.372	-1.683	-1.497	-0.660	-1.197	-1.269	-1.346	-2.030	-1.437
Intercept Click: RT Outbound	-0.736	0.244	-1.235	-0.634	-0.317	-0.603	-0.816	-0.440	-0.773	-0.701	-0.347	-0.643	-0.520	-0.818	-0.926	-0.737
Intercept Click: RT Inbound	0.087	1.278	-0.731	-0.087	0.474	0.148	-0.428	0.467	-0.293	-0.006	0.733	0.667	0.672	-0.317	0.582	0.047
Control for whether product was clicked before	1.257	1.180	1.411	1.810	0.635	1.673	0.791	1.729	1.711	1.039	1.856	1.067	1.096	1.048	0.946	1.280
Purchase intercept																
Intercept Purchase	-5.550	-2.909	-6.684	-4.426	-3.898	-5.382	-6.125	-3.428	-5.834	-6.177	-4.754	-4.878	-4.435	-6.285	-6.566	-5.207
Product attribute preferences																
Price	-0.567	-0.740	-0.480	-0.602	-0.691	-0.633	-0.667	-0.622	-0.494	-0.680	-0.558	-0.495	-0.690	-0.535	-0.436	-0.554
Length of trip (hours)	-0.737	-0.605	-0.751	-0.837	-0.627	-0.901	-0.608	-0.941	-0.956	-0.555	-1.339	-0.678	-0.622	-0.615	-0.659	-0.729
Number of stops: Non stop	0.023	0.039	0.032	0.241	0.306	0.170	0.129	-0.027	0.061	0.291	-0.535	-0.225	-0.012	0.083	-0.157	0.304
Number of stops: 2+ stops	-1.621	-0.610	-1.979	-1.174	-0.487	-2.053	-1.534	-1.102	-1.422	-1.206	-1.028	-1.529	-0.560	-2.778	-2.414	-2.850
Alliance: Skyteam (Delta)	-0.564	-0.192	-0.324	-0.380	-0.457	-0.431	-0.532	-0.032	-0.446	-0.361	-0.327	-0.733	-0.291	-1.661	-1.086	-0.252
Alliance: Star Alliance (United)	-0.367	-0.093	-0.231	-0.351	-0.336	-0.213	-0.312	-0.032	-0.262	-0.069	-0.694	-0.549	-0.174	-1.729	-0.774	-0.477
Alliance: Alaska Airlines	-0.497	-0.513	-0.416	-0.402	-0.098	-0.426	-0.321	-0.386	-0.243	-0.536	-0.304	-0.195	-0.256	-1.120	-0.709	-1.791
Alliance: Spirit	-0.667	-0.636	-0.357	-0.780	-0.990	-1.624	-0.389	-0.462	-0.176	-0.450	-1.651	-0.307	-0.376	-0.661	-0.849	-3.928
Alliance: JetBlue	-0.097	-0.378	0.017	-0.656	0.213	-0.003	0.027	0.108	-0.257	-0.127	0.099	-0.237	-0.073	0.043	0.017	-0.013
Alliance: Frontier	-0.130	0.398	0.161	-0.205	0.446	-0.372	-0.014	-0.645	-0.225	-0.508	-0.661	0.147	-0.149	0.048	-0.236	-0.535
Alliance: Other – No alliance	-0.228	-0.064	-0.068	-0.336	-0.904	-0.106	-0.370	-0.012	0.194	-0.232	-0.492	-0.175	-0.120	-0.878	-0.263	-0.051
Alliance: Multiple alliances	-1.542	-0.525	-1.638	-1.971	-1.508	-1.584	-0.974	-1.457	-1.634	-0.879	-2.573	-1.411	-0.701	-3.102	-1.700	-1.604
Outbound dep. time: Early morning (0:00am - 4:59am)	-0.638	-0.406	-0.421	-0.103	0.077	0.041	-0.513	0.256	-0.214	-0.249	-2.954	-1.062	-0.324	-1.667	-0.854	-1.347
Outbound dep. time: Afternoon (12:00pm - 5:59pm)	-0.162	-0.042	-0.053	-0.225	-0.098	-0.166	-0.332	-0.018	0.117	-0.066	0.157	-0.319	0.092	-0.593	-0.352	-0.279
Outbound dep. time: Evening (6:00pm - 11:59pm)	-0.226	-0.248	-0.048	-0.224	-0.161	-0.392	-0.471	-0.153	-0.020	-0.285	-0.225	-0.461	-0.058	-0.298	-0.239	-0.239
Outbound arr. time: Early morning (0:00am - 4:59am)	-0.835	-0.675	-1.038	-0.386	-1.168	-0.713	-0.683	-0.369	-1.021	-0.693	-0.776	-0.658	-0.917	-1.322	-0.265	-0.375
Outbound arr. time: Afternoon (12:00pm - 5:59pm)	-0.160	-0.101	-0.265	0.240	0.145	0.269	-0.009	0.265	-0.098	-0.150	0.080	-0.565	-0.109	-0.480	-0.348	-0.486
Outbound arr. time: Evening (6:00pm - 11:59pm)	-0.213	-0.346	-0.457	0.228	-0.049	0.269	-0.116	0.147	-0.227	-0.164	-0.110	-0.151	-0.152	-0.229	-0.527	-0.486
Inbound dep. time: Early morning (0:00am - 4:59am)	-0.964	-0.301	-1.054	-0.539	-5.904	-2.102	-0.405	0.359	-1.080	-0.181	-3.157	-1.423	0.183	-0.859	-1.035	-1.551
Inbound dep. time: Afternoon (12:00pm - 5:59pm)	-0.146	-0.163	-0.188	0.283	0.181	0.899	0.090	-0.031	-0.143	0.148	-0.331	0.067	0.193	-0.267	-1.531	0.459
Inbound dep. time: Evening (6:00pm - 11:59pm)	-0.486	-0.212	-0.331	-0.081	-0.836	0.048	-0.342	-0.205	-0.058	-0.119	-0.381	-1.495	0.060	-0.198	-1.886	0.232
Inbound arr. time: Early morning (0:00am - 4:59am)	-0.886	-0.496	-0.212	-0.164	0.189	-1.022	-0.602	-0.423	-0.344	-0.374	-1.746	-2.449	-0.579	-0.897	-2.090	-1.545
Inbound arr. time: Afternoon (12:00pm - 5:59pm)	-0.665	-0.063	-0.097	0.095	0.455	-1.008	0.103	-0.070	-0.192	0.125	-0.626	-1.744	-0.515	-0.409	-3.163	-0.509
Inbound arr. time: Evening (6:00pm - 11:59pm)	-0.078	-0.005	0.532	0.315	0.705	-0.206	0.454	0.374	0.155	0.283	-0.416	-1.792	-0.280	0.432	-0.905	-0.355

Table B2: Posterior mean of location click and purchase parameters