



Marketing Science

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Model-Based Purchase Predictions for Large Assortments

Bruno J.D. Jacobs, Bas Donkers, Dennis Fok

To cite this article:

Bruno J.D. Jacobs, Bas Donkers, Dennis Fok (2016) Model-Based Purchase Predictions for Large Assortments. Marketing Science 35(3):389-404. <https://doi.org/10.1287/mksc.2016.0985>

Full terms and conditions of use: <http://pubsonline.informs.org/page/terms-and-conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2016, INFORMS

Please scroll down for article—it is on subsequent pages



INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Model-Based Purchase Predictions for Large Assortments

Bruno J.D. Jacobs

Department of Business Economics and Econometric Institute, Erasmus School of Economics, Erasmus University Rotterdam,
3000 DR Rotterdam, Netherlands, jacobs@ese.eur.nl

Bas Donkers

Department of Business Economics, Erasmus School of Economics, Erasmus University Rotterdam,
3000 DR Rotterdam, Netherlands, donkers@ese.eur.nl

Dennis Fok

Econometric Institute, Erasmus School of Economics, Erasmus University Rotterdam,
3000 DR Rotterdam, Netherlands, dfok@ese.eur.nl

An accurate prediction of what a customer will purchase next is of paramount importance to successful online retailing. In practice, customer purchase history data is readily available to make such predictions, sometimes complemented with customer characteristics. Given the large product assortments maintained by online retailers, scalability of the prediction method is just as important as its accuracy. We study two classes of models that use such data to predict what a customer will buy next, i.e., a novel approach that uses latent Dirichlet allocation (LDA), and mixtures of Dirichlet-Multinomials (MDM). A key benefit of a model-based approach is the potential to accommodate observed customer heterogeneity through the inclusion of predictor variables. We show that LDA can be extended in this direction while retaining its scalability. We apply the models to purchase data from an online retailer and contrast their predictive performance with that of a collaborative filter and a discrete choice model. Both LDA and MDM outperform the other methods. Moreover, LDA attains performance similar to that of MDM while being far more scalable, rendering it a promising approach to purchase prediction in large product assortments.

Data, as supplemental material, are available at <http://dx.doi.org/10.1287/mksc.2016.0985>.

Keywords: model-based predictions; large scale purchase prediction; scalability; purchase history data; latent Dirichlet allocation; mixture of Dirichlet-Multinomials

History: Received: December 23, 2013; accepted: January 27, 2016; Pradeep Chintagunta, Dominique Hanssens, and John Hauser served as the special issue editors and Daniel Goldstein served as associate editor for this article. Published online in *Articles in Advance* April 18, 2016.

1. Introduction

The ability to predict what a customer will purchase next is valuable in many marketing applications. This is especially true for online retailing. Adequate predictions for the next products to be purchased enable online retailers to implement a product recommendation system, determine the positions of products in the result of a customer's search query, optimize the collection of products to be displayed on a personalized landing page or suggest products to complement the contents of a customer's shopping basket.

Examples of personalization in practice are Amazon's "Customers Who Bought This Item Also Bought" section, Apple's iTunes Genius, and the Netflix recommendation system. There is also clear evidence that such personalized configurations influence behavior (Ghose et al. 2014, Pan et al. 2007, Salganik et al. 2006). These applications have in common that they require a personalized selection of products from a potentially large product assortment. Ideally, the selection contains those products that are most likely to be of interest

to the customer. Moreover, the selection should be relatively small as the available space to show products is often limited.

The effectiveness of personalization attempts crucially depends on the accuracy of the predictions. A complicating factor in purchase prediction is the fact that the typical online retailer sells items from a very broad assortment to an even larger customer base. Hence predictions should not only be accurate, but the prediction method should scale to large applications as well (Naik et al. 2008). Additionally, to be useful in an online setting predictions should be available in near real-time. Obtaining predictions, and updating them as new information comes in, should therefore be fast.

The typical data available at an online retailer for purchase prediction are the customer purchase histories. In some cases additional customer characteristics (e.g., demographics) are also available. However, on the product level characteristics are often absent and if descriptions are available, it is not obvious how to extract useful predictors from this information. In this

paper we focus on predicting purchase behavior based on purchase history data, possibly complemented with customer characteristics.

Many online retailers predict a customer's next purchase using collaborative filtering algorithms, for example, by relying on counts of the co-occurrence of items in purchase history data (Jannach et al. 2011, Liu et al. 2009). In such a count-based approach a decision must be made on how to measure the co-occurrence of items, as one can count pairs, triplets or even higher-order product combinations. A choice for small sets of items results in information loss, i.e., purchase patterns that span many products might not be easily identified. On the other hand, for large combinations of products the matrix of co-occurrence counts becomes sparse, resulting in predictions based on just a few matches in the customer base. Another challenge in collaborative filtering algorithms is incorporating customer characteristics. One possible approach is to partition the customer base using such characteristics. However, this can only be done for a few variables with a limited number of levels; otherwise sample sizes per subgroup become too small.

By contrast, model-based approaches to predict individuals' choices have a long history in marketing (Guadagni and Little 1983, McFadden 1986, Wagner and Taubes 1986, Fader and Hardie 1996) and such methods are well suited to include customer characteristics. However, the usual implementations of these models tend to break down in the typical online retail setting, where a wide variety of products is sold to a large number of customers (Naik et al. 2008). One way to make methods more scalable is to consider only a subset of the data in terms of customers and/or products (Zanutto and Bradlow 2006). Clearly, this is not a viable solution if the aim is to predict purchase behavior for each customer across the entire product assortment.

In this paper we try to bridge the gap between retail practice and marketing academia by discussing model-based prediction methods that do work in the context of large assortments. By developing such methods we open an avenue for future research on marketing interventions in large scale assortments, for example on the effectiveness of product recommendations, extending the work of Bodapati (2008). Note that this would not be feasible with the heuristic, count-based approaches currently used in practice. We consider two model-based approaches. In addition, we present (an implementation of) a count-based collaborative filter and a scalable version of a discrete choice model that will serve as benchmarks. We compare the methods on their (i) heterogeneity assumptions, (ii) estimation complexity, (iii) memory requirements for real-time online predictions, and (iv) predictive performance.

The first method we present is a novel approach inspired by topic models used in the text modeling literature. Traditionally, a topic model describes a document by relating the words in the text to latent topics. We adapt this class of models to the purchase prediction context: Words become product purchases, a document is a customer's purchase history, and a topic represents a certain preference for products in the assortment. Given that the word "topic" is incongruous in a retailing context, we refer to topics as motivations.¹ Naturally, customers can have more than one motivation, just as a document can cover multiple topics. This idea leads to a class of models that can describe and predict customer purchase behavior in large assortments.

The most frequently used topic model is latent Dirichlet allocation (LDA) by Blei et al. (2003). This model has been used to analyze very large text corpora (Ramage et al. 2010, Mimno et al. 2012), showing that LDA provides the necessary scalability. By contrast to the text modeling literature, where documents tend to have many words, customers often have only a couple of purchases or they might even be new to the retailer. Given such limited information per customer, we need to formally estimate the population-level a priori probabilities of particular motivations. This extends the LDA text modeling implementation, where these probabilities are typically considered to be known or at best calibrated using heuristics (Wallach et al. 2009, Asuncion et al. 2009).

To account for observed heterogeneity, we extend LDA by relating customer characteristics to the a priori motivation probabilities. This can capture heterogeneity related to variables such as referrer site, demographics or other customer characteristics. Most likely this increases the model's predictive power in particular for customers with few or no observed purchases. We refer to this model as LDA-X.

The next method we consider is a mixture of Dirichlet-Multinomials (MDM) (Jain et al. 1990). MDM specifies individual-specific probability vectors that contain a customer's purchase probabilities over all products in the assortment. In turn, these probability vectors follow a discrete mixture of Dirichlet distributions. MDM has previously been applied in marketing (Jain et al. 1990), but to our knowledge never to a large product assortment. Although, in theory, customer characteristics can also be included in MDM we will argue that the resulting model will no longer be feasible in terms of estimation complexity, given the setting of our application.

The predictive performance of LDA(-X) and MDM is compared to that of a count-based collaborative filter and a discrete choice model. We assess predictive

¹ While intuitively plausible, we do not claim that the actual decision process is driven by these motivations.

performance using data from an online retailer. For each method, we create customer-specific prediction sets that contain the products most likely to be purchased. These sets are next matched with hold-out purchase data. For further insight into the differences between the methods, we also consider the predictive performance for groups of customers who differ in the length of their observed purchase history. Furthermore, in a setting where repeat purchases are frequent, e.g., fast moving consumer goods, performing well by correctly predicting frequently purchased products or repeat purchases might not be too difficult. Such recommendations might even be perceived as trivial or boring (Fleder and Hosanagar 2009). Therefore we also study the predictive performance for *unexpected* products, which we define as products that have not previously been purchased by the customer and which are in the tail of the assortment.

The remainder of this article proceeds as follows: In §2 we present the methods used in this research and discuss their heterogeneity assumptions and scalability. Technical details are available in the appendices. Subsequently, we apply the methods to data of an online retailer. An overview of this data is provided in §3 and the results are reported in §4. To conclude, we summarize our findings and suggest directions for future research in §5.

2. Methods

In this section we present our prediction methods. First, we introduce two model-based prediction methods, LDA(-X) and MDM, that infer latent customer traits from purchase data. We compare these methods on their heterogeneity assumptions and estimation complexity. Next, the two benchmark methods are introduced, i.e., a set of collaborative filters (CF) and a discrete choice model (DCM) that captures customer heterogeneity through constructed, but *observed* predictor variables. Subsequently, all methods are compared on their suitability to update predictions in a real-time setting. Finally, we discuss our assessment of the quality of predictions.

All methods share the following notation: The products from the J -dimensional assortment are numbered $j = 1, \dots, J$. For each customer $i = 1, \dots, I$ we observe n_i product purchases from this assortment. The purchase history of customer i is denoted by the vector $\mathbf{y}_i = [y_{i1}, \dots, y_{in_i}]$, where $y_{in} \in \{1, \dots, J\}$ represents the n th purchase of customer i . In addition we have customer-level characteristics coded in the K -dimensional column vector $\mathbf{x}_i = [x_{i1}, \dots, x_{iK}]'$. We combine the purchase histories in $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_I\}$ and the predictor variables in $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_I\}$.

2.1. Latent Dirichlet Allocation

Our first model is inspired on topic models. The key idea underlying our application of these models to the context of purchase history data is that customer purchases are driven by a (small) set of latent motivations (the topics). Each motivation then drives preferences for a subset of products in the assortment, for example, a preference for eco-friendly products, for low-fat products or for products for sensitive skin.

In general, customers are likely to be driven by different motivations over time and even within a single purchase occasion. Additionally, the same product purchased by different customers may be driven by different underlying motivations: A movie can be purchased by a fan of the lead actor or by a customer that is fond of the movie's genre. These features are embedded in topic models, in which customers may have multiple motivations, and products may be associated with more than one motivation.

The basis for our method is latent Dirichlet allocation (LDA) introduced by Blei et al. (2003). LDA has been proven to scale to applications well beyond the dimensions of a typical online retailer. For example, it has been used to analyze over eight million posts on Twitter that contain words from a vocabulary of more than five million entries (Ramage et al. 2010), or for the analysis of 1.2 million out-of-copyright books (Mimno et al. 2012). Below, we first present the details of our adaptation of LDA in the context of predicting customer purchase behavior. Next, we extend LDA by including customer-level predictor variables.

In LDA each latent motivation $m = 1, \dots, M$ is represented by a probability vector $\boldsymbol{\phi}_m$ over the complete J -dimensional assortment. Given that a purchase is driven by motivation m , the probability of buying product j is simply ϕ_{mj} . The motivation-specific probability vectors are distributed as $\boldsymbol{\phi}_m | \boldsymbol{\beta} \sim \text{Dirichlet}_J(\boldsymbol{\beta})$.

A priori there is no reason to favor one product over another in a motivation. This is reflected in the parameterization of $\boldsymbol{\beta}$, where we set each element equal to a common value β_0 . This value determines whether the Dirichlet distribution tends to favor more narrow (β_0 close to zero) or more broad (large β_0) motivations (Wallach et al. 2009).

Even though each purchase is driven by a single motivation, a customer's entire purchase history may be driven by multiple motivations. This variation is described by an individual-specific discrete mixture $\boldsymbol{\theta}_i$ over the M motivations. The probability that a product purchase of customer i is driven by motivation m is then given by θ_{im} . These probabilities differ across customers and are modeled as $\boldsymbol{\theta}_i | \boldsymbol{\alpha} \sim \text{Dirichlet}_M(\boldsymbol{\alpha})$.

Here, $\boldsymbol{\alpha}$ is an M -dimensional vector that captures the relevance of each motivation across the customer base.

The expected value of the probability that motivation m drives a purchase equals

$$\mathbb{E}[\theta_{im} | \alpha] = \frac{\alpha_m}{\sum_{l=1}^M \alpha_l}. \quad (1)$$

Therefore, the larger the value of α_m , the more likely it is that a customer will make a purchase driven by motivation m .

The last step is to link motivations to actual purchases. We denote by $z_{in} \in \{1, \dots, M\}$ the actual motivation that drives purchase y_{in} . As motivations are latent, we must account for all possible motivations to obtain the marginal probability that customer i will purchase product j , resulting in

$$\begin{aligned} \Pr[y_{in} = j | \{\phi_l\}_{l=1}^M, \theta_i] \\ &= \sum_{m=1}^M \Pr[y_{in} = j | z_{in} = m, \{\phi_l\}_{l=1}^M] \Pr[z_{in} = m | \theta_i] \\ &= \sum_{m=1}^M \phi_{mj} \theta_{im}. \end{aligned} \quad (2)$$

In the topic modeling literature it is common practice to determine the parameters of the Dirichlet distributions α (for θ_i) and β_0 (for ϕ_m) by means of heuristics, rather than formally inferring their values from available data (Wallach et al. 2009), for example, imposing $\alpha = 50/M$ (Griffiths and Steyvers 2004) and $\beta_0 = 0.01$ (Steyvers and Griffiths 2007), or by applying a grid search for α and β_0 (Asuncion et al. 2009). These heuristics are not directly applicable in our setting as they have been designed for text modeling. Given that purchase histories tend to be much shorter than documents, we expect the LDA predictions to be more sensitive to the values of α and (to a lesser degree) of β_0 . We therefore extend the common LDA model and place proper prior distributions on both parameters and formally estimate α and β_0 in a Bayesian setting.

We specify a log-normal distribution for α_m , that is, we define $\log(\alpha_m) = \gamma_m$, and set a normal prior for γ_m . We set the mode of the log-normal distribution equal to M^{-1} , which is within the range of values frequently used in the text modeling literature, and place 10% of its probability mass above 1.² This prior specification favors θ_i -vectors that allocate the majority of the probability mass to a small number of motivations, while still allowing for more uniformly distributed θ_i -vectors. Similarly, we place a log-normal distribution on β_0 with its mode equal to 0.01, and 10% of its probability mass above 1. This specification supports ϕ_m -vectors where only a few products from the assortment receive significant probability mass, representing fairly specific

motivations. Still, this prior is rather uninformative; broader motivations that spread the probability mass more equally over the assortment remain quite likely.

These prior specifications also allow us to easily extend LDA by including customer characteristics, coded in \mathbf{x}_i . Such variables are likely to improve the predictive performance of the model. We extend the log-linear specification for α_m to α_{im} as follows: $\log(\alpha_{im}) = \gamma_m + \mathbf{x}_i' \delta_m$. This links customer preferences, represented by the likelihood of each of the motivations, to the additional customer-level information, resulting in LDA-X. To illustrate the effect of this specification on the distribution of θ_i , consider the expected value of θ_{im} , which gives the probability that a typical customer with characteristics \mathbf{x}_i makes a purchase driven by motivation m

$$\mathbb{E}[\theta_{im} | \alpha_i] = \frac{\alpha_{im}}{\sum_{l=1}^M \alpha_{il}} = \frac{\exp(\gamma_m + \mathbf{x}_i' \delta_m)}{\sum_{l=1}^M \exp(\gamma_l + \mathbf{x}_i' \delta_l)}. \quad (3)$$

The δ_m parameters capture the dependence of the probability that motivation m is used, on the customer-specific variables \mathbf{x}_i . The prior distribution of γ_m and δ_m can only be sensibly determined if the level and scale of the \mathbf{x}_i variables are known. We therefore standardize the customer-level variables such that they have mean zero and unit variance. Given this scale, we assume that all elements in δ_m are normally distributed with zero mean and variance equal to 0.04. This corresponds to a prior 95% confidence interval that is approximately equal to $[-0.4, +0.4]$. Note that this prior distribution is chosen to be relatively narrow on purpose, as the effect of δ_m on α_{im} is exponential. As \mathbf{x}_i is mean-centered, we use the same prior for γ_m as in LDA.

To obtain customer-specific predictive distributions, we condition on the model structure of LDA. In particular, given the model parameters α , β_0 and the latent purchase assignments \mathbf{Z} , the predictive distribution for a new purchase \tilde{y}_{in} can be shown to equal (Griffiths and Steyvers 2004)

$$\begin{aligned} \Pr[\tilde{y}_{in} = j | \mathbf{Z}, \alpha, \beta_0, \mathbf{Y}] \\ &= \sum_{m=1}^M \Pr[\tilde{y}_{in} = j | \tilde{z}_{in} = m, \mathbf{Z}, \beta_0, \mathbf{Y}] \Pr[\tilde{z}_{in} = m | \mathbf{z}_i, \alpha] \\ &= \sum_{m=1}^M \mathbb{E}[\phi_{mj} | \mathbf{Z}, \beta_0, \mathbf{Y}] \mathbb{E}[\theta_{im} | \mathbf{z}_i, \alpha] \\ &= \sum_{m=1}^M \frac{\beta_0 + c_{mj}^{\text{MJ}}}{\beta_0 + \sum_{p=1}^J c_{mp}^{\text{MJ}}} \cdot \frac{\alpha_m + c_{im}^{\text{IM}}}{\sum_{l=1}^M \alpha_l + c_{il}^{\text{IM}}}, \end{aligned} \quad (4)$$

where c_{mj}^{MJ} is the number of times a purchase of product j is driven by motivation m and c_{im}^{IM} is the number of purchases made by customer i that are driven by motivation m . To obtain the predictive distribution for the LDA-X model one simply replaces α with α_i in (4).

² These two conditions implicitly identify the two parameters of the log-normal distribution.

2.2. Dirichlet-Multinomial Models

The Dirichlet-Multinomial (DM) model (Jeuland et al. 1980, Goodhardt et al. 1984) is a known model-based approach to capture heterogeneity in purchase behavior. Applications of this model can be found in Grover and Srinivasan (1987), Fader (1993), and Fader and Schmittlein (1993). In this model, each customer is given an individual-specific vector ϕ_i containing the purchase probabilities for each product in the J -dimensional assortment, where $\sum_{p=1}^J \phi_{ip} = 1$. The probability that customer i purchases product j at a specific purchase occasion n is given by $\Pr[y_{in} = j | \phi_i] = \phi_{ij}$.

Large values for the purchase probability ϕ_{ij} imply that customer i is likely to buy product j . In the DM model the customer-specific ϕ_i -vectors are assumed to arise from a single Dirichlet distribution: $\phi_i | \beta \sim \text{Dirichlet}_J(\beta)$. The β -vector describes the overall purchase behavior in the customer base: If product j is frequently purchased, β_j will have a large value relative to the other values in β and vice versa.

The original DM model has been extended such that the probability vectors ϕ_i originate from a finite mixture of Dirichlet distributions (Jain et al. 1990), not from a single Dirichlet distribution. This extension is known as a mixture of Dirichlet-Multinomials (MDM). In MDM, each customer is assigned to one of M segments and each segment is characterized by its own Dirichlet distribution. Given that customer i is allocated to segment m , denoted by $s_i = m$, the customer's purchase probabilities ϕ_i are distributed as $\text{Dirichlet}_J(\beta_m)$. The β_m -vectors are segment specific, describing the distribution of the purchase probability vectors for customers in segment m . Customers are hence expected to be similar, though not identical, within a segment, but different across segments.

Segment membership in MDM is described by an M -dimensional categorical distribution with probability vector π . The element π_m gives the a priori probability that a customer is a member of segment m , that is, $\Pr[s_i = m | \pi] = \pi_m$.

As we consider MDM in the Bayesian paradigm we also specify prior distributions over π and the β_m -vectors. For π it is natural to favor no segment over any other a priori, therefore we use a uniform distribution over the $(M - 1)$ -dimensional simplex. For each β_{mj} we use a log-normal prior distribution with its mode at 0.01 and 10% of the probability mass above 1. This specification allows for ϕ_i -vectors that allow many products to be purchased with a large probability, but it favors customer segments who purchase from a more limited subset of the assortment.

Similar to the approach in LDA(-X), we obtain customer-specific predictive distributions of a new purchase \tilde{y}_{in} conditional on the data, parameters, and segment allocations. In MDM this requires a prediction of segment membership of the customer, combined

with the purchase probabilities, conditional on segment membership

$$\begin{aligned} \Pr[\tilde{y}_{in} = j | \mathbf{s}^i, \{\beta_l\}_{l=1}^M, \mathbf{y}_i] \\ &= \sum_{m=1}^M \Pr[\tilde{y}_{in} = j | s_i = m, \beta_m, \mathbf{y}_i] \Pr[s_i = m | \mathbf{s}^i, \{\beta_l\}_{l=1}^M, \mathbf{y}_i] \\ &= \sum_{m=1}^M \mathbb{E}[\phi_{ij} | s_i = m, \beta_m, \mathbf{y}_i] \Pr[s_i = m | \mathbf{s}^i, \{\beta_l\}_{l=1}^M, \mathbf{y}_i] \\ &= \sum_{m=1}^M \left(\frac{\beta_{mj} + c_{ij}^{\text{II}}}{\sum_{p=1}^J \beta_{mp} + c_{ip}^{\text{II}}} \right) \Pr[s_i = m | \mathbf{s}^i, \{\beta_l\}_{l=1}^M, \mathbf{y}_i], \quad (5) \end{aligned}$$

where $\Pr[s_i = m | \mathbf{s}^i, \{\beta_l\}_{l=1}^M, \mathbf{y}_i]$ is specified in Online Appendix 1 (available as supplemental material at <http://dx.doi.org/10.1287/mksc.2016.0985>) (see Equation (A.9)) and c_{ij}^{II} equals the number of times customer i has purchased product j . If i is a new customer $c_{ij}^{\text{II}} = 0$ for all j by definition. Note that both components in (5) depend on the customer's purchase history, unlike LDA(-X) where only the motivation probabilities are customer specific.

2.3. Model Inference

The predictive distributions specified above are conditional on the number of segments/motivations M , the model parameters, and segment/motivation allocations to customers/purchases. For a given number of M , we rely on Bayesian methodology to infer the model parameters and latent variables of the models. Direct inference on the posterior distribution is not tractable; therefore, we derive Markov Chain Monte Carlo (MCMC) methods to generate samples from the posterior distribution. Specifically, we use a random walk Metropolis–Hastings within Gibbs sampler to draw samples from the target posterior distribution. The predictive distributions can then be obtained by averaging over these draws.

The full posterior of LDA(-X) is given by

$$p(\mathbf{Z}, \{\phi_l\}_{l=1}^M, \beta_0, \{\theta_l\}_{l=1}^I, \gamma, \{\delta_l\}_{l=1}^M | \mathbf{Y}, \mathbf{X}),$$

where $\{\delta_l\}_{l=1}^M$ is only relevant when customer characteristics \mathbf{X} are included. Straightforward use of a Gibbs sampler for this posterior distribution is very inefficient. This is the result of a strong dependence between the latent motivation assignments \mathbf{Z} on one hand and the parameters ϕ_m and θ_i on the other hand. A Gibbs sampler would therefore require an excessive number of draws to properly explore this posterior. Instead, we take advantage of the fact that the Dirichlet distribution is the conjugate prior for a categorical random variable. This allows us to marginalize over the ϕ_m and θ_i parameters, while retaining closed-form expressions for the conditional distributions of the other parameters in LDA(-X). By doing so we substantially improve

the mixing properties of the Gibbs sampler (Griffiths and Steyvers 2004). Hence, we examine the so-called *collapsed* posterior distribution of LDA(-X), defined as: $p(\mathbf{Z}, \beta_0, \gamma, \{\delta_i\}_{i=1}^M | \mathbf{Y}, \mathbf{X})$. The elements of \mathbf{Z} are sampled using a Gibbs sampler, while for the other parameters we implement a random walk Metropolis–Hastings sampler.

The set-up for inference in MDM is very similar to LDA(-X). The complete posterior distribution is given by: $p(\mathbf{s}, \{\varphi_i\}_{i=1}^I, \{\beta_i\}_{i=1}^M, \pi | \mathbf{Y})$. Again, we marginalize over the discrete distributions φ_i and π , resulting in a collapsed posterior distribution of MDM: $p(\mathbf{s}, \{\beta_i\}_{i=1}^M | \mathbf{Y})$. Here the segment allocations \mathbf{s} can be sampled in a Gibbs step, while the β_i parameters require a random walk Metropolis–Hastings sampler.

LDA(-X) and MDM are members of the general class of mixture models. This class of models is well known to be susceptible to end up in an area around a local maximum of the posterior distribution. As is common in this literature, this risk is reduced by using multiple random starts (Wedel and Kamakura 2000, Train 2009). For each value of M , we consider 250 different random starts. We reduce the computational burden of this approach by evaluating each random start at several intermediate steps of the estimation routine. At each step, we continue only with the best performing candidates. Performance is measured by the likelihood that results from the model's predictive distributions, averaged over purchases in a model-selection data set. This measure is closely related to the goal of predicting a new purchase as accurately as possible.

The same performance measure is also used to determine the number of motivations (for LDA(-X)) or segments (for MDM). In particular, for each model we increase the value of M until we find a decrease in the average predictive likelihood of the model-selection data.³ More details on the estimation routines, including pseudo-code, are provided in Online Appendix 1.

2.4. Comparison of LDA(-X) and MDM

Although the structures of LDA(-X) and MDM might appear quite similar at first, these models differ fundamentally on various grounds. In this subsection we first discuss this difference in terms of customer heterogeneity. Next, we consider the estimation complexity of the LDA(-X) and MDM models.

2.4.1. Heterogeneity Assumption. MDM assumes that heterogeneity in purchase behavior can be described by segmenting the customer base into groups of customers. Customers across segments are expected to be dissimilar, while customers in a segment are

expected to be rather similar. Hence, similarity between customers is mainly driven by segment membership. In LDA(-X) purchase behavior is described by motivations, where each motivation represents a preference for certain products in the assortment. Heterogeneity in purchase behavior is described by customer-specific probabilities for these motivations. This leads to a model where the purchases of a single customer are driven by *multiple* motivations. Here similarity between customers is motivation specific. Customers can have very similar purchase behavior for one set of products, corresponding to a shared motivation, and be very different for a set of products that belong to another motivation.

Which heterogeneity structure fits best depends on the specific situation. If customers typically have one or very few motivations, grouping customers in segments might be beneficial. If many combinations of motivations are present, the continuous mixture of motivations in LDA(-X) will be more parsimonious. Therefore, if a retailer has many different (latent) subcategories in its product assortment, and preferences across these subcategories vary independently across individuals, it is likely that the heterogeneity can be specified more parsimoniously by LDA(-X).

Although MDM assumes a hard clustering of customers into segments, one will use posterior segment probabilities to eventually make predictions. This will typically lead to a form of soft clustering, where a weighted combination of different segments is used. This brings the heterogeneity structure of MDM closer to that of LDA(-X). As we observe more purchases, the posterior segment probabilities in MDM will of course become more and more extreme. In the end this converges to strictly assigning a customer to a single segment.

2.4.2. Estimation Complexity. The different heterogeneity assumptions underlying LDA(-X) and MDM have a large impact on estimation complexity through the number of customer-specific parameters. In MDM each customer is given a distribution over the J -dimensional assortment, while in LDA(-X) a customer is described by a probability distribution over the M motivations, where M is generally much smaller than J . Even though we marginalize over these customer-specific distributions, this still affects the scalability of the models. Table 1 summarizes, for each model, the parameters that must be sampled to infer the model structure after marginalization. We differentiate between the sampling technique required, as Gibbs steps tend to be much faster and have better mixing properties than Metropolis–Hastings steps (Damien et al. 1999).

In LDA(-X) we need as many motivation allocations as purchases (N in total), whereas for MDM we only need to sample one segment allocation per customer

³ To validate this approach we also consider the models for larger values of M . The predictive performance stabilized at the values obtained with the selected value of M .

Table 1 Parameters to Sample in the MCMC Estimation Procedures Across Different Models

Model	Gibbs sampler		Metropolis–Hastings sampler	
	Parameters	No. parameters	Parameters	No. parameters
MDM	\mathbf{s}	I	$\{\beta_i\}_{i=1}^M$	$M \times J$
LDA	\mathbf{Z}	N	β_0, γ	$1 + M$
LDA-X	\mathbf{Z}	N	$\beta_0, \gamma, \{\delta_i\}_{i=1}^M$	$1 + M \times (1 + K)$

Notes. I , number of customers; M , number of segments/motivations; J , assortment size; K , number of predictor variables in \mathbf{x}_i ; N , total number of purchases.

(I in total). Although the number of allocations is larger in LDA(-X), this does not imply that the total allocation in LDA(-X) is computationally more demanding. The sampling step for each motivation assignment in LDA(-X) involves only elementary arithmetic operations, while for each segment allocation in MDM we have to evaluate complex Gamma functions. It is difficult to quantify the difference in computational complexity as it also depends on the (latent) structure in the data, but we anticipate that MDM will be slightly more complex for these Gibbs sampling steps.⁴

The remaining model parameters are sampled using Metropolis–Hastings steps, each of which is computationally demanding. For LDA we sample $1 + M$ parameters; for LDA-X this increases to $1 + M \times (1 + K)$ parameters. These numbers are in sharp contrast to MDM in which $M \times J$ parameters are sampled. This renders MDM much more demanding in terms of estimation time, as the assortment size J is large. This is the price that must be paid for the many degrees of freedom per customer. The number of Metropolis–Hastings steps in LDA(-X) is largely insensitive to the size of the assortment, the number of customers, and the number of purchases. In MDM, on the other hand, the number of Metropolis–Hastings steps linearly increases with the assortment size. This limits the scalability of MDM, which is why we can only extend LDA by including observed heterogeneity through \mathbf{x}_i .

2.5. Benchmark Methods

In this section we present the main ideas underlying the two benchmark methods to which we will compare the predictive performance of LDA(-X) and MDM. Details are available in Appendix A. The first benchmark is a collaborative filter; the second is built on standard discrete choice modeling.

2.5.1. Collaborative Filtering. A collaborative filter is a deterministic algorithm that predicts purchases by matching customers to each other based on purchase histories. There are many ways to implement a collaborative filter. Details of the actual implementations

used in industry are not common knowledge. Below we develop our own implementation of a collaborative filter.

Ideally, a focal customer is matched to customers who purchased the focal customer’s previously purchased products and at least one additional item. However, such a matching on the complete purchase history is in general not feasible due to the curse of dimensionality; the larger the purchase history, the less likely it matches other customers’ histories.

We alleviate this curse of dimensionality by matching on parts of the purchase history. First, for each customer i we replace the complete purchase history vector \mathbf{y}_i by the set of unique sorted subvectors of length k that can be created from \mathbf{y}_i . We denote this set of vectors by H_i^k . For example, for $k = 2$ a customer’s purchase history is replaced by all of the unique sorted pairs that can be formed using the purchase history, so $\mathbf{y}_i = [1, 1, 1, 2, 3]$ would be reduced to the set H_i^2 containing the pairs (1, 1), (1, 2), (1, 3), and (2, 3).⁵ Next, for each subvector in this set we match the focal customer against all customers. If k is relatively small, this will result in many more matches compared to a matching on the complete purchase history. This solves the curse of dimensionality problem at the cost of a loss of information.

Using these matches for a customer, we create individual-specific product scores that can be used to construct rankings over the product assortment. This ranking obviously depends on k . In our application we consider collaborative filters with two combination sizes, $k = 1$ and $k = 2$, denoted by CF-1 and CF-2, respectively. Using $k = 1$, customers are matched on the presence of single products in their purchase history. For $k = 2$ customers are matched on the presence of pairs of products in their purchase history. Larger product combinations are not desirable in our application, in terms of computational feasibility and the degree of sparseness in these larger combinations.

2.5.2. Discrete Choice Models. Random utility-based multinomial choice models (Maddala 1983, McFadden 1986) have been extensively used in marketing to model discrete choices from a set of given alternatives. Implementing a traditional discrete choice model that directly uses purchase history data from a large assortment as predictor variables, however, is infeasible. Such a model would have to predict purchases for J products based on J predictor variables, where each predictor variable describes whether a product was purchased by the customer in the past. This model specification would require simultaneous estimation of $J(J - 1)$ parameters, which is infeasible from a computational perspective and will also likely result in

⁴ More details on the required sampling steps can be found in Online Appendix 1.

⁵ Use of unique sorted pairs implies that (1, 1) occurs in H_i^2 only once and that H_i^2 contains the pair (1, 2) and not (2, 1).

identification issues due to sparse data. Hence, traditional discrete choice models do not scale well when the number of products J becomes large.

The benchmark discrete choice model that we propose resolves these problems by smartly constructing predictor variables, thus enabling a huge reduction in the number of parameters to be estimated. We build on the idea that in the binary logit model, the probability that customer i purchases product j is specified by

$$\Pr[y_{in} = j] = \frac{\exp(\theta_{ij})}{1 + \exp(\theta_{ij})}, \quad (6)$$

where θ_{ij} represents the log odds of customer i having purchased product j . Our goal is to capture the variation of θ_{ij} across customers and products through predictor variables, without introducing customer- or product-specific parameters.

In the simplest model we could relate θ_{ij} to a constant, the log of the number of products purchased by customer i , and the log of the observed odds ratio for product j . However, in this model the implied product ranking will not differ across customers. To obtain personalized rankings we need to replace the observed population odds by odds ratios that are specifically relevant to customer i .

To this end we first apply k -means clustering to the purchase history data to obtain clusters of customers. Using these clusters, we can define cluster-specific odds ratios for each product and add these as predictor variables to our model. To capture the heterogeneity across customers we also include interaction between the cluster-specific odds ratios and a measure of the similarity between customer i and the cluster. As in LDA(-X) and MDM we can vary the number of identified clusters. For similarity with these methods we denote the number of clusters in this method by M .

Finally, the model can straightforwardly be extended using customer characteristics. We label this model as DCM, short for discrete choice model. Technical details are provided in Appendix A.

2.6. Real-Time Online Predictions

For each of the prediction methods, it is straightforward to construct a product ranking over the assortment for each customer. In the context of online retailing it is important to continuously update this ranking based on the customer's new purchases. Re-estimating the (population-level) parameters can be done offline after a substantial amount of new data has been collected. However, updating the predictions for a specific customer should be feasible online. This allows the retailer to update predictions while customers select products during a shopping trip. For all methods, the real-time update step itself consists of simple arithmetic operations using the details provided in Online Appendix 2. A possible bottleneck could be the amount of data that must be available, retrieved, and processed to enable the updates. In the top half of Table 2 we display the number of elements needed to update a single customer's product ranking in real-time, for each new product purchase that is observed. The bottom half of the table provides information on the amount of data that must be stored for the entire customer base to enable the aforementioned real-time update step.

The first row in Table 2 mimics the context of our application, i.e., a medium-sized online retailer with an assortment of 500 products, 10,000 customers, and on average 10 purchases per customer. The number of segments/motivations/clusters (M) is set to 20, which is slightly larger than our empirical findings in this paper. We consider our implementation of a collaborative filter with combination size $k = 2$. In this

Table 2 Comparison of Memory Requirements for Real-Time Updating

No. of selected data elements for each real-time update step							
Retailer context							
I	J	n_i	M	LDA(-X)	MDM	CF-2	DCM
10,000	500	10	20	$1.00 \cdot 10^4$	$1.00 \cdot 10^4$	$5.51 \cdot 10^3$	$1.01 \cdot 10^4$
100,000	5,000	20	40	$2.00 \cdot 10^5$	$2.00 \cdot 10^5$	$1.05 \cdot 10^5$	$2.00 \cdot 10^5$
1,000,000	50,000	40	80	$4.00 \cdot 10^6$	$4.00 \cdot 10^6$	$2.05 \cdot 10^6$	$4.00 \cdot 10^6$
No. of stored data elements for the real-time update step							
Retailer context							
I	J	N/I	M	LDA(-X)	MDM	CF-2	DCM
10,000	500	10	20	$2.10 \cdot 10^5$	$3.10 \cdot 10^5$	$6.77 \cdot 10^7$	$1.10 \cdot 10^5$
100,000	5,000	20	40	$4.20 \cdot 10^6$	$6.20 \cdot 10^6$	$6.30 \cdot 10^{10}$	$2.20 \cdot 10^6$
1,000,000	50,000	40	80	$8.40 \cdot 10^7$	$1.24 \cdot 10^8$	$6.26 \cdot 10^{13}$	$4.40 \cdot 10^7$

Notes. I , number of customers; M , number of segments/motivations/clusters; J , assortment size; n_i , number of purchases made by customer i ; N , total number of purchases.

context, the number of elements that must be selected for the real-time update step is of the same order of magnitude across the prediction methods. The storage requirements, on the other hand, are of a different order of magnitude, i.e., millions for the collaborative filter versus thousands for the model-based approaches. However, for these settings all methods can easily be used in practice.

To illustrate the scalability of the various methods we increase the size of the assortment and customer base by a factor of 10 and we double the average purchase history size and M . Naturally, all memory requirements increase in this setting, but the rate of growth differs significantly. For the collaborative filter the storage requirements increase approximately by a factor of 1,000, while the model-based approaches only increase by a factor of 20. The same holds for the third context, in which we again increase the dimensions. This illustrates that the dimension reduction achieved by the model-based approaches ensures that they are suitable for real-time predictions in large scale applications, even if the number of underlying dimensions grows with the amount of available data. In addition, it is infeasible to use a combination size larger than $k = 2$ in our implementation of a collaborative filter, as in that case the storage requirements would increase even faster. For very large applications, one might even need to rely on the simpler CF-1, which only matches purchase histories on the presence of single products.⁶

2.7. Performance Measures

To evaluate the methods for a range of different customization applications, we consider *prediction sets* of different sizes. A prediction set of size S contains the S highest ranked products for a customer. To recommend a single product, the prediction set of size 1 is most relevant. However, when customizing a page with search results the prediction set of size 10 may be more relevant. We assess the quality of a prediction set by matching its contents against hold-out purchase data. These purchases are denoted by \mathbf{y}'_i for customer i ; the number of unique purchased products in \mathbf{y}'_i is given by u'_i .

We denote a complete ranking of all J products for customer i by the vector \mathbf{r}_i . The first element, r_{i1} , is the product with the highest predicted purchase probability for the model-based rankings, the highest product score for the collaborative filters, and the highest odds for DCM. The quality of a prediction set of size S can be measured by the number of products in the prediction set that overlap with the hold-out purchases: $\sum_{s=1}^S \mathbb{I}[r_{is} \in \mathbf{y}'_i]$. This number should be seen relative to the maximum number of hits possible to obtain a hit

rate that may be compared across prediction sets of different sizes. This maximum is bounded by S , the size of the prediction set, and the number of unique hold-out purchases u'_i . Hence, the hit rate for customer i could be defined as: $\sum_{s=1}^S \mathbb{I}[r_{is} \in \mathbf{y}'_i] / \min(S, u'_i)$.

If a prediction set is presented to a customer in an application, such as a recommendation list, the positions within the set are also of importance (Xu and Kim 2008). We incorporate this notion in our hit rate by weighing the hits according to their ranks. For the s th ranked product in a prediction set of size S this weight is specified as: $w(s, S) = 1 - (s - 1)/S$. Combining the above, we obtain our final performance measure, the weighted hit rate

$$h_i(\mathbf{r}_i, S) = \frac{\sum_{s=1}^S \mathbb{I}[r_{is} \in \mathbf{y}'_i] w(s, S)}{\sum_{s=1}^{\min(S, u'_i)} w(s, S)}. \quad (7)$$

3. Data

We apply the prediction methods to purchase data from a medium-sized online retailer in the Netherlands.⁷ The data starts at the launch of the retailing platform and covers a period of approximately 67 weeks. The product assortment primarily consists of nonfood fast-moving consumer goods, such as detergents, deodorants, and shampoo. The assortment is complemented with a small selection of high turnover products for infants and toddlers, such as diapers and baby food. Consequently, the data contains many repeat purchases.

Initially, the data contains 3,226 unique product IDs. These IDs correspond to a very fine-grained classification, e.g., different package sizes of the same product each receive a unique ID. We opt for a more coarse-grained classification and combine products on the category-brand level. For example, different fragrances of the same deodorant brand are aggregated to one category-brand combination. This approach results in a total of 440 unique category-brand combinations. Additionally, this aggregation step is applied to the customer orders: If an order contains multiple products from the same category-brand, we consider this a single purchase from this category-brand. Finally, the category-brands that are purchased five times or fewer across all purchases are removed from the data. Below we refer to the category-brand combinations as “products.” After the aggregation steps the data contains 95,208 product purchases of 394 products made by 11,783 distinct customers.

We chronologically split the data in two parts: The first 80% of the purchases are used as in-sample data; the hold-out data comprises the last 20% of the purchases. We use the hold-out data to assess the predictive

⁶ In our application, this simpler collaborative filter performs systematically worse than CF-2.

⁷ The authors wish to thank Christian van Someren, former Managing Director of Truus.nl, for kindly providing us this data.

Table 3 Characteristics of the Three Subsets of the Purchase Data

Subset	Customers	Unique products	Purchases
Full data	11,783	394	95,208
Estimation data	8,831	393	71,346
Model-selection data	4,820	323	4,820
Hold-out data	3,745	369	19,042

performance of the methods. This division mimics the setting of predicting future purchase behavior. Subsequently, we split the in-sample data into an estimation and a model-selection subset. We randomly select half of the customers from the in-sample data; for each of these customers, a single product purchase is randomly selected as model-selection data. We use the remaining in-sample data to estimate LDA(-X), MDM, and DCM, and to create the collaborative filters. The model-selection data is used to determine the number of motivations/segments/clusters (M) in LDA(-X), MDM, and DCM, respectively. Table 3 summarizes the three data subsets in terms of number of customers, unique products, and number of product purchases.

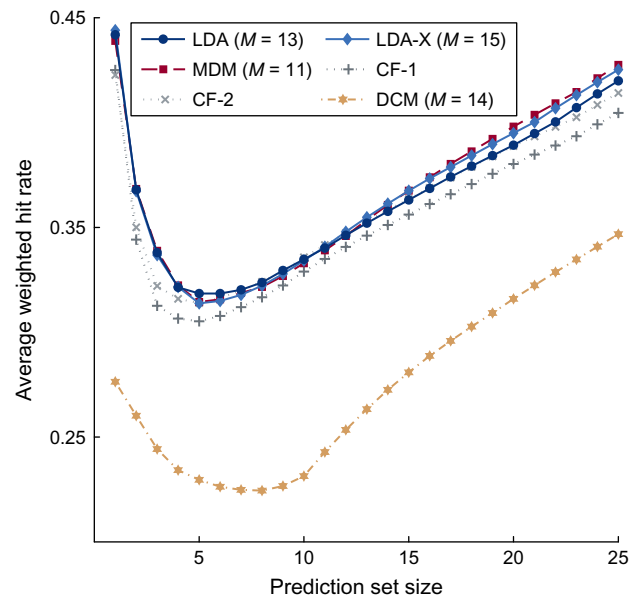
It is likely that the type of customer acquired by the retailer changes over time, for example due to (a shift in) brand awareness or the mix of advertising channels used. Therefore, we investigate whether the customer's time of adoption at the retailer systematically shifts customer preferences. We define the time of adoption as the number of days between a customer's first order, and the starting date of the retailing platform. We take the natural logarithm of this variable to allow for larger shifts in the preferences of customers acquired during the early stages of the retailing platform. Finally, this variable is standardized using the mean and variance of the in-sample data.

4. Predictive Performance

In this section we report on the predictive performance of the methods considered in this paper. First, for LDA(-X), MDM, and DCM we determine M , the number of motivations, segments, and clusters, respectively. For each model we select M using the model's average predictive likelihood for the model-selection data. For details refer to Appendix B. We select $M = 13$ for LDA, $M = 15$ for LDA-X, $M = 11$ for MDM, and $M = 14$ for DCM.

To assess a method's predictive performance we evaluate its weighted hit rate for the hold-out data, see (7). In the weighted hit rate, each hit receives a weight that depends on the rank assigned to the prediction. A better (numerical lower) rank receives a larger weight than a worse (numerical higher) rank.

Figure 1 presents the average weighted hit rate for each method, obtained by averaging the individual-specific hit rates across all customers in the hold-out

Figure 1 (Color online) Predictive Performance for the Complete Hold-Out Data, as a Function of Prediction Set Size

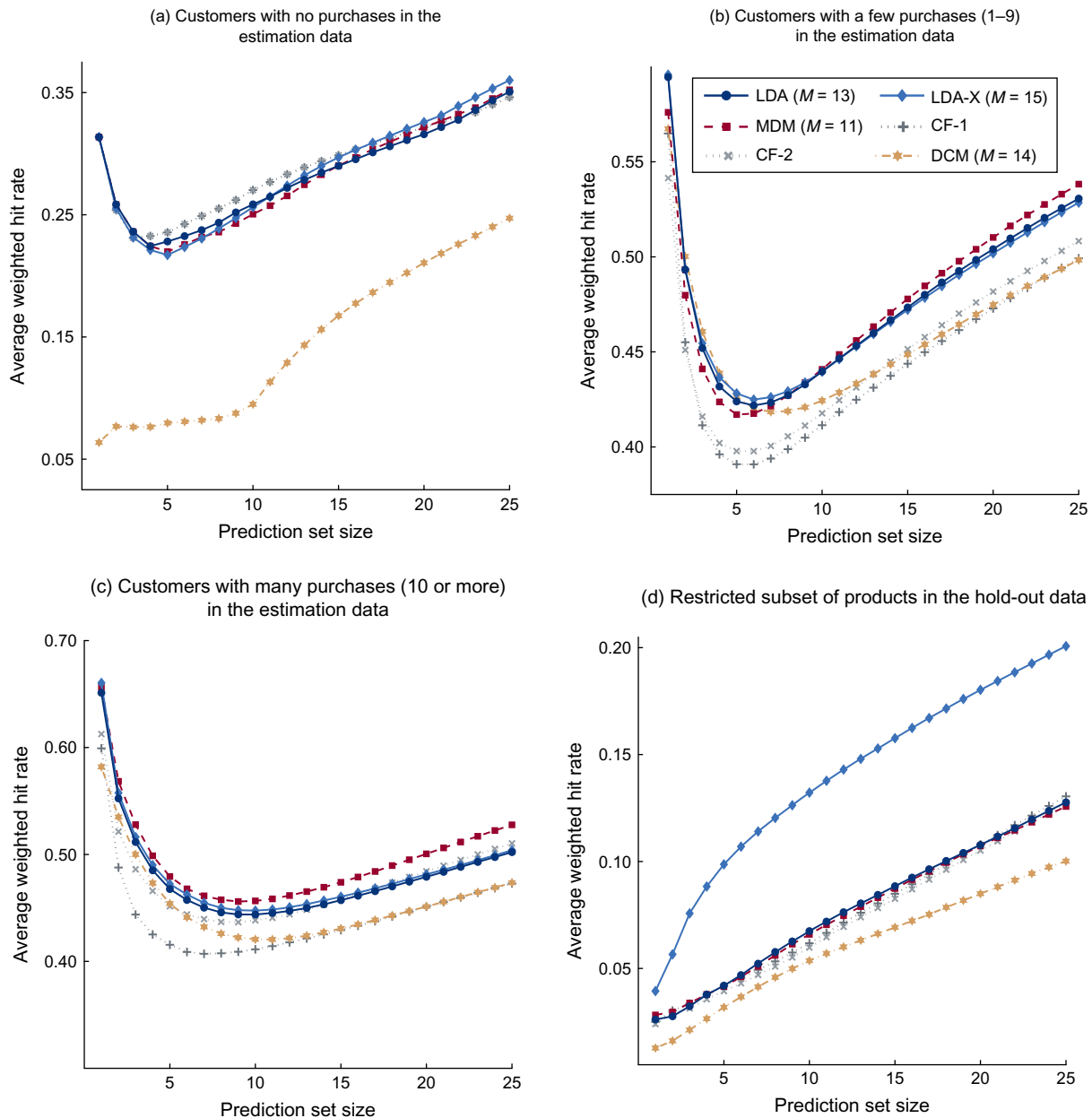
data. In case we predict only a single product for each customer, i.e., a prediction set of size one, LDA-X has the best performance with a hit rate close to 0.45. For most prediction set sizes, LDA(-X) and MDM outperform the collaborative filters. The best performing collaborative filter is CF-2, which matches customers on the presence of pairs of products in their purchase history. Given the decent predictive likelihoods generated by DCM (see Appendix B), its performance in terms of ranking the purchased products is unexpectedly poor.

Note that the average hit rate declines for the first few prediction set sizes. This is a direct consequence of the denominator in the definition of the hit rate in (7), which divides the total number of hits by the maximum number of hits possible for a given customer and prediction set size. This number increases with the size of the prediction set until it reaches the number of unique products purchased by the customer. As the average number of unique purchases per customer in the hold-out data is almost 5, the hit rates increase beyond that value for most methods.

We study the difference in performance for the prediction methods in more detail by separately considering specific groups of customers and products. In particular, we first divide the customers in the hold-out data into three groups based on the number of purchases in the estimation data: (i) 2,185 customers with no prior observed purchases (Figure 2(a)); (ii) 809 customers with a moderate amount (1–9) of purchases (Figure 2(b)); and (iii) 751 customers with many (10 or more) purchases (Figure 2(c)).

The most apparent performance difference between these groups is visible in the range of the y -axis. If we observe many purchases for a customer the average hit

Figure 2 (Color online) Predictive Performance for Different Groups of Customers/Products in the Hold-Out Data



rates are twice as large for the smaller prediction sets, compared to those for customers with no purchases in the estimation data. This is exactly according to our expectations, and provides empirical evidence that purchase history data is very informative about a customer's future purchases.

By examining Figure 2(a) we see that for customers without previous purchases the collaborative filters perform very well (particularly for moderate-sized prediction sets). Note that for this specific group of customers the collaborative filters rank the products according to their market penetration in the customer base. Also for LDA and MDM there is no information that can be used to make a personalized prediction.

LDA-X uses the time of adoption, although this does not seem to shift the baseline predictions a lot. Hence, the performance differences between LDA(-X) and MDM are small.

In the absence of a purchase history, the similarity of a customer to each of the M clusters, used to create predictor variables in DCM, is meaningless. As a result, the DCM's predictive power is low for these customers. In fact, a large part of the performance gap on the complete hold-out data between DCM and the other prediction methods is driven by the poor performance for the group of customers without a purchase history.

We observe a different pattern for customers with a moderate number of past purchases in Figure 2(b),

where LDA(-X) and MDM consistently outperform the collaborative filters. This indicates that these model-based methods are better able to learn from a customer's previous purchases than the collaborative filters. Comparing the methods, LDA(-X) attains the highest overall performance and performs best when we predict only a single product; MDM performs better for larger prediction sets. The performance of DCM is competitive for the smaller prediction sets, although its relative performance drops substantially for larger prediction set sizes.

The final group of customers that we consider are those who made many purchases, displayed in Figure 2(c). The general conclusion is similar to that of the customers with a moderate number of purchases. However, in this case MDM obtains the highest performance for prediction sets with more than one product. This result, combined with the previous findings, may be explained by the flexibility of the customer-level heterogeneity structure. In MDM preferences are modeled by a customer-specific probability vector over the product assortment. On the other hand, in LDA(-X) a customer's individual preferences are described by a lower-dimensional probability vector over the M motivations. Both models learn from previous purchases, but in MDM this learning is directly incorporated in the preferences over the assortment, while in LDA(-X) it is done indirectly through the probabilities for the motivations. Consequently, MDM has more degrees of freedom at the level of the individual customer as the assortment size J is much larger than the number of motivations M . This additional flexibility pays off when many purchases are observed for a customer.

The results above highlight the performance of the methods for the complete assortment. However, many of the highly ranked products are frequently purchased or products that have been previously purchased by the focal customer. Customers can easily anticipate such recommendations and might even be bored by them (Fleder and Hosanagar 2009). It is therefore interesting to evaluate the performance of the methods when predicting products that may be more *unexpected*.

To assess the performance of the methods for predicting such unexpected products, we evaluate the predictive performance for a restricted subset of the product assortment. This subset is constructed as follows: First, we remove 20% of the products in the assortment that are most frequently purchased in the estimation data. Second, we create a customer-specific restriction by removing the products that have previously been purchased by this customer. Subsequently, for each customer, we only consider the predictions and hold-out purchases for products in this restricted subset of the assortment. As customers are less likely to be aware of these products, performing well on this

aspect could increase the *cross-selling* performance of marketing actions that are based on such predictions.

The predictive performance for the restricted set of products is displayed in Figure 2(d). LDA and MDM perform better than the collaborative filters and DCM, but LDA-X clearly outperforms all of the other prediction methods. This remarkable performance difference primarily arises for the highly ranked products. By examining these products, we find that the product "Slimming nutrition—Weight Care" appears in the top of many of the LDA-X customer-level prediction sets. The prediction sets resulting from the other methods, however, do not contain this product. In fact, "Slimming nutrition—Weight Care" is the most frequently purchased product in the hold-out data. Its purchase frequency has shifted from 0.04% in the estimation data to 4.88% in the hold-out data. LDA-X captures this shift through the time of adoption variable.⁸ This shows that the inclusion of predictor variables has merit in the context of purchase prediction, even though the time of adoption variable in general does not add much explanatory power. The reason that we do not see a similar shift for DCM can be explained by the way the predictor variables enter the model. In LDA-X, it directly influences the likelihood of a certain motivation, in effect boosting a motivation that is relevant for customers who adopted later in time. In this case, it boosts the motivation with products that are purchased more frequently later in the observation period, including the period of the hold-out predictions. By contrast, in DCM the clusters are determined "outside" the model, using the k -means algorithm. The performance of the clustering algorithm does not benefit from selecting a cluster that is linked to the other prediction variables, as the predictor variables are not included when constructing the clusters. In the absence of such clusters of customers, inclusion of the predictor variables cannot shift the importance of these products, as they are not in a separate cluster.

5. Conclusion

In this paper we have evaluated several methods for purchase prediction in large product assortments using purchase history data. Inspired by the text modeling literature, we have introduced a novel model-based approach that uses latent Dirichlet allocation (LDA(-X)) to predict purchases. In addition, we have considered mixtures of Dirichlet-Multinomials (MDM), a framework well known in the brand-choice modeling literature. The performance of these model-based approaches has been contrasted against two benchmarks, i.e., a set of count-based collaborative filters,

⁸ We acknowledge that there can be many external influences that drive this shift in purchase behavior. Our predictor variable (time of adoption) likely serves as a proxy for the actual causes.

in which customers are matched on the contents of their purchase history, and a scalable implementation of a discrete choice model (DCM), that does not break down when used with a large product assortment. All methods can construct customer-specific product rankings over the assortment that can be used for purchase prediction.

Naturally, the prediction methods differ in their heterogeneity assumptions, estimation complexity, and memory requirements. In MDM purchase heterogeneity is specified at the customer level by segmenting the customer base. In LDA(-X), on the other hand, this heterogeneity is specified at the motivation level, which groups products, not customers. These heterogeneity assumptions also affect the estimation complexity of the models. MDM has more flexibility to model a customer's purchase behavior than LDA(-X), but this comes at the price of increased estimation complexity as more parameters must be estimated. The estimation complexity of the logit part of the DCM is relatively low, but it does depend on customer clusters from an external method (i.e., the k -means algorithm). The collaborative filter has as advantage that no (latent) model structure must be estimated, but its storage requirements for generating real-time online predictions rapidly increase for large applications. By contrast, the model-based approaches require less storage and this grows much slower with the size of the application.

The performance of the methods was assessed based on purchase prediction sets derived from the product rankings, and comparing these sets to actual hold-out purchases. In general, LDA(-X) and MDM perform best and, even though these two models are conceptually different, their predictive performance is comparable. In addition, we have considered the setting where we focus on the predictive performance for products in the tail of the assortment that have not yet been purchased by the customer. In this case LDA-X clearly outperforms the other methods, which can be attributed to the time of adoption variable included in LDA-X. Although DCM also includes this predictor variable, its dependence on the k -means algorithm prevents it from effectively using the additional information to generate better predictions.

In summary the LDA(-X) prediction method that we have introduced in this paper is the most promising approach to purchase prediction, particularly in the context of large online retailers. Its predictive performance is very competitive compared with the other methods and it scales well with the size of the application. Finally, it is a self-contained prediction method that can readily accommodate additional information available to the retailer. In our application we only had access to a weak predictor, but the potential benefits of including stronger predictors of customer preferences into the model could be great.

To conclude, LDA(-X) can be readily used as a stepping stone for further model-based research that quantifies and optimizes the impact of marketing interventions in large scale retailing environments. For example, one could optimize a recommendation system that differentiates between the likelihood of purchasing a product and the added benefits from recommending that product (Bodapati 2008, Wagner and Taudes 1986); this is difficult to implement in a count-based method such as a collaborative filter. Such extensions are an interesting avenue for further research.

Supplemental Material

Supplemental material to this paper is available at <http://dx.doi.org/10.1287/mksc.2016.0985>.

Acknowledgments

The authors thank the review team for their suggestions and remarks that helped us to improve this paper. The computations for this work were carried out on the Dutch national e-infrastructure with the support of SURF Cooperative.

Appendix A: Technical Details for the Benchmark Methods

A.1. Collaborative Filtering

We refer to a subvector of a customer's purchase history as a product combination, denoted by \mathbf{h} , and $c(\mathbf{h})$ gives the number of customers who purchased product combination \mathbf{h} , that is

$$c(\mathbf{h}) = \sum_{i=1}^I \mathbb{I}[\mathbf{h} \in H_i^{\dim(\mathbf{h})}], \quad (\text{A1})$$

where $\dim(\mathbf{h})$ denotes the dimension of \mathbf{h} and $\mathbb{I}[A]$ equals 1 if condition A is true and 0 otherwise. To obtain purchase predictions for customer i , using product combinations of size k , we score all products in the assortment based on their co-occurrence with each of the product combinations in H_i^k . For product combination $\mathbf{h} \in H_i^k$ the prediction score for product j equals the number of customers who purchased j and the products in \mathbf{h} , normalized by the sum of the score for \mathbf{h} and any product $p = 1, \dots, J$. This normalization ensures that each product combination $\mathbf{h} \in H_i^k$ receives the same weight, independent of the prevalence of \mathbf{h} in other customers' purchase histories. The final product score is the sum of the normalized scores across all $\mathbf{h} \in H_i^k$. Formally, for combination size k , the overall score of product j for customer i equals

$$s_{ij}^k = \sum_{\mathbf{h} \in H_i^k} \frac{c(\langle \mathbf{h}, j \rangle)}{\sum_{p=1}^J c(\langle \mathbf{h}, p \rangle)}, \quad (\text{A2})$$

where the arguments between angle brackets represent a single product combination of size $k+1$.⁹ Hence, to obtain product scores s_{ij}^k , by matching customers based on purchase histories reduced to combinations of size k , we need

⁹ For $k > 0$ it is possible that a product combination \mathbf{h} is never purchased with another product, i.e., for all p we have $\sum_{p=1}^J c(\langle \mathbf{h}, p \rangle) = 0$ in (A2). If a customer's purchase history contains such a combination, we regress to a lower value of k for this customer.

the summary of all purchase histories reduced to product combinations of size $k + 1$. So, matching customers on pairs of products requires counts over triplets of products as input for the purchase predictions. The product ranking for each customer is constructed by sorting the products on the product score defined above.¹⁰

A.2. Discrete Choice Models

In the binary logit model, the probability that customer i purchases product j is specified by

$$\Pr[y_{in} = j] = \frac{\exp(\theta_{ij})}{1 + \exp(\theta_{ij})}. \quad (\text{A3})$$

Here, θ_{ij} represents the log odds of having purchased product j . Ignoring heterogeneity among customers for the moment, these odds will largely be driven by the log of the number of (unique) products purchased by customer i , denoted by u_i , and the relative attractiveness of product j . We capture the relative attractiveness of product j using the log odds of this product based on the *observed* product-purchase frequencies in the purchase data at the customer-base level. This leads to the following expression for the log odds of customer i buying product j :

$$\theta_{ij} = \alpha + \beta \log(u_i) + \gamma \log(\text{odds}_j). \quad (\text{A4})$$

The product ranking resulting from this specification will be the same for all customers as the product attractiveness is defined at the customer-base level, not the customer level. To obtain predictions that differ across customers, we introduce heterogeneity in the model. To achieve this without resorting to a model with unobserved heterogeneity, as in LDA(-X) or in MDM, or requiring an excessive number of parameters, as in a regular choice model implementation, we construct variables at the customer-product level that characterize the attractiveness of product j for customer i using the available purchase history data.

The first step is to characterize customers based on their purchase history. We describe each customer's purchases by \mathbf{v}_i , a J -dimensional vector containing the proportions of each product in the customer's purchase history, with $\sum_{p=1}^J v_{ip} = 1$.¹¹ We then perform k -means clustering on these proportion vectors using M clusters. Customer heterogeneity can now be characterized by a customer's similarity with respect to each of the cluster means. We define the similarity of customer i with cluster m as

$$w_{im} = \frac{1}{1 + \|\mathbf{v}_i - \bar{\mathbf{v}}^{(m)}\|},$$

where $\|\mathbf{v}_i - \bar{\mathbf{v}}^{(m)}\|$ measures the Euclidean distance between customer i 's proportion vector and the m th cluster mean $\bar{\mathbf{v}}^{(m)}$.

¹⁰ In the rare case that two or more products receive the same score, they are ranked according to their order in the data set, which is alphabetical.

¹¹ For smoothing purposes we add one pseudo observation to each customer's purchase history that is equal to the relative market shares of each product.

We can now parsimoniously introduce customer-level heterogeneity by combining the cluster-level product attractiveness and the similarity measures w_{im} that capture the relevance of each cluster for each customer

$$\theta_{ij} = \alpha + \beta \log(u_i) + \sum_{m=1}^M \log(o_{mj})(\gamma_{1m} + \gamma_{2m}w_{im}), \quad (\text{A5})$$

where o_{mj} denotes the odds for product j that corresponds to the purchase proportions in cluster mean $\bar{\mathbf{v}}^{(m)}$. Note that in this model specification, the parameters are not product specific, as the relative attractiveness of each product is captured through the summary of the purchase behavior of the various clusters.¹²

Maximum likelihood estimation of this parsimonious discrete choice model (DCM) is relatively straightforward and including the other available predictor variables is therefore feasible. To do so, we extend the specification in (A5) to include interactions with the customer-specific predictor variables in \mathbf{x}_i , resulting in

$$\theta_{ij} = \alpha + \beta \log(u_i) + \sum_{m=1}^M \log(o_{mj}) \cdot \left(\gamma_{1m} + \gamma_{2m}w_{im} + \sum_{k=1}^K x_{ik}(\delta_{1km} + \delta_{2km}w_{im}) \right). \quad (\text{A6})$$

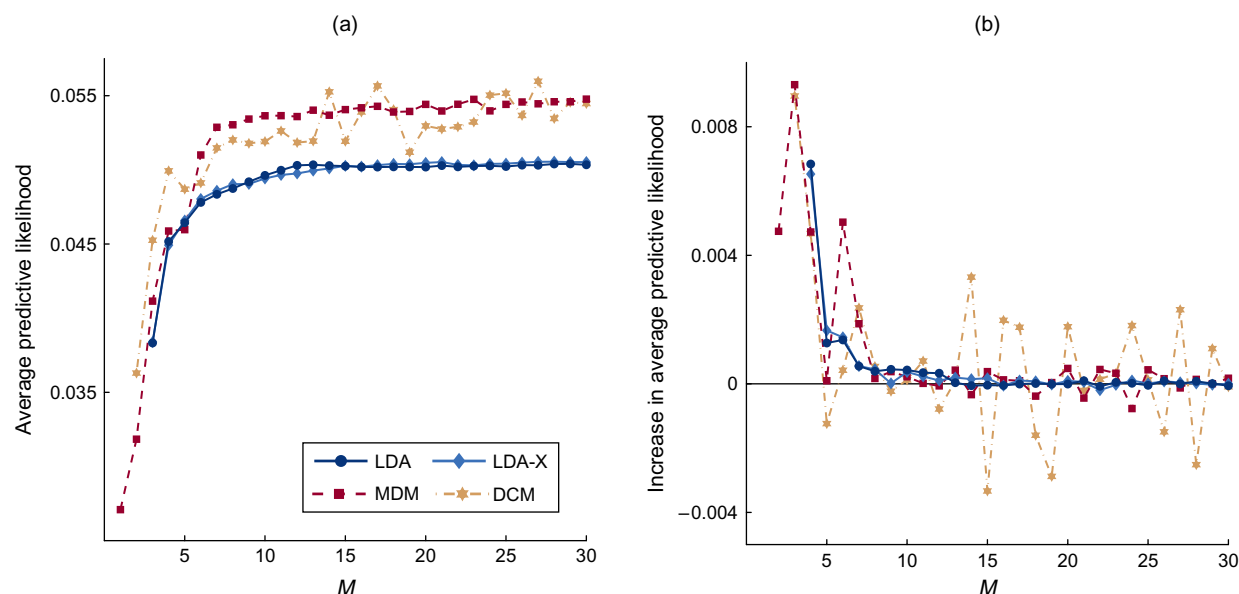
Appendix B: Model Selection

In all model-based approaches we must determine M the number of motivations, segments, and clusters. We evaluate LDA(-X) for $M = 3, \dots, 30$ and MDM for $M = 1, \dots, 30$, where MDM with $M = 1$ corresponds to the DM model. For each of these model configurations (choice of model plus a value of M) we use 250 different random starts to avoid local maxima. Throughout the estimation procedure the performance of each random start is measured by the average predictive likelihood for the model-selection data. As discussed in §2.3, at several points during the procedure we drop the worst-performing starting values (see Online Appendix 1). At the end of the estimation routine we use the parameter estimates that result from the random start with the highest average predictive likelihood. We evaluate DCM for $M = 2, \dots, 30$. To avoid local maxima in the k -means algorithm used in DCM, we use 1,000 different random cluster initializations. For each value of M , the clustering that obtains the lowest in cluster sum-of-squares is selected.

The average predictive likelihoods for the model-based approaches are shown in Figure B.1(a). We find that for each method the average predictive likelihood steeply increases for the first few values of M and then levels off for larger values of M . This result indicates that choosing M too small likely impedes performance more than choosing M too large. The average predictive likelihood of LDA and LDA-X is similar, reaching a value of approximately 0.05 for the larger values of M . MDM performs slightly better, reaching a value close to 0.055. DCM performs similarly and in between LDA(-X) and MDM, although its performance

¹² Model specifications where the coefficients were allowed to be product specific suffered from severe identification problems in our application as the number of parameters is increased by a factor J .

Figure B.1 (Color online) Average Predictive Likelihood for the Model-Selection Data as a Function of M



fluctuates across values of M . Note that the average predictive likelihood is merely an indicator for the actual predictive performance in our application, as we will consider the rank assigned to purchased products to evaluate the predictive performance and not actual purchase likelihoods.

To determine the number of motivations and segments in LDA(-X) and MDM, we select the first value of M for which the average predictive likelihood decreases when M is increased by 1, i.e., we select the first local maximum. As the graphs in Figure B.1(a) stabilize after their first local maximum, this approach results in a parsimonious, yet high performing model specification. Figure B.1(b) shows the differences in performance between subsequent values of M . The first negative value, corresponding to a decrease in performance, is obtained at $M = 14$ for LDA, $M = 16$ for LDA-X, and $M = 12$ for MDM. Hence, we select $M = 13$ for LDA, $M = 15$ for LDA-X, and $M = 11$ for MDM. The average predictive likelihood is more volatile across values of M for DCM, resulting in the first local maximum for $M = 4$. In the spirit of our M selection criterion for LDA(-X) and MDM, we instead select the smallest value of M that corresponds to a local maximum in the range of the values of M where the predictive likelihood has leveled off. For DCM, this happens at $M = 14$.

References

- Asuncion A, Welling M, Smyth P, Teh YW (2009) On smoothing and inference for topic models. Bilmes J, Ng A, eds. *Proc. Twenty-Fifth Conf. Uncertainty Artificial Intelligence* (AUAI Press, Arlington, VA), 27–34.
- Blei DM, Ng AY, Jordan MI (2003) Latent Dirichlet allocation. *J. Machine Learn. Res.* 3:993–1022.
- Bodapati AV (2008) Recommendation systems with purchase data. *J. Marketing Res.* 45(1):77–93.
- Damien P, Wakefield J, Walker S (1999) Gibbs sampling for Bayesian non-conjugate and hierarchical models by using auxiliary variables. *J. Roy. Statist. Soc. Ser. B (Statist. Methodology)* 61(2): 331–344.
- Fader PS (1993) Integrating the Dirichlet-multinomial and multinomial logit models of brand choice. *Marketing Lett.* 4(2):99–112.
- Fader PS, Hardie BGS (1996) Modeling consumer choice among SKUs. *J. Marketing Res.* 33(4):442–452.
- Fader PS, Schmittlein DC (1993) Excess behavioral loyalty for high-share brands: Deviations from the Dirichlet model for repeat purchasing. *J. Marketing Res.* 30(4):478–493.
- Fleder D, Hosanagar K (2009) Blockbuster culture's next rise or fall: The impact of recommender systems on sales diversity. *Management Sci.* 55(5):697–712.
- Ghose A, Ipeirotis PG, Li B (2014) Examining the impact of ranking on consumer behavior and search engine revenue. *Management Sci.* 60(7):1632–1654.
- Goodhardt GJ, Ehrenberg ASC, Chatfield C (1984) The Dirichlet: A comprehensive model of buying behaviour. *J. Roy. Statist. Soc. Ser. A (General)* 147(5):621–655.
- Griffiths TL, Steyvers M (2004) Finding scientific topics. *Proc. Natl. Acad. Sci. USA* 101(Suppl 1):5228–5235.
- Grover R, Srinivasan V (1987) A simultaneous approach to market segmentation and market structuring. *J. Marketing Res.* 24(2): 139–153.
- Guadagni PM, Little JDC (1983) A logit model of brand choice calibrated on scanner data. *Marketing Sci.* 2(3):203–238.
- Jain D, Bass FM, Chen Y-M (1990) Estimation of latent class models with heterogeneous choice probabilities: An application to market structuring. *J. Marketing Res.* 27(1):94–101.
- Jannach D, Zanker M, Felfernig A, Friedrich G (2011) *Recommender Systems: An Introduction* (Cambridge University Press, New York).
- Jeuland AP, Bass FM, Wright GP (1980) A multibrand stochastic model compounding heterogeneous Erlang timing and multinomial choice processes. *Oper. Res.* 28(2):255–277.
- Liu D-R, Lai C-H, Lee W-J (2009) A hybrid of sequential rules and collaborative filtering for product recommendation. *Inform. Sci.* 179(20):3505–3519.
- Maddala GS (1983) *Limited-Dependent and Qualitative Variables in Econometrics* (Cambridge University Press, New York).
- McFadden D (1986) The choice theory approach to market research. *Marketing Sci.* 5(4):275–297.
- Mimno D, Hoffman MD, Blei DM (2012) Sparse stochastic inference for latent Dirichlet allocation. Langford J, Pineau J, eds. *Proc. Twenty-Ninth Internat. Conf. Machine Learn.* (ACM, New York), 1599–1606.
- Naik P, Wedel M, Bacon L, Bodapati A, Bradlow E, Kamakura W, Kreulen J, Lenk P, Madigan DM, Montgomery A (2008)

- Challenges and opportunities in high-dimensional choice data analyses. *Marketing Lett.* 19(3–4):201–213.
- Pan B, Hembrooke H, Joachims T, Lorigo L, Gay G, Granka L (2007) In Google we trust: Users decisions on rank, position, and relevance. *J. Comput.-Mediated Comm.* 12(3):801–823.
- Ramage D, Dumais S, Liebling D (2010) Characterizing microblogs with topic models. Cohen W, Gosling S, eds. *Proc. Fourth Internat. AAAI Conf. Weblogs Soc. Media* (AAAI Press, Menlo Park, CA), 130–137.
- Salganik MJ, Dodds PS, Watts DJ (2006) Experimental study of inequality and unpredictability in an artificial cultural market. *Science* 311(5762):854–856.
- Steyvers M, Griffiths T (2007) Probabilistic topic models. Landauer TK, McNamara DS, Dennis S, Kintsch W, eds. *Handbook of Latent Semantic Analysis* (Psychology Press, London), 424–440.
- Train KE (2009) *Discrete Choice Methods with Simulation*, 2nd ed. (Cambridge University Press, New York).
- Wagner U, Taubes A (1986) A multivariate Polya model of brand choice and purchase incidence. *Marketing Sci.* 5(3):219–244.
- Wallach HM, Mimno D, McCallum A (2009) Rethinking LDA: Why priors matter. Bengio Y, Schuurmans D, Lafferty JD, Williams CKI, Culotta A, eds. *Advances in Neural Information Processing Systems*, Vol 22 (Curran Associates, Red Hook, NY), 1973–1981.
- Wedel M, Kamakura WA (2000) *Market Segmentation: Conceptual and Methodological Foundations*, 2nd ed. (Kluwer Academic Publishers, Norwell, MA).
- Xu Y(C), Kim H-W (2008) Order effect and vendor inspection in online comparison shopping. *J. Retailing* 84(4):477–486.
- Zanutto EL, Bradlow ET (2006) Data pruning in consumer choice models. *Quant. Marketing Econom.* 4(3):267–287.