

E-Customization

Asim Ansari

Columbia University

Carl F. Mela*

Duke University

December 2000

Revised September 2001

Revised June 2002

* Asim Ansari (email: maa48@columbia.edu, telephone: 212-854-3476, fax: 212-854-7647) is Associate Professor at the Columbia Business School, 517 Uris Hall, 3022 Broadway, New York, NY 10027. Carl F. Mela (email: mela@duke.edu, telephone: 919-660-7767, fax: 919-681-6245) is Associate Professor of Marketing, The Fuqua School of Business, Duke University, Durham, North Carolina, 27708. Authors are listed alphabetically. The authors would like to thank Bill Boulding, Andrew Gershoff, Don Lehman, Kamel Jedidi, John Lynch, Richard Staelin and Seenu Srinivasan and the seminar participants at Columbia and Duke Universities for their helpful comments. Special thanks to Rajeev Kohli for many insightful discussions and help on the optimization.

ABSTRACT
E-customization

Customized communications have the potential to reduce information overload and aid customer decisions, and the highly relevant products that result from customization can form the cornerstone of enduring customer relationships. In spite of such potential benefits, few models exist in the marketing literature to exploit the Internet's unique ability to design communications or marketing programs at the individual level. We develop a statistical and optimization approach for customization of information on the Internet. We use clickstream data from users at one of the top ten most tracked websites to estimate the model and to optimize the design and content of such communications for each individual. Our model is applied to the context of permission based email marketing, where the objective is to customize the design and content of the email in order to increase website traffic. Our analysis suggests that our content targeting approach can potentially increase the expected number of click-throughs by 62%.

Keywords: Internet Marketing, Mass-customization, Targeting, Information Design, Hierarchical Bayes, Dirichlet Process Priors.

INTRODUCTION

Marketers have long realized the value of targeting and customization. Customized products and communications attract customer attention and foster customer loyalty and lock-in. Targeted communications aid customer decisions and reduce information overload, whereas highly relevant products yield satisfied customers. The customer loyalty that results from such personalization and targeting can translate into increased cash-inflows and enhanced profitability. Customized marketing solutions are useful for both customer acquisition and retention and can engender successful, long-term relationships. However, customization has often proven difficult due to implementation challenges, insufficient customer information and other factors.

The web makes mass-customization eminently possible. Firms can exploit the capabilities afforded by digitization and networking to provide unique content of direct relevance to each customer. Moreover, such tailoring of information can be done quickly and at low cost. Customization is possible in part because of the interactivity afforded by the Web. Firms can collect and update preference information of customers from onsite surveys and from the traces that customers leave as they navigate through a web site. This knowledge can then be seamlessly integrated with algorithms and software to automatically customize content for individuals. Indeed, customized design (serving different variants of content to different users at different points in time) represents one of the key features that differentiates the Web from more “traditional” media.

In the online world, content sites such as C-net and Yahoo can leverage a loyal customer base to increase readership and therefore advertising revenues.¹ Various surveys have extolled

¹ Jupiter Communications forecasts online advertising spending will reach \$11.5 billion by 2003, surpassing spending in traditional media.

the rapid growth in advertising dollars on the Web. Given the large magnitude of expected revenue involved, content providers are increasingly turning towards customization strategies to increase their share of advertising income. Similarly e-commerce sites such as Amazon.com and Dell can customize content, (e.g., information, digital products like software, advertising, promotions and other incentives) to increase repeat purchases and cross-selling.²

In spite of this potential, few models exist to help firms actually implement one-to-one marketing on the internet. Therefore, given the substantial potential arising from e-customization, it is our objective to develop a statistical and optimization approach for customization of information on the Internet. A secondary goal is to model clicking behavior on the Internet. Our procedures allow firms to customize both the content (what and how much information) and design (rendition) to their users. Specifically, we develop an approach that enables websites to customize permission-based e-mail communications in order to increase website traffic (although the approach can be more broadly applied to the issue of website customization). As such, we contribute to the growing literature on one-to-one marketing (e.g., Rossi, McCullouch and Allenby 1996; Shaffer and Zhang, 1995).

However, e-customization differs from some other contexts inasmuch as i) the “product” (or the E-mail/Web site itself) is complex, and requires explicit optimization for customization, ii) the product can be constructed dynamically, thereby making customization truly possible, and iii) the media is highly addressable, thereby facilitating targeting. Moreover, while it is possible to target products or marketing tactics in other contexts, the costs of doing so (e.g., printing individual catalogs and the cost of point of sale coupons) can often be exceedingly high.

E-customization requires knowledge of individual-level preferences, and it is therefore important to accommodate unobserved sources of preference heterogeneity (by allowing model

² The on-online market now exceeds \$171 billion and is growing rapidly.

parameters to vary across individuals). Moreover, the efficacy of content and design is also likely to differ across implementations. While it is possible to identify and account for some of the factors affecting the response to content and design in a typical application, it is unlikely that a Web marketer can enumerate all the variables that affect consumer response. Accordingly, our approach also models sources of unobserved contextual heterogeneity across e-mail content and design. This enables Web marketers to better predict potential responses to a particular type of content on a particular e-mail (or, in the case of an on-site targeting strategy, a particular type of content on a particular Web page design).

In the marketing literature, sources of heterogeneity are typically modeled using a random coefficients approach. The computer science literature on web customization suggest the use of collaborative filtering approaches to model heterogeneity (Breese et al., 1998). Collaborative filtering systems use data from users with similar preferences to recommend new items. Such systems form the basis for most commercial recommendation systems (e.g., Netperceptions and Macromedia). We show how model-based collaborative filtering can be implemented using a Bayesian semi-parametric model. Specifically, we show how a Mixture of Dirichlet Process (MDP) probit model can be used to perform collaborative filtering and to flexibly accommodate different sources of unobserved heterogeneity. The Bayesian literature in marketing has predominantly used normal distribution to capture differences across customers (Rossi, McCulloch and Allenby 1996, Allenby and Rossi 1999, Ansari, Essegaiier and Kohli, 2000). The normal distribution has limited flexibility as it is unimodal, has thin tails and does not accommodate skewness. These sources of inflexibility of the normal distribution could result in misleading inferences and inaccurate individual-level estimates (Escobar 1994). The MDP model, in contrast is flexible enough to accommodate deviations from normality and depending

upon the date, can automatically adjust to mimic either a finite mixture of support points, or a continuous distribution for heterogeneity, whichever is appropriate. Moreover, because of its discrete representation of heterogeneity, it can mimic a collaborative filtering representation. In this paper, we explicitly compare results from a MDP specification to those obtained from models that use normal distributions to capture heterogeneity.

Given a set of parameter estimates, customization then requires the construction of customized e-mails for each individual. We therefore develop an optimization procedure for customized e-mail design. The optimization procedure uses as input the individual-level parameter estimates from the hierarchical Bayesian statistical model and allows firms to a) select relevant information to include in an e-mail and b) configure the content to enhance the probability of site visits. Although the permutations in design can be quite large, we provide an exact solution to the design problem. Our design problem shows how the promise of targeting that is offered by hierarchical Bayesian models can be brought to fruition.

The rest of the paper proceeds as follows. In Section 2, we detail various approaches to customization on the Internet. Section 3 describes our statistical model for clicking behavior. In Section 4, we provide the details of our application. We overview the data, specify the model, discuss parameter estimates and perform model comparisons. In Section 5, we present the e-mail design optimization problem and the optimization approach for obtaining optimal configurations. Finally, in Section 6, we offer conclusions and suggest future research directions.

CUSTOMIZATION APPROACHES

Web sites can use a combination of onsite and external customization approaches to manage customer relationships. Both approaches are useful in enhancing site loyalty because they increase switching costs for users. When faced with a decision to switch sites, it is quite

feasible that users will be somewhat reticent to invest the time to begin “training” another firm. As noted by Alba, Lynch, Weitz, Janiszewski, Lutz, Sawyer, and Wood (1997), consumers might expect to experience switching costs and a decrease in customer service were they to switch to another site.

Onsite customization

In this approach, companies either customize the Web site to appeal to users, or enable the users themselves to customize the content. For example, portal sites such as Netscape and Altavista allow users to self-customize the site. Users of such sites can specify keywords of interest to filter news stories, can also provide lists of stocks for which they require regular information, or manipulate the page views themselves. Such user-initiated customization has obvious advantages as it elicits user preferences and gives control to users in defining what they want.

However, in many instances, such a user-initiated approach may not be completely successful as users may not be able to fully or accurately self-explicate their preferences. Many novice users may not feel confident about performing such customization actions. Moreover, preferences are dynamic and users may be reluctant to continually provide information (or find it cumbersome to do so). In such situations, company initiated customization based on revealed preferences data may be more useful.

On many Web sites, company initiated onsite customization occurs in the form of recommendation systems. For example, firms such as Amazon.com and Kraft, use recommendation systems to recommend products (Ansari et al. 2000, Gershoff and West, 1999) or content to customers.³ Most commercial recommendation systems (e.g., Netperceptions and

³ Apart from content recommendations, some companies also provide navigational aids to help customers interact with the information provided at their Web sites. Perkowitz and Etzioni (1997) show how Web sites can use logfile

Macromedia) use techniques such as collaborative or content filtering on customer ratings data to determine customers' product preferences.⁴ Recommendations can also be made using attribute-based approaches. For example, Ansari et al. (1999) show that hierarchical Bayesian models are best suited for recommendation systems because they incorporate different sources of heterogeneity and provide individual-level estimates even in sparse data environments.

These recommendation systems have typically been oriented towards *suggesting* a new product (e.g., a movie) or service rather than *designing* Web pages or e-mails. While these approaches could be adapted to the customization of content (by extending them to consider multiple, concurrent recommendations), they are more difficult to adapt to issues of customized design, because design is a large-scale (many control variable) optimization problem (e.g., how to design a new product as opposed to recommending an existing one). Our approach therefore generalizes these earlier works.

External Customization

In an external customization approach, the interest is in bringing users to a Web site. Typically, e-mails, banner advertisements, affiliate sites, or other communication media herald site content that may be of interest to site users. For example, companies such as Amazon.com, Morningstar, and New York Times regularly send e-mails containing hypertext links to the content of their Web sites to registered users. E-mails intended to attract customers to a site typically contain i) brief summaries of editorial content and ii) a link (or a set of links) to the site where more detailed information can be found. After reading the summaries, users can click on the link listed in the e-mail. By learning user preferences from their clicking histories and

analysis to suggest index pages of links to users. Such link prediction can also be used to prefetch documents while the user is reading a page. Alternatively, Web sites can suggest navigational shortcuts predicated on most popular navigational patterns. Sarukkai (2000) describes such a tour generation procedure based on Markov chains.

⁴ Collaborative filtering systems use data from users with similar preferences to recommend new items.

demographics, Web sites can tailor the messages in the e-mail to the user's interests. The greater the history of user's information, the more likely a firm can learn (and thus match) the interests of the user.

E-mail is one of the most popular Internet applications and is rapidly being adopted for e-commerce. Between 1999 and 2000, e-mail revenues increased 270% to \$342 million, and is expected to grow to \$1 billion by 2003 (The Aberdeen Group, 2001). As click-through rates on banner ads continue to drop, e-mail is becoming the instrument of choice for business to consumer communication. Customization of e-mails to suit the preferences of users is therefore of paramount importance. A number of companies such as Doubleclick, Clickaction.com, Netperceptions and Macromedia have developed e-mail marketing systems to assist companies in outbound e-mail marketing. The specific details of their implementations, unfortunately, are not publicly available.

We focus on an external customization application for several reasons (our particular application involves personalizing permission-based e-mail design and content in order to attract the e-mail recipients to a Web site). First, external customization is relatively easy to implement as firms do not have to create a significant number of alternative Web site designs (according to a 1999 Gartner Group study, an average site costs \$1 million and 5 months of implementation time). Second, internal customization strategies rely on the visitor to come to the site, which can take some time. In contrast, external customization strategies can be effective immediately, as communications are sent directly to the user by e-mail, post, or advertising. Direct communications, like e-mails, enable firms to more actively entice users to their site and are therefore useful for both acquisition and retention activities. Third, onsite customization of the Web site can be risky if users have become familiar with the interface; changes to the home page

can confuse loyal users. Fourth, as noted above, e-mail targeting is an important and growing application in its own right.

It should be noted that the application of our algorithm to e-mail does not preclude its use on web page customization. Indeed, often e-mails are sent as web pages opened directly by a web browser. Yet porting our analysis to the re-design of web sites would require careful deliberation. First, dynamic websites may confuse users, suggesting that greater benefit may accrue to varying content than design. Second, successful website customization is incumbent upon reliably identifying site visitors and may thus be of limited use in instances where dynamic ip addresses manifest, or cookies are disabled. Third, as e-mails are served daily (or less frequently), there is ample time to estimate models and serve new content between points of contact with the user. Applications such as ours are not likely to scale well to on-the-fly customization and more research is needed in that area. Nonetheless, it is possible to update web sites on a daily or weekly basis, and substantial benefits might obtain even at this frequency of customization.

MODELING APPROACH

Both aggregate and disaggregate statistical models can be used to assess the effect of e-mail design and content on click-through rates. Aggregate models relate population-level response rates (e.g., percent of the population that clicked on a link) to the design characteristics of the e-mail. As aggregate data combines each individuals' response into a single measure of population response, aggregate data greatly reduces the data processing requirements involved in e-mail design. Aggregate models will likely predict as well as individual level models for new, as opposed to existing users given the lack of historical data for such users. Moreover, aggregate models are often easier to estimate. However, as we show in this paper, aggregate models are

incapable of exploiting unobserved individual-level differences in preferences for content and design (our results indicate that an aggregate-level model yields about one-quarter the potential improvement in click-through rates realized by an individual-level model).

As our interest is in using parameter estimates to custom design e-mails, we develop an individual-level model for estimating the probability of clicking on links within e-mails. Given that our objective is to customize content and design, our modeling approach consists of two phases.

- 1) Phase I: In the first phase, we specify and estimate a probability model that correlates content and design characteristics to individual clicking likelihoods. The input to this model is individual-level clickstream data of past responses to content links included in e-mails. The output is a probability function and a set of individual-level and e-mail-level parameters that represent the preference structures of the users and the differences across e-mails.
- 2) Phase II: In the second phase, we use the probability model and the individual-level preferences as inputs to an optimization model. The optimization model recommends the optimal e-mail configuration (content and layout) for each recipient on each occasion.

As the optimization model is predicated on the results obtained from the econometric model, we first describe our statistical model and its results.

In the statistical model, we use an attribute-based approach and model the customer responses in terms of e-mail design attributes, content descriptors and user characteristics. The database contains click-through responses of users for content-links delivered over different e-mails. Let $i = 1$ to I index users, $j = 1$ to J represent e-mails and $k = 1$ to K indicate the distinct

links for which user response data is available. Each customer i provides binary responses y_{ijk} for n_{i2} links over n_{i1} e-mails. Let $E_i = \{j_1, j_2, \dots, j_{n_{i1}}\}$ be the index set of e-mails sent to user i and let $L_i = \{k_1, k_2, \dots, k_{n_{i2}}\}$ be the index set of links for which user i 's responses are available. Users differ in the number of observations (links) and similarly, e-mails differ in the number of links they contain, thus yielding a highly unbalanced data set. The total number of observations in the data set are given by $N = \sum n_{i2}$.

The observed binary responses (clicks) y_{ijk} can be modeled using a random utility framework. Users click on a particular link when the utility for exploring the content associated with the link exceeds a threshold. The relation between the observed response and the latent utility of clicking can be written as

$$(1) \quad y_{ijk} = \begin{cases} 0 & \text{if } u_{ijk} \leq 0 \\ 1 & \text{if } u_{ijk} > 0 \end{cases}$$

We model the latent utility u_{ijk} for link k of e-mail j for user i as the function of observed and unobserved e-mail and link characteristics. Specifically, the utility function can be written as:

$$(2) \quad u_{ijk} = \mathbf{x}'_{ijk} \boldsymbol{\mu} + \mathbf{z}'_{jk} \boldsymbol{\lambda}_i + \mathbf{w}'_{ik} \boldsymbol{\theta}_j + \gamma_k + e_{ijk},$$

where $e_{ijk} \sim N(0,1)$. The vector \mathbf{x}_{ijk} contains observed user, e-mail (design variables) and link-level (content) variables. The coefficients in $\boldsymbol{\mu}$ contain the ‘‘fixed effects’’ and describe the population level impacts of the independent variables. The remaining terms specify the random effects and are used for capturing different sources of heterogeneity. Note that the error specification in equation (2) assumes that errors are independent across the different links.

While it is possible to relax this assumption, we refrain from doing so for several reasons. First, such a solution would not be scalable, especially when the number of content categories is large.

The introduction of dependencies via a multivariate probit specification leads to a substantial increase in computational complexity of the estimation and optimization algorithms, making broader implementation of our model difficult. Second, the introduction of covariance terms between the utilities of links makes customized design difficult, as customization would require the knowledge of the correlations between each new link and all others, and the evaluation of a multivariate integral. Third, our model predicts response quite well even though we assume independence (we predict 649 e-mail clicks in our data and observe 639). Thus, in our view, the costs of this approach outweigh the benefits.

It is unlikely that all factors that impact responses can be isolated in any given application. It is therefore crucial to allow for multiple sources of heterogeneity. For example, users may differ in their preferences for content and in their propensities to click on links in different portions of the e-mail. Moreover, these differences in preferences may be unrelated to user demographics and other observed user variables. To allow for such unobserved preference heterogeneity across users, we introduce the term $\mathbf{z}'_{jk} \boldsymbol{\lambda}_i$ in the utility function. The vector \mathbf{z}_{jk} can contain a subset of the variables in \mathbf{x}_{ijk} such as the content descriptors and design variables pertaining to the link k and the e-mail j . The individual-specific coefficients in $\boldsymbol{\lambda}_i$ then indicate how the content and design preferences for individual i differ from the population average, $\boldsymbol{\mu}$.

Apart from modeling preference heterogeneity, in applications involving responses to information products, it is also important to model contextual heterogeneity. E-mails being information products are complex objects and therefore are not completely amenable to simplistic feature based renditions. For instance, some e-mails may be better designed than others and the design features may interact in intricate patterns, thus making it difficult to code an e-mail using few observable attributes. Thus, e-mails may differ in terms of both observed

and unobserved attributes. It is therefore desirable to account for these contextual influences using a random effects approach. Accordingly, we use the term $\mathbf{w}'_{ik} \boldsymbol{\theta}_j$ to capture the differential impact of e-mail j on the utility function. The vector \mathbf{w}'_{ik} can contain link-level as well as individual-level variables. The coefficients in $\boldsymbol{\theta}_j$ indicate how e-mail j differs from the average e-mail in terms of the impact of link-level and individual-level variables on click-through.

A broad categorization of the content of a link is possible using a few features. However, it is likely that the content remains only partially explained in terms of the observed variables. For example, while a news item can be broadly classified as a “Business News” item, there can still be considerable variation in content among all “Business News” items. We therefore use a random effect γ_k to accommodate unobserved content heterogeneity.

Dirichlet Process Priors

In this section we show how semi-parametric distributional assumptions on the population distribution of the random effects can yield a principled approach for model-based collaborative filtering. The random effects are assumed to come from a population distribution. The marketing literature has used either finite mixture distributions (Wedel and Kamakura, 2000) or continuous distributions (Ansari, Jedidi and Jagpal 2000) to represent population heterogeneity. Finite mixtures allow flexibility, but in complex models, it is difficult to determine the appropriate number of components. Moreover, it is not straightforward to incorporate multiple sources of heterogeneity in finite mixture models. Alternatively, continuous distributions are used as part of hierarchical Bayesian models to capture heterogeneity. In this case, typically a normal population distribution is used to represent the variation in random effects. While the choice of the normal distribution is made for tractability and conjugacy reasons, this assumption may not necessarily hold in reality. The normal distribution provides

limited flexibility because it is unimodal, has thin tails and does not accommodate skewness. If the population distribution is not normal, then misleading inferences about the magnitude of effects and the nature of heterogeneity are possible. Researchers have used finite mixtures of normal components (Allenby, Arora and Ginter 1998) to circumvent these problems, but in such models the difficulty of determining the number of components remains.

In this paper we show how a Mixture of Dirichlet Process (MDP) model (Escobar 1994 and MacEachern 1994) can be used to model heterogeneity in a flexible yet structured manner. The MDP model allows us to capture the uncertainty about the functional form of the population distribution using a semi-parametric approach. This model avoids the typical assumption of a parametric population distribution such as the normal and instead uses an unknown distribution F to model heterogeneity. As this population distribution is assumed to be random, in the MDP model, a Dirichlet Process Prior (Ferguson 1973, 1974; Blackwell and MacQueen 1973) is placed on the population distribution F . The Dirichlet process provides a mechanism of placing a probability distribution on the space of distributions.

The Dirichlet process prior $F \sim D(F | F_0, \alpha)$, is described by two parameters. F_0 is a parametric baseline distribution that defines the “location” of the Dirichlet process prior, and α is a positive scalar precision parameter that determines the concentration of the prior for F about the baseline distribution F_0 . The baseline distribution F_0 can be considered as a prior “guess” for the population distribution. In this paper, we use a normal distribution as the baseline distribution. The precision parameter α determines how close the non-parametric distribution F is to the baseline distribution F_0 . When α is large, a randomly sampled population distribution F is very similar to F_0 . Thus, if the baseline distribution F_0 is normal and the precision parameter, $\alpha \rightarrow \infty$, then the population distribution is a discrete distribution

that mimics a normal distribution. On the other hand, when α is small ($\alpha \rightarrow 0$), the sampled population distribution has its mass concentrated on a few points and is therefore similar to a finite mixture distribution.

As the precision parameter is inferred from the data, the MDP specification allows flexible incorporation of heterogeneity. If the nature of the heterogeneity is consistent with a normal distribution then the precision parameter is automatically adjusted to be large and the MDP yields a population distribution that mimics a normal. On the other hand, if the data come from a non-normal population distribution, then the MDP model allows enough flexibility because of its semi-parametric nature and like finite mixture models, accommodates deviations from normality. In addition, the number of “segments” is automatically determined by the MDP algorithm as outlined below. Thus the MDP model places fewer restrictions on the shape of the population distribution and as the population distribution is well approximated, it can have beneficial consequences for the accuracy of individual-level estimates (see Escobar, 1994).

In the context of our model, we assume that the user, e-mail and link specific random effects come from population distributions arising from different Dirichlet process priors. Specifically, we use

$$(3) \quad \begin{aligned} \lambda_i &\sim F_1 \\ F_1 &\sim D(F_1 | N(0, \Lambda), \alpha_1) \end{aligned}$$

to model the user-specific random effects. The above assumes that the user-specific random effects come from an unknown population distribution F_1 . The population distribution in turn comes from a Dirichlet process prior with a multivariate normal baseline distribution having a mean $\mathbf{0}$ and a unknown covariance matrix Λ . The precision parameter of the Dirichlet process

prior, α_1 , controls how close the sampled population distribution is to the baseline normal distribution.

Similarly, the e-mail random effects can be modeled as

$$(4) \quad \begin{aligned} \boldsymbol{\theta}_j &\sim F_2 \\ F_2 &\sim D(N(0, \boldsymbol{\Theta}), \alpha_2) \end{aligned}$$

where the baseline distribution is a normal with $\boldsymbol{\Theta}$ as the covariance matrix. Finally, the link-specific random effects can be modeled as

$$(5) \quad \begin{aligned} \gamma_k &\sim F_3 \\ F_3 &\sim D(N(0, \tau), \alpha_3) \end{aligned}$$

where τ represents the variance of the associated univariate baseline distribution.

A Bayesian approach is needed for inference regarding the unknowns parameters of the MDP model. The unknown quantities in our model include $\{\{u\}, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_i\}, \{\boldsymbol{\theta}_j\}, \{\gamma_k\}, \boldsymbol{\Lambda}, \boldsymbol{\Theta}, \tau, \alpha_1, \alpha_2, \alpha_3\}$. The priors over the hyperparameters are described in the appendix. The joint posterior distribution cannot be written in closed form and therefore Markov chain Monte Carlo (MCMC) methods (see Doss 1994, Bush and MacEachern 1996 and West, Muller and Escobar 1994) are needed to sample from the posterior distribution. MCMC methods involve sampling iteratively from the full conditional distributions which are described in detail in the appendix.

To understand further how the MDP model implements model based collaborative filtering and to explicate how the MDP model differs from the popular approach of using normal population distributions, we contrast here, the full conditional distribution of the user-specific random effects $\boldsymbol{\lambda}_i$ obtained from these models. The full conditional expresses the uncertainty about $\boldsymbol{\lambda}_i$ given the values of the other unknowns. In contrasting these full conditional distributions, let $\tilde{u}_{ijk\lambda} = u_{ijk} - \mathbf{x}'_{ijk} \boldsymbol{\mu} - \mathbf{w}'_{ik} \boldsymbol{\theta}_j - \gamma_k$, represent the adjusted utility for an observation.

Then $\tilde{u}_{ijk\lambda} \sim N(\mathbf{z}'_{jk}\boldsymbol{\lambda}_i, 1)$. Form the vector $\tilde{\mathbf{u}}_i$ by stacking the adjusted utilities $\tilde{u}_{ijk\lambda}$ for all the observations of the user and the matrix \mathbf{Z}_i by stacking row by row all the row vectors \mathbf{z}'_{jk} for the observations belonging to user i . When $\boldsymbol{\lambda}_i$ is assumed to be distributed normal $N(0, \boldsymbol{\Lambda})$ (as is the case for one of our null models), it is well known that the full conditional distribution for the random effects $\boldsymbol{\lambda}_i$ is multivariate normal and can be written as

$$(6) \quad p(\boldsymbol{\lambda}_i \mid \{\mathbf{u}_{ijk}\}, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\boldsymbol{\gamma}_k\}, \boldsymbol{\Lambda}) \sim N(\hat{\boldsymbol{\lambda}}_i, \mathbf{V}_i)$$

where the posterior precision is given by $\mathbf{V}_i^{-1} = \boldsymbol{\Lambda}^{-1} + \mathbf{Z}'_i \mathbf{Z}_i$, and the posterior mean is given by

$$\hat{\boldsymbol{\lambda}}_i = \mathbf{V}_i + \mathbf{Z}'_i \tilde{\mathbf{u}}_i.$$

In contrast, for the MDP model, the full conditional for $\boldsymbol{\lambda}_i$ is given by the following mixture of a normal distribution and a mass point distribution

$$(7) \quad p(\boldsymbol{\lambda}_i \mid \{\boldsymbol{\lambda}_k, k \neq i\}, \{\mathbf{u}_{ijk}\}, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\boldsymbol{\gamma}_k\}, \boldsymbol{\Lambda}) \sim q_{p0} F_b(\boldsymbol{\lambda}_i \mid \cdot) + \sum_{l \neq i} q_{pl} \delta_{\boldsymbol{\lambda}_l}$$

where

- $F_b(\boldsymbol{\lambda}_i \mid \cdot)$ is the baseline posterior distribution given in equation (6)
- The weight associated with the normal component is $q_{p0} \propto \alpha_1 f_i$ where f_i is the marginal density of the adjusted utilities for user i under the multivariate normal baseline prior density $N(0, \boldsymbol{\Lambda})$. The marginal density is obtained by integrating out the user random effects (i.e., $\int f(\tilde{u}_i \mid \boldsymbol{\lambda}_i, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\boldsymbol{\gamma}_k\}) N(0, \boldsymbol{\Lambda}) d\boldsymbol{\lambda}_i$) and is multivariate normal $N(\mathbf{0}, \mathbf{Z}_i \boldsymbol{\Lambda} \mathbf{Z}'_i + \mathbf{I}_{n_i})$, where \mathbf{I}_{n_i} is the identity matrix. The quantity f_i is obtained by evaluating at $\tilde{\mathbf{u}}_i$, the marginal density.

- $q_{pl} \propto f(\tilde{\mathbf{u}}_i | \boldsymbol{\lambda}_l, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\gamma_k\})$, the normal density of the utilities for user i evaluated using user l 's parameters, i.e., each q_{pl} is proportional to the multivariate normal density of $\tilde{\mathbf{u}}_i \sim N(\mathbf{Z}_i \boldsymbol{\lambda}_l, \mathbf{I}_{n_i})$.

The weights q_{pk} are standardized to sum to 1. δ_s is a degenerate distribution with point mass at s . Thus, in equation (7), with probability proportional to q_{p0} we sample $\boldsymbol{\lambda}_i$ from the full conditional under the baseline population distribution, i.e., using equation (6) and with probability proportional to q_{pl} we select from the degenerate distribution $\delta_{\boldsymbol{\lambda}_l}$, which means that we set $\boldsymbol{\lambda}_i = \boldsymbol{\lambda}_l$, i.e., we set person i 's parameters to be the same as person l 's. This results in a mixture with one component being a normal distribution and all other components are point masses on the parameters of other persons.

Intuitively, the above mixing scheme implies that in any iteration of the MCMC scheme, if the likelihood of observing user i 's data is relatively large using user l 's parameters, then the random effect $\boldsymbol{\lambda}_l$ for person l is more likely to be chosen as user i 's random effect. In this instance, users i and l would be in the same cluster. On the other hand, if the likelihood of observing i 's data is relatively low when user l 's random effect is used for user i , the more likely it is that the i 's random effect is a new value (generated from the baseline distribution or from a different user, l' , whose parameters generate a higher likelihood for i 's data). This mixing scheme results in a clustering of the random effects because users share common random effect parameters on any given iteration. However, the number of "clusters" and the allocation of individuals to the clusters changes from one iteration to other. Thus, this mechanism for generating the user-specific random effects is similar to what can be called model-based collaborative filtering on parameter space, as information from similar individuals is used to

predict users' preferences. In contrast to the standard model, where the estimates for a user depend on the data for other users only through the population mean and variance, here, the posterior of λ_i heavily depends on the data for the user and the data of "nearest neighbors".

APPLICATION

Data

The data for this project are furnished by one of the Web's leading Internet sites. As a condition for its use, the sponsoring firm wishes to remain anonymous and any identifying aspects of the data are therefore disguised. The organization derives the majority of its revenue from selling advertising space on its Web site to client firms and is therefore highly interested in increasing site usage. To accomplish that, they send e-mails to registered users inviting them to visit the site. These e-mails include a synopsis of several articles on the site. Below each article summary is a link to the article (the link is a URL address that readers can click on to send them to the site). It is the response to these links that forms our dependent measure.

To further clarify our exposition of the site and its content we will depict the site as an automotive news and reviews site (although the content is not automotive). We will also categorize the site's content (including the links it sends in its e-mails) much like an automotive site can be categorized as cars or trucks. These categorizations are denoted content areas and are established by the sponsoring firm based on its knowledge of the industry. On this site, a particular content area (e.g., the car portion of the site) sends a daily e-mail to users who register to receive the e-mail. These e-mail links pertain to content types within the content area. For example, the link types for a car area may include car reviews, car pricing, car specifications, automotive news, etc. Upon receiving this e-mail, a user might click on a particular type of link. If they do, this information is recorded. There are two key components to the data set provided

by this firm: i) the user log files that record usage history for a given respondent, and ii) the e-mail files that provide the date, content and design of the e-mails.

User Log Files. Each time the visitor to the site clicks on an e-mails or Web site link, a record is generated and stored. The record contains the time, the origin of the click (the IP address), and information about the link or Web page that was clicked upon. These records, collected between June and August of 1999, form one of the two key portions of the data we analyze. The origin of the click (the clicker's id) is determined via "cookies," or a record placed on the visitor's (clicker's) hard drive as well as the users IP address. Upon further visits to the site, the content provider makes note of the machine's cookie to ascertain who visited the site.

E-mail files. Visitors to the content provider's site can request to receive e-mails pertaining to information on the site. Only persons who registered receive e-mails, thus all recipients were registered.⁵ In addition to the e-mail addresses of the users, the e-mail files contain the dates of the e-mails, the links listed on the e-mails and a synopsis of information contained on those links. This information is used to code the design and content of the links provided in the e-mail, and also to infer which links were not clicked. The design of the e-mail includes i) the amount (number) of links listed in the e-mail, ii) the order of the links, and iii) the e-mail type (html or text). The link content in our data is coded into content categories provided by the firm. For example, an article reviewing an automobile would be coded "review," while automotive news would be coded "news." All in all, there are 12 of these content categories. Finally, information regarding the links in the e-mails can be merged with the database on user clicking histories to determine which links (if any) the user clicked.

⁵ It should be noted that our model applies equally well to contexts in which users are not registered. In this instance, e-mail addresses are obtained from the purchase of e-mail lists from third party vendors.

Sample Size. All totaled, there are three months of email file data and two months of log file data. The number of emails averaged about five per week and there are 1048 users in our sample.⁶ The number of links per email averaged about 5.6 and ranged from two to eight. The average response rate across links is about 7%. We used data from a random sample of 100 users to estimate the models and assigned 60% of the e-mails (11,436 observations) into an estimation sample and 40% of the e-mails into validation sets.

We created two such validation data sets. In the first, we randomly held out observations across all e-mails, persons and links – so the validation set included extant e-mails and extant links. Thus, predictions regarding a users' likelihood of clicking on a given link could be made with information regarding how others clicked on that link. The size of this validation set, which we denote EXTANT, was 3,735 observations. In the second validation data set, the sample was comprised entirely of new links and e-mails. Thus, only information regarding persons' past behavior could be used to forecast the likelihood of clicking on a link. This validation data set, which we denote NOVEL, consisted of 3,900 observations.

Ansari et al. (2000) outline the benefits of creating multiple validation data sets. In contexts wherein extant communications are to be targeted to additional individuals who have yet to receive them, the EXTANT data set is more relevant. Such may be the case after a test e-mail, or when an e-mail is sent to additional customers of a site. In contexts wherein new content is to be targeted, the NOVEL data set may be more relevant. Such may be the case when new editorial content is generated, as in the case of daily news.

⁶ Only users who responded to at least one e-mail per month are included in the sample. Persons who never or rarely respond represent 24% of the sample, but only 2% of the clicks. Due to their lack of response, there is little information available to improve this number. Thus, customization is likely to be inefficacious for these persons so we omit them from our analysis. We note that the practice of including only respondents has an analog in scanner data research pertaining to category incidence (Manchanda, Ansari, and Gupta 1999; Mela, Jedidi, and Bowman 1998).

Model Specification

In Section 3 we developed a general random effects model that can be applied across a number of contexts. Here we describe the specific instantiation for our application. We include several observed link variables, e-mail design variables, and user variables in our specification. Moreover, in addition to these observed effects, we also allow for different sources of unobserved heterogeneity across users, e-mails, and links. Our specification can be described as follows:

$$\begin{aligned} u_{ijk} = & \mu_1 + \mu_2 \text{Content}_1 + \dots + \mu_{12} \text{Content}_{11} + \mu_{13} \text{Position} + \mu_{14} \text{Num - Items} + \mu_{15} \text{Text} + \mu_{16} \text{Since} \\ & + \lambda_{i1} + \lambda_{i2} \text{Content}_1 + \dots + \lambda_{i,12} \text{Content}_{11} + \lambda_{i,13} \text{Position} + \lambda_{i,14} \text{Num - Items} + \lambda_{i,15} \text{Since} \\ & + \theta_{j1} + \theta_{j2} \text{Content}_1 + \dots + \theta_{j,12} \text{Content}_{11} + \theta_{j,13} \text{Position} + \theta_{j,14} \text{Since} \\ (8) \quad & + \gamma_k + e_{ijk} \end{aligned}$$

Link Variables. The link variables characterize both the editorial content of the link and its position within the e-mail. The likelihood of clicking on a link will certainly be a function of its content (e.g., news vs. reviews) because users are expected to differ in their preference for content categories. The editorial content within the e-mails can be categorized into twelve categories (or content types). These content types are included using a set of eleven dummy variables (*Content1-Content 11*).

As is the case with traditional media, the placement of the link may affect response (click-through). Hanssens and Weitz (1980) find that the later the advertisement is presented in a magazine, the less likely it is to be seen or read. Hoque and Lhose (1999) have replicated this result in electronic media. They argue that the impact of placement is magnified in electronic media because it is more difficult to read online and because of the effort involved in scrolling. The ordinal position of the link within the e-mail is represented by a variable (*Position*).

E-mail Variables. We specify two e-mail design variables. The first represents the number of links within the e-mail (Num-Items). Houston and Scott (1984) show that increasing the number of advertisements within a print environment lowers the readership of any particular advertisement on a page. We expect a similar result to exist in electronic settings as an increase in the number of links exacerbates clutter and therefore increases the cognitive costs of perusing the e-mail. For some users, the cost of having to scan a large number of links may exceed the potential value of reading the e-mail. The second e-mail-level variable (*Text*) characterizes whether the e-mail is a text e-mail or an html e-mail.⁷

Person Variable. As the duration of time increases since the last link clicked, the likelihood of a subsequent click may change. Thus, we include a covariate for days since the last click (*Since*). We use one observation to initialize this variable. In customer relationship management (CRM) applications, an increase in time since last purchase leads to a decrease in subsequent purchase likelihood (Schmittlein, Morrison and Colombo 1987). Thus, we expect that the coefficient for (*Since*) will typically be negative as persons who have not clicked in some time are less likely to be active.

The descriptive information for the link, e-mail, and person variables is presented in Table 1.

[Insert Table 1 About Here]

User Heterogeneity. We model unobserved user heterogeneity by specifying a random effects model for both the intercepts and slopes in our model. Thus, we specify i) user specific random intercepts (which captures differences across users in their propensities to click on links, and ii) user specific random slopes (which capture differences across users in their response to link

⁷ We tested for an interaction between *NumItems* and *Position*. This did not enhance model fit. We also tested for non-linear effects for these variables and found none.

content, e-mail variables, and time since last click). We note that user and random effects are identified only for those variables that are not fixed across persons (similarly, e-mail random effects are only identified for variables that are not fixed across e-mails).

E-mail Heterogeneity. We capture e-mail heterogeneity by allowing i) an e-mail-specific random intercept which captures e-mail “attractiveness” and ii) e-mail-specific random slopes which capture differences across e-mails. The unobserved e-mail variables (captured via θ_j in Equation 8) interact with the link-level variables such as the content of the link and position of the link within the e-mail. These interactions allow us to model contextual heterogeneity. As it is not possible to completely describe e-mails using observed attributes, incorporating e-mail-level heterogeneity is crucial for modeling click-through probabilities.

Link Heterogeneity. For parsimony, we capture link heterogeneity by including a single link-specific random effect. This term captures the impact of unobserved link effects (e.g., the particular editorial (news or review) content of a link type on a given day) that is left unexplained by the content variables describing the link.

RESULTS

We estimated three models on the data. The first model is a simple model (Model S) that includes no heterogeneity. The second model (Model N) includes all observed variables and accounts for person, e-mail and link specific unobserved sources of heterogeneity, using normal population distributions for the random effects. The third model (Model DP) uses Dirichlet process priors to account for user, e-mail and link heterogeneity.

The models were estimated using Markov chain Monte Carlo methods. The full conditional distributions used in MCMC sampling for the DP model are described in Appendix 1. For each model, the chain was run for 30,000 iterations. The results reported are based on a

sample of 20,000 draws from the posterior distribution, after discarding 10,000 burn-in draws. Convergence was ensured by monitoring the time series of draws.

Model Selection

Model Fit. We use the pseudo-Bayes factor (PsBF) (Geisser and Eddy 1979, Gelfand 1996) to compare different models. The PsBF is based on the cross-validation predictive density which can be very conveniently computed for our models using the MCMC draws. Let \mathbf{y} be the observed data, let y_{il} represent the l th observation for individual i and let $\mathbf{y}_{(il)}$ represent the data with the observation l for individual i deleted. The cross-validation predictive density can be written as

$$(9) \quad \pi(y_{il} | \mathbf{y}_{(il)}) = \int \pi(y_{il} | \boldsymbol{\beta}, \mathbf{y}_{(il)}) \pi(\boldsymbol{\beta} | \mathbf{y}_{(il)}) d\boldsymbol{\beta}$$

where $\boldsymbol{\beta}$ is the vector of all parameters in the model. The PsBF for comparing two models (M1 and M2) is expressed in terms of the product of cross-validation predictive densities and can be written as

$$(10) \quad \text{PsBF} = \prod_{i=1}^I \prod_{l=1}^{n_i} \frac{\pi(y_{il} | \mathbf{y}_{(il)}, M1)}{\pi(y_{il} | \mathbf{y}_{(il)}, M2)}.$$

The PsBF summarizes the evidence provided by the data for M1 against M2 and its value can be interpreted as the number of times model M1 is more (or less) probable than model M2.

The PsBF for our model can be calculated easily from a sample of d MCMC draws $\{\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_d\}$. As $\boldsymbol{\beta}$ is the vector of all parameters, including all the random effects, the binary responses y_{il} , $i = 1$ to I , $l = 1$ to n_i are conditionally independent given $\boldsymbol{\beta}$. Thus, a Monte Carlo estimate of $\pi(y_{il} | \mathbf{y}_{(il)})$ can be obtained as

$$(11) \quad \hat{\pi}(y_{il} | \mathbf{y}_{(il)}) = \left[\frac{1}{d} \sum_{t=1}^d \frac{1}{(p_{il}^{(t)})^{y_{il}} (1 - p_{il}^{(t)})^{1-y_{il}}} \right]^{-1}$$

where $p_{il}^{(t)}$ is the probability $\Pr(y_{il} = 1 | \beta^{(t)}) = 1 - \Phi(u_{il}^{(t)} - \mathbf{x}'_{il} \boldsymbol{\mu}^{(t)} - \mathbf{z}'_{j(l)k(l)} \boldsymbol{\lambda}_i^{(t)} - \mathbf{w}'_{ik(l)} \boldsymbol{\theta}_{j(l)}^{(t)} - \gamma_{k(l)}^{(t)})$,

where the superscript t represents the t -th draw of the MCMC sampler and $j(l)$ denotes the e-mail associated with the l th observation and $k(l)$ represents the link associated with that observation.

Gelfand (1996) provides the derivation for the above expression for the cross-validation predictive distribution. The estimates from equation (11) can be used to calculate the logarithms of the numerator and denominator of the PsBF. These quantities can be considered as a surrogate for the log-marginal data likelihoods for the two competing models. Based on the MCMC output, the log marginal data likelihood for the Dirichlet Process Model (Model DP) is -2290.25 , the log marginal data likelihood for the normal model (Model N) is -2340.45 and for the non-heterogeneous model is -2807.22 . Accordingly, the pseudo-Bayes factor implies an $\exp(50.2)$ improvement for model DP over model N, and an $\exp(517.0)$ improvement for model DP over model S. Thus the pseudo-Bayes factors provide support for the Dirichlet process model.

Model Predictions

The predictive ability of the three competing models can be assessed using the validation data. The link-level probabilities for the links within an e-mail can be used to compute an e-mail-level probability of at least one click-through from the e-mail. To compute this probability, we first use the observed data to determine the likelihood of clicking on each link in the e-mail, \hat{p}_{ijk} . The complement of these predictions yields the likelihood of the link is not clicked $1 - \hat{p}_{ijk}$. Under the assumption of independence, the product of these link-level non-click likelihoods yields the joint probability that the e-mail was not clicked. The complement of this probability,

$\hat{p}_{ij} = (1 - \prod_{k=1}^K (1 - \hat{p}_{ijk})),$ yields the probability that at least one link was clicked. The

independence assumption seems plausible as it predicts 649 e-mail clicks when compared to 639 e-mail clicks in our data. The link-level estimated probabilities of click-through can be compared against a classification threshold to predict the response on any given observation in the validation data. For example, a classification threshold of 0.5 means that if the estimated clicking probability for a observation is greater than 0.5 then the observation is classified as a click, otherwise it is classified as a non-click. Similarly, the e-mail-level probabilities can be compared against a classification threshold to predict whether a given e-mail will elicit a click from a given individual.

We use Receiver Operating Characteristic curves (ROC) to compare the predictive performance of our models for a range of thresholds. ROC analysis is used in psychology and in medical statistics for signal detection purposes (e.g., Egan 1975; Swets, 1979; Metz, 1986). It illustrates the tradeoffs between two types of errors. For any threshold that is used for predicting observations as clicks or non-clicks, two types of errors can occur. False negatives are errors that occur when actual clicks are predicted as non-clicks. False positives are errors that occur when non-responses, (i.e., non-clicks) are predicted as responses (i.e., clicks). The proportion of actual click-throughs that are predicted as non-clicks is called the False Negative Fraction (FNF) and the proportion of actual non-clicks that are predicted as clicks is called the False Positive Fraction (FPF).

An ROC curve (also called a Lorenz diagram) is constructed by plotting the True Positive Fraction (1-FPF) against the false positive fraction (FPF) for a range of possible classification thresholds. The resulting plot is represented over a unit square. Different points on the ROC curve correspond to different classification threshold values used for prediction. The area under

the ROC curve, A_z provides a summary measure of the quality of the model. A model with an ROC curve that tracks the 45° line would be worthless as it would not separate the two classes of observations at all (i.e., there would be as many false positives as true positives). Such a curve has $A_z = 0.5$. In contrast a perfect model would have a ROC curve that follows the two axes and would have $A_z = 1$. The ROC curve, as it spans across all possible thresholds, yields a more complete picture of predictive accuracy than measures predicated solely upon a single classification threshold (e.g., a hit rate using a 50% classification threshold).

[Insert Figure 1 About Here]

Figure 1 compares the ROC curves for the three models using the EXTANT data set. The top panel shows the ROC curves based on the link-level probability estimates and the bottom panel shows the curves obtained from the e-mail-level probabilities. In each of these panels, the dark solid curve pertains to the DP model, the normal model is represented by the grey curve and the non-heterogeneous model is depicted by the dashed curve. In both panels we observe that the two models with heterogeneity (Model DP and Model N) have similar predictive abilities. Moreover, the ROC curves for the two heterogeneous models dominate the ROC curve for the non-heterogeneous model (Model S) at all values of the classification threshold.

Comparing the areas under the link-level ROC curves, we find that $A_z = 0.85$, for Model DP, $A_z = 0.84$ for Model N and $A_z = 0.76$ for Model S. For the e-mail-level ROC curves, $A_z = 0.81$ for Model DP, $A_z = 0.80$ for Model N and $A_z = 0.74$ for Model S. These results suggest that for our data set, Model DP performs slightly better than Model N and that both substantially outperform a model with no heterogeneity.

For the NOVEL data set, the conclusions are similar. We find that $A_z = 0.77$, for Model DP, $A_z = 0.75$ for Model N and $A_z = 0.68$ for Model S. The email-level ROC curves also provide similar conclusions. Specifically, $A_z = 0.69$ for Model DP, $A_z = 0.66$ for Model N and $A_z = 0.56$ for Model S. As with EXTANT, Model DP predicts slightly better than Model N while both models are superior to the non-heterogeneous model. These findings reinforce the importance of modeling heterogeneity for customization purposes.

In summary, the PsBF, and the ROC analysis indicate that the models that account for heterogeneity are preferable to the model that ignores sources of difference in parameters. The fact that unobserved heterogeneity matters indicates that there could be sufficient gains from customization. Moreover, we find that the DP model is superior to Model N, based on PsBF and evidences slightly better predictive performance on the validation data.

Parameter Estimates

We now report the parameter estimates for the Dirichlet Process model. As reported earlier, these estimates are based on a post burn-in sample of 20,000 MCMC draws. Table 2 reports the parameter estimates. These parameter estimates are discussed below.

[Insert Table 2 About Here]

Design Variables. The response rates for links within text e-mails appear to be no greater than the response rates for links within html e-mails. As anticipated, the effect of link order is negative, indicating that the effectiveness of links decreases as the link appears later in the e-mail. In contrast to our expectations, we observe no population effect for the number of links within an e-mail. In our data, the actual number of links never exceeded eight, and this may have been too few to generate a negative effect of clutter. We find that there is considerable design heterogeneity; persons react differently to different designs. Thus, the number of items

and the link position can influence the responses of at least some individuals in our sample. The magnitude of the e-mail heterogeneity also indicates that some e-mails evidence a greater effect of link order than others. The finding that some respondents are more likely to respond to text, while others respond better to html may reflect the influence of bandwidth. Given the html is bandwidth-intensive, those recipients with lesser bandwidth may prefer text-based e-mails.

Content Variables. At the population level, content types differ in their ability to elicit click-throughs. Moreover, the standard deviations associated with the normal baseline distribution for the unobserved user random effects clearly reveal that there is considerable user preference heterogeneity in the data. People differ in their preferences for different content types. There is a sizable degree of heterogeneity in content preference across e-mails. This heterogeneity presumably arises from editorial and design differences across the e-mails (e.g., a review in one e-mail may be of more interest than a review in another e-mail because of how it interacts with unobserved contextual variables). Taking into account the fixed and random effects, it is clear that the content variables play an important role in predicting click-throughs.

Heterogeneity. The Since variable has a negative impact on response. This suggests that the greater the duration since the previous click, the less likely it is that a user will click on a link within the e-mail. The precision parameters α_1, α_2 and α_3 associated with the DP priors suggest that there is greater clustering in the random effects than what is evidenced by normal population distributions. For example, $\alpha_1 \approx 103$ implies an average of 61 clusters for the users. Similarly, $\alpha_2 \approx 115$ and $\alpha_3 \approx 383$ imply on average 65 clusters for the e-mails and 383 clusters for the link random effects. This coupled with the fact that the PsBF favors the DP model, indicate that the population distributions deviate from normality and justifies the need for a semi-parametric approach.

It is informative to compare the different sources of heterogeneity. A variance decomposition of the random terms in the utility function shows that person heterogeneity accounts for 28.37%, the email heterogeneity accounts for 38.99%, the link-heterogeneity for 1.85% and the residual errors account for 31.77% of the total random variation. This implies substantial improvements in model performance can be realized by modeling multiple sources of heterogeneity.

CUSTOMIZED E-MAIL DESIGN

The customer and e-mail specific parameter vectors obtained from the statistical model are used to forecast potential customer reactions to proposed changes in e-mail content and configuration. These forecasts can be used to determine the optimal design and content of an e-mail for each customer. Given an objective function, combinatorial optimization is used to customize the e-mail design for each e-mail and each person.

In e-mail marketing situations, management's primary goal is to maximize the expected number of click-throughs within e-mails. A secondary objective might be to maximize the probability that at least one link is clicked within an e-mail. Given the objective function, design optimization involves 1) selecting from the set of available content the specific content (links) to be included for each person and each e-mail and 2) configuring the e-mail layout to maximize objectives. It is important to note that in our sample, even though number of links within the e-mail (*Num-Items*) does not have an effect at the population level, there is considerable heterogeneity at the user level and therefore, content selection can be important for at least some individuals.

Let n be the total number of content links available to be included in a particular e-mail. For the aforementioned objective functions, the design problem of selecting from the n available

links and then ordering the included links can be solved in two stages. In the first stage, n linear assignment subproblems are solved. We index each of these n subproblems by k and let $k \in \{1, \dots, n\}$. In the second stage the solution for that first stage subproblem which yields the largest objective function is chosen as the optimal solution for the original problem.

In the first stage, let us consider the k th subproblem of assigning the n available links to k contiguous positions within an e-mail. Let x_{ij} be a binary variable that is equal to 1 when content link i is present in ordinal position j ($j = 1$ for the first position and k for the lowest position) within an e-mail, and is 0 otherwise. Similarly, let $p_{ij|k}$ be the probability that the person clicks on link i if it is placed in position j , when the total number of included items is k . When interest is in maximizing the expected number of click-throughs, the problem of assigning n links to k positions is given by the linear assignment specification:

$$\text{Maximize } \sum_{i=1}^n \sum_{j=1}^k p_{ij|k} x_{ij}$$

subject to

$$\sum_{j=1}^k x_{ij} \leq 1, \quad \text{for } i = 1, 2, \dots, n$$

$$\sum_{i=1}^n x_{ij} = 1, \quad \text{for } j = 1, 2, \dots, k$$

The first constraint in the above specification ensures that each content link i can be in at most one of the k destinations. The second constraint ensures that each of the k destinations (ordinal position) can have at most one link. Each of the n subproblems, corresponding to $k = 1$ to n can be solved using the Hungarian method (Foulds, 1984, pp 72-76) which provides an efficient and simple approach. The solution to the k th subproblem yields an assignment \mathbf{x}_k^* of the content

links to k contiguous positions, and the optimal value of the objective function V_k . In each problem, where $k \neq n$, $n - k$ links are left unassigned and are hence not included in the e-mail.

In the second stage, the solutions to the n subproblems are compared to determine the subproblem that yields that maximum value for the objective function. That is, we determine the subproblem m such that $V_m = \max\{V_1, \dots, V_n\}$. The solution \mathbf{x}_m^* is chosen as the solution of the original problem. This yields the optimal content and configuration of the e-mail.

Similarly, when the interest is in maximizing the probability of at least one click from an e-mail the k th second stage subproblem of assigning n links to k positions can be written as:

$$\text{Minimize } \sum_{i=1}^n \sum_{j=1}^k \log(1 - p_{ij|k}) x_{ij}$$

subject to

$$\sum_{j=1}^k x_{ij} \leq 1, \quad \text{for } i = 1, 2, \dots, n$$

$$\sum_{i=1}^n x_{ij} = 1, \quad \text{for } j = 1, 2, \dots, k$$

Minimizing the objective function in the above specification is the same as maximizing the probability that at least one link is clicked, $1 - \prod_{ij} (1 - p_{ij|k})^{x_{ij}}$, and therefore, the problem is of the linear assignment type. The optimal solution can be ascertained as before, by comparing the n first stage subproblems in terms of the objective function.

Optimization Results

In this section we report the optimization results, which are based upon the validation data. In particular, we compare and contrast the results from the optimization procedure detailed

above (**Optimal**) with those obtained using two sub-optimal, but simpler procedures. The first sub-optimal procedure (**Ordering**) uses a single linear assignment algorithm (instead of n) for each individual to merely reorder the existing links within an e-mail. This procedure therefore ignores the content selection aspect of the optimization and is expected to do well when clutter does not significantly influence the probability of click-through. The second sub-optimal procedure (**Greedy**) uses a greedy heuristic to reorder the links. In this heuristic, links are assigned sequentially from the first position within the e-mail to the last position. To determine which link resides in which position for a given individual, the link having their highest probability of click-through (among the set of unassigned links) is assigned to the highest (uppermost) remaining position. Each algorithm yields a customized solution for each individual and e-mail, predicated on their link level clicking likelihoods. Below, we report the results for each validation set described in the data section.

EXTANT VALIDATION DATA. Table 3 reports the optimization results for the two objective functions and the three optimization procedures using the estimates from the Dirichlet process model (Model DP) and the Normal Model (Model N). Columns 1 and 2 display the results when the objective is to maximize the probability of at least one click from an e-mail. The entries in the first two columns contain the mean probability of at least one click, across all the e-mails in the validation sample. Columns 3 and 4 report the results when the objective is to maximize the expected number of click-throughs in an e-mail. The entries in these columns give the mean of the expected number of clicks across all the e-mails in the validation sample. Columns 1 and 3 present the results for the Model DP, while Columns 2 and 4 indicate the Model N results. As the Model DP and Model N results are similar, we discuss only the Model DP results.

Table 3
EXTANT OPTIMIZATION RESULTS

| Configurations | Prob (At-Least One Click) | | Expected Number of Clicks | |
|----------------|---------------------------|--------|---------------------------|--------|
| | DP | Normal | DP | Normal |
| Original | 0.23 | 0.23 | 0.34 | 0.32 |
| Greedy | 0.34 | 0.36 | 0.51 | 0.51 |
| Ordering | 0.35 | 0.36 | 0.53 | 0.52 |
| Optimal | 0.36 | 0.39 | 0.55 | 0.57 |

The first row gives the predicted results when the original configuration within the data is e-mailed and therefore serves as a benchmark for assessing the performance of the various optimization approaches. For example, the table shows that mean probability of at least one click from an e-mail is 0.23 when the original e-mails (as designed by the site) are sent. The second row gives the predicted results for the Greedy procedure, the third row gives the results arising from the Ordering procedure, while the last row gives the results using the Optimal two stage procedure. Thus, the table indicates that, for the Model DP, the mean probability of at least one click can increase to 0.36 if the optimal e-mails are sent.

Thus, comparing the entries in the first column, we see that Optimal procedure can potentially yield a 56% improvement over the original configuration in the mean probability of at least one click-through. In contrast, the improvement for the Ordering procedure is 52%. Decomposing the total potential improvement into the portion arising from re-ordering and the portion arising from content selection suggests that re-ordering results in 92% of the total improvement, while content selection constitutes the balance. Further analysis suggests that Optimal algorithm improves the likelihood of at least one click over Ordering algorithm for 43% of the e-mails sent out in the validation sample. A majority of the gains are small and arise from the users who react adversely to clutter (i.e., have a negative coefficient for the variable, number of items). Were users more averse to clutter, content selection would matter even more.

The Greedy procedure results in a 48% increase in predicted e-mail click rates (from 0.23 to 0.34). Thus, the algorithm performs nearly as well as the Ordering algorithm. Nonetheless, one should be wary of predicting similar improvements in different data. In particular, when subjects have a positive coefficient for the positive variable (i.e., they tend to scroll to the bottom of the e-mail and then click on one of those links at the bottom), then the Greedy algorithm we use will not do well.

The entries in the second column show that the Optimal optimization procedure can potentially yield a 62% improvement over the original configuration in the expected number of clicks. In contrast, the improvements for the Ordering procedure and the Greedy algorithm are 53% and 50% respectively. Furthermore, for about 42% of the e-mails sent out in the validation sample, Optimal made an improvement over and above Ordering. Similarly, for 58% of the e-mails, Optimal was better than Greedy, albeit, the magnitude of improvement was small in most cases.

The Optimal procedure is better than the other two procedures and leads to improvements in response rates for e-mails, especially when clutter adversely affects the probability of clicking. In addition, the linear assignment problem which forms the basis for the optimal procedure can be solved quickly and efficiently, even for large problems. The greedy solution, although simpler, performs poorly for users who have a propensity to click at the bottom of the e-mail. In contrast, the Ordering algorithm can do relatively well in this situation yet performs poorly when clutter matters. In general, the Ordering procedure and the Greedy procedure are marginally computationally less demanding. The choice between these depend upon the balance between accuracy, speed and complexity desired by managers.

When applying a similar analysis using the parameters from Model S (the non-heterogeneous Model), the optimal two-stage optimization procedure yields a predicted improvement in the mean probability of at least one click-through of only 12.5% and an improvement of 15.4% for the second objective function. This relative lack of improvements predicted by the simpler model when compared to the heterogeneous models arise because it has limited flexibility in differentiating among individuals and e-mails. It therefore generates e-mail configurations that are optimal only for the average individual or e-mail.

NOVEL VALIDATION DATA. Table 4 reports the optimization results for the NOVEL validation data set. Column 2 indicates that the mean probability of at least one click increases 20% from 0.45 when the original emails are sent to 0.54 if the optimal emails are sent. The improvement for the Ordering procedure is similar, i.e., 18%. Decomposing the total potential improvement into the portion arising from re-ordering and the portion arising from content selection suggests that re-ordering results in 90% of the total improvement, while content selection constitutes the balance.

Table 4
NOVEL OPTIMIZATION RESULTS

| Configurations | Prob (At-Least One Click) | | Expected Number of Clicks | |
|----------------|---------------------------|--------|---------------------------|--------|
| | DP | Normal | DP | Normal |
| Original | 0.45 | 0.45 | 0.61 | 0.60 |
| Greedy | 0.53 | 0.53 | 0.75 | 0.74 |
| Ordering | 0.53 | 0.53 | 0.75 | 0.74 |
| Optimal | 0.54 | 0.54 | 0.75 | 0.75 |

For the second objective function, the results obtained from the DP model estimates in the fourth column show that the Optimal optimization procedure can potentially yield a 23% improvement over the original configuration in the expected number of clicks. The

improvements for the Ordering procedure and the Greedy algorithm are also 23%. The results for the Normal model are similar and are shown in the fifth column.

Model S yields a predicted improvement in the mean probability of at least one click-through of 14.2% and an improvement of 15.8% for expected clicks. However, given that heterogeneous models predict significantly better than the simple model with no heterogeneity, we place greater credence on the optimization results from the heterogeneous models.

Finally, we note that the improvement in hit rates using the NOVEL data, while substantial, is not as great as the improvement indicated by the EXTANT data. The difference arises because the information regarding the clicking behavior of other persons on a particular link is informative about the targeting of links. When possible, content providers should seek to “test market” information, as this enhances targeting efficacy.

CONCLUSIONS

The advent of the Internet has enhanced the ability of marketers to personalize communications and engender relationships with consumers. By enabling the right content to reach the right person at the right time, the Web can yield substantial dividends to Web marketers and can enhance the quality of service to consumers. However, this promise is contingent upon learning more about consumer preferences and developing techniques that enable marketers to fulfill these preferences. Our objective has been to facilitate that task.

Accordingly, we describe an approach to harness the potential afforded by the Web to determine individual-level preferences, and then develop an algorithm to customize content predicated on those preferences. In the context of targeting and customizing e-mails that herald content in a Web magazine, we develop a customization system that uses a Mixtures of Dirichlet

Processes probit model coupled with an optimization model to personalize communications on the Internet.

Our adaptation of the Dirichlet Processes model is unique from previous implementations as it incorporates multiple sources of heterogeneity, uses the probit framework, and is applied to a large-scale data application. In contrast to the normal model, the MDP model predicates an individual's choice behavior on that of the individual's "nearest neighbors". As such, a comparison of the MDP and normal models provides some insight into the additional predictive power of model based collaborative filtering. Our model comparison results indicate that the MDP model is preferable to the normal model based on the pseudo-Bayes factor. The predictive performance of the MDP model is slightly better than that of the normal model. We leave a thorough comparison of these alternative approaches for future research, but do note that our approach is tailored to the problem, as opposed to the data. By its virtue of flexibility, there may be instances in which the Dirichlet Processes model substantially outperforms the normal model (i.e., non-normal heterogeneity). The converse is unlikely to be true, as normal models do not adjust well to non-normal heterogeneity.

Given that the additional programming demands inherent in the Dirichlet Processes model (over that of the normal model) are negligible, the trade-off between the approaches is an issue of flexibility and scalability. The computational demands of the normal model are simpler, and therefore we recommend the normal model when scalability of the model is a major concern. Moreover, the scalability of either model to the demands of sending emails for many users requires a careful decomposition of the overall requirements into offline and online (i.e., real time) components. For example, aggregate features of the models, such as the population distribution, can be estimated offline based on a sample of users. Moreover, the population

distribution can be updated periodically (say weekly) as new data arrives. Once the population distribution is known, obtaining the estimates for a particular individual is not very computationally intensive and can be done relatively quickly. The optimization component is very quick as it is based on the linear assignment problem that has quick and exact solutions.

Our approach adds to the targeting and customization literature in marketing by integrating heterogeneous choice models with optimization techniques to personalize content. Specifically, describe an optimization algorithm based on the assignment algorithm to optimize the design and content of electronic communications. We believe that such a general approach (combining choice models with optimization models (Tellis and Zufryden 1995; Rossi, McCulloch and Allenby 1996)) has utility beyond e-mail customization and can be potentially used in the design of tailored services, custom-designed catalogs, and bundling of goods.

The results of our model indicate that the design of the e-mail is crucial in affecting click-through probabilities. For example, we find that the order of content matters and that there exists a great deal of heterogeneity across persons in their preferences, and a great deal of heterogeneity across links and e-mails in terms of their effectiveness in design and content. Capitalizing on these results, we demonstrate that design and content can indeed be optimized. We find that response rates (expected click-throughs) could be increased by 62% by customizing the e-mail's design.

Finally, we propose that our analysis be extended along two dimensions – modeling other behaviors in the Internet environment and using our methodological approach in other contexts. With regard to other Internet applications, it would be desirable to customize Web content in an effort to increase the frequency of site visits and clicks per visit. Similarly, our general approach could be adapted to the design of e-commerce sites and personal agents. It is also possible to use

our underlying methodology to target advertising content. Another pressing problem centers about optimal contact strategies for e-mail communications. Excessive contact, or “over-touching”, can lead to unsubscribe decisions. Infrequent contact could lead to few responses. Moreover, these effects could vary by user. With respect to other contexts, the proposed model could be extended beyond e-commerce models to more traditional models of targeting purchase opportunities, such as direct mail marketing or product customization. Our design approach could be adapted to conjoint tasks. The conjoint domain may prove especially promising as affective measures can be integrated with behavioral ones to enhance the predictive capability of the model, and the researcher may have more latitude in the design of the stimulus set. It is our hope that this analysis will encourage further research along these dimensions.

APPENDIX: PRIOR DISTRIBUTIONS AND FULL CONDITIONAL DISTRIBUTIONS

Priors

The model as implemented in our application can be written succinctly as $u_{ijk} = \mathbf{x}'_{ijk} \boldsymbol{\mu} + z'_{jk} \boldsymbol{\lambda}_i + \mathbf{w}'_{ik} \boldsymbol{\theta}_j + \gamma_k + e_{ijk}$, where $e_{ijk} \sim N(0,1)$, $\boldsymbol{\lambda}_i \sim DP(N(0, \boldsymbol{\Lambda}), \alpha_1)$, $\boldsymbol{\theta}_{jl} \sim DP(N(0, \Theta_{ll}), \alpha_2)$ for $l = 1$ to p , and $\gamma_k \sim DP(N(0, \tau), \alpha_3)$. We set $\text{var } e_{ijk} = 1$ to identify the probit. Notice that the random effects in the vector $\boldsymbol{\theta}_j$ are assumed to be independent in order to obtain a parsimonious specification. Thus Θ_{ll} is a univariate quantity. Similarly, the link-level random effect γ_k is a scalar and τ represents the across links heterogeneity.

We specify independent, diffuse but proper priors on $\{\boldsymbol{\mu}, \boldsymbol{\Lambda}, \boldsymbol{\tau}, \{\Theta_{ll}\}, \alpha_1, \alpha_2, \alpha_3\}$. The prior for $\boldsymbol{\mu}$ is multivariate normal $N(\boldsymbol{\eta}, \mathbf{C})$. The covariance matrix \mathbf{C} is diagonal with large values for the variances to reflect uncertainty. We use $\boldsymbol{\eta} = \mathbf{0}$ and $\mathbf{C} = 1000I$, where I is the identity matrix. The precision matrix $\boldsymbol{\Lambda}^{-1}$ associated with the population distribution, $\boldsymbol{\lambda}_i \sim DP(N(\mathbf{0}, \boldsymbol{\Lambda}), \alpha_1)$, is a $(m+1) \times (m+1)$ positive definite matrix (where m is the number of coefficients that vary across people). We assume Wishart priors: $W(\iota, (\iota \mathbf{L})^{-1})$ for the precision matrix $\boldsymbol{\Lambda}^{-1}$. The matrix \mathbf{L} can be considered as the expected prior variances of the $\boldsymbol{\lambda}_i$'s. Smaller values for ι correspond to more diffuse prior distributions. We set $\iota = m + 1$, and $\mathbf{L} = \text{diag}(1)$. As Θ is a diagonal matrix, the coefficients in $\boldsymbol{\theta}_j$ are assumed to be independent. Independent inverse gamma priors $IG(3,1)$ are specified over the diagonal elements Θ_{ll} , $l = 1$ to p , where p represents the number of coefficients that vary across e-mails. We use a single random effect in γ_k . The prior for the variance τ is chosen to be inverse gamma $IG(a,b)$ with

$a = 3$ and $b = 1$. Finally, following Escobar and West, (1997), we assume independent gamma priors $Ga(a,b)$, where $a = 1$ and $b = 1$ for the precision parameters α_1, α_2 and α_3 .

Gibbs Sampler

An Gibbs sample involves iteratively sampling from the full conditional distributions of the unknowns. Draws after a burn-in period are used for inference. The full conditional can be written as follows:

- (a) As the utilities are distributed normal and as the prior for $\boldsymbol{\mu}$ is $N(\boldsymbol{\eta}, \mathbf{C})$, the full conditional distribution for $\boldsymbol{\mu}$ is multivariate normal and can be written as

$$(i) \quad p(\boldsymbol{\mu} | \mathbf{u}_i, \{\boldsymbol{\lambda}_i\}, \{\boldsymbol{\theta}_j\}, \gamma_k) \sim N(\hat{\boldsymbol{\mu}}, \mathbf{V}_\mu)$$

where $\mathbf{V}_\mu^{-1} = \mathbf{C}^{-1} + \mathbf{X}'\mathbf{X}$, and $\hat{\boldsymbol{\mu}} = \mathbf{V}_\mu (\mathbf{X}'\tilde{\mathbf{u}} + \mathbf{C}^{-1}\boldsymbol{\eta})$. The matrix \mathbf{X} is obtained by stacking row by row all the row vectors \mathbf{x}'_{ij} . The vector $\tilde{\mathbf{u}}$ is obtained by stacking all the elements $\tilde{u}_{ijk} = u_{ijk} - \mathbf{z}'_{jk}\boldsymbol{\lambda}_i - \mathbf{w}'_{ik}\boldsymbol{\theta}_j - \gamma_k$, for all the observations in the sample.

- (b) The respondent random effects $\boldsymbol{\lambda}_i$ are assumed to come from a random distribution F_1 .

According to the MDP specification, the uncertainty about F_1 is captured by using a Dirichlet Process prior, i.e., $F_1 \sim D(N(0, \boldsymbol{\Lambda}), \alpha_1)$. Blackwell and MacQueen (1973) show that is one integrates over the random distribution F_1 , A Polya urn representation of the Dirichlet process results. Thus the prior for $\boldsymbol{\lambda}_i$ can be described in terms of a series of conditional distributions of the form

$$\boldsymbol{\lambda}_i | \boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_{i-1} \sim \frac{1}{i-1 + \alpha_1} \sum_{m=1}^{i-1} \delta_{\boldsymbol{\lambda}_m} + \frac{\alpha_1}{i-1 + \alpha_1} F_0$$

where F_0 is the baseline distribution $N(0, \Lambda)$. The full conditional distribution for λ_i is obtained by combining the likelihood and the prior conditional on λ_i given above and can be written as

$$(ii) \quad p(\lambda_i | \{\lambda_l, l \neq i\}, \{\mathbf{u}_{ijk}\}, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\gamma_k\}, \Lambda) \sim q_{p0} F_b(\lambda_i | \cdot) + \sum_{l \neq i} q_{pl} \delta_{\lambda_i}$$

where

- $F_b(\lambda_i | \cdot)$ is the baseline posterior distribution given by $N(\hat{\lambda}_i, \mathbf{V}_i)$. The posterior precision for this baseline distribution is given by $\mathbf{V}_i^{-1} = \Lambda^{-1} + \mathbf{Z}'_i \mathbf{Z}_i$, and the posterior mean is given by $\hat{\lambda}_i = \mathbf{V}_i \mathbf{Z}'_i \tilde{\mathbf{u}}_i$. The matrix \mathbf{Z}_i is obtained by stacking row by row all the row vectors \mathbf{z}'_{jk} for the observations belonging to user i . The vector $\tilde{\mathbf{u}}_i$ is obtained by stacking the adjusted utilities $\tilde{u}_{ijk\lambda} = u_{ijk} - \mathbf{x}'_{ijk} \boldsymbol{\mu} - \mathbf{w}'_{ik} \boldsymbol{\theta}_j - \gamma_k$, for all the observations of the user.
- $q_{p0} \propto \alpha_1 f_i$ where $f_i = \int f(\tilde{\mathbf{u}}_i | \lambda_i, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\gamma_k\}) N(0, \Lambda) d\lambda_i$ is the marginal density of the adjusted utilities for user i under the multivariate normal baseline prior density $N(0, \Lambda)$. The marginal density is obtained by integrating out the user random effects and is the density at $\tilde{\mathbf{u}}_i$ of the multivariate normal distribution $N(\mathbf{0}, \mathbf{Z}_i \Lambda \mathbf{Z}'_i + \mathbf{I})$, where \mathbf{I} is the identity matrix.
- $q_{pl} \propto f(\tilde{\mathbf{u}}_i | \lambda_l, \boldsymbol{\mu}, \{\boldsymbol{\theta}_j\}, \{\gamma_k\})$, the normal density of the utilities for user l evaluated using user l 's parameters, i.e., each q_{pl} is proportional to the multinormal density of $\tilde{\mathbf{u}}_i \sim N(\mathbf{Z}_i \lambda_l, \mathbf{I})$.

The weights q_{pl} , $\forall l$ are standardized to sum to 1. δ_s is a degenerate distribution with point mass at s . Thus, in the above full conditional distribution, with probability proportional to

$$\alpha_l |\mathbf{V}_i|^{1/2} |\mathbf{\Lambda}|^{-1/2} \exp\left\{\frac{1}{2} \tilde{\mathbf{u}}_i' \mathbf{G}_i \tilde{\mathbf{u}}_i\right\}$$

we sample λ_i from the full conditional under the baseline population distribution. In the above expression, $\mathbf{G}_i = (\mathbf{Z}_i \mathbf{V}_i \mathbf{Z}_i' - I_{n_i})$ using the binomial inverse theorem (see Press, 1971, page 23). With probability proportional to

$$\exp\left\{-\frac{1}{2} (\tilde{\mathbf{u}}_i - \mathbf{Z}_i \lambda_i)' (\tilde{\mathbf{u}}_i - \mathbf{Z}_i \lambda_i)\right\}$$

we select from the distribution δ_{λ_i} , which means that we set $\lambda_i = \lambda_l$. This results in a mixture with one component being a normal distribution and all other components are point masses.

- (c) After choosing the random effects for each user in step (b), because of the discrete nature of the population distribution, the users are naturally grouped into clusters of users with the same λ_i 's. In any iteration, there are some number I^* , ($0 < I^* \leq I$) of unique values for the user-specific λ_i 's. Denote the unique user parameters as λ_m^* , $m = 1 \dots I^*$. These can be interpreted as cluster-specific coefficients. In addition, let S_m represent the set of users with common random effects λ_m^* . Bush and MacEachern (1996) suggest remixing of these cluster-specific parameters to aid the convergence for the Gibbs sampler. After

determining the cluster structure of the user-specific parameters, the cluster-level parameters λ_m^* are recomputed from the conditional density

$$(iii) \quad p(\lambda_m^* | \cdot) = N(\mathbf{Q}_m \sum_{i \in S_m} \mathbf{Z}'_i \tilde{\mathbf{u}}_i, \mathbf{Q}_m)$$

where $\mathbf{Q}_m^{-1} = \Lambda^{-1} + \sum_{i \in S_m} \mathbf{Z}'_i \mathbf{Z}_i$. Notice that this full conditional distribution results because the unique values of the random effects come from the base-line distribution. Because of the conjugacy of the normal distribution with the normal utilities, the full conditional distribution is normal. Selecting a new value for the cluster parameter λ_i^* changes the $\lambda_{i,s}$ for the users within that cluster.

(d) The e-mail random effects θ_j can be generated from the mixture distribution

$$(iv) \quad p(\theta_j | \{\theta_l, l \neq j\}, \mathbf{u}_j, \boldsymbol{\mu}, \{\lambda_i\}, \{\gamma_k\}, \Theta) \sim q_{e0} F_b(\theta_j | \cdot) + \sum_{l \neq i} q_{el} \delta_{\theta_l}$$

where

- $F_b(\theta_j | \cdot)$ is the baseline posterior distribution given by $N(\hat{\boldsymbol{\theta}}_j, \mathbf{V}_j)$ where the posterior precision is given by $\mathbf{V}_j^{-1} = \Theta^{-1} + \mathbf{W}'_j \mathbf{W}_j$, and the posterior mean is given by $\hat{\boldsymbol{\theta}}_j = \mathbf{V}_j \mathbf{W}'_j \tilde{\mathbf{u}}_j$. The matrix \mathbf{W}_j is obtained by stacking row by row all the row vectors \mathbf{w}'_{ik} for the observations belonging to e-mail j . The vector $\tilde{\mathbf{u}}_j$ is obtained by stacking the adjusted utilities $\tilde{u}_{ijk\theta} = u_{ijk} - \mathbf{x}'_{ijk} \boldsymbol{\mu} - \mathbf{z}'_{jk} \lambda_i - \gamma_k$, for all the observations of the e-mail.
- $q_{e0} \propto \alpha_2 \int f(\mathbf{u}_j | \theta_j, \boldsymbol{\mu}, \{\lambda_i\}, \{\gamma_k\}) N(0, \Theta) d\theta_j$, just α_2 times the marginal density of the utilities for e-mail j under the multivariate normal baseline prior density $N(0, \Theta)$. This marginal density is obtained by integrating out the e-mail random effects.

- $q_{el} \propto f(\tilde{\mathbf{u}}_j | \boldsymbol{\theta}_l, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_j\}, \{\gamma_k\})$, the normal density of the utilities for e-mail j evaluated using e-mail l 's parameters.

The weights $q_{el}, \forall l$ are standardized to sum to 1.

- (e) As for the user effects, a remixing step is used for the e-mail-specific parameters. The e-mails can be grouped into clusters with the same $\boldsymbol{\theta}_{jS}$. There are some number

J^* , ($0 < J^* \leq J$) of unique values for the e-mail-specific $\boldsymbol{\theta}_{jS}$. Denote the unique e-mail parameters as $\boldsymbol{\theta}_m^*$, $m = 1 \dots J^*$. In addition, let C_m represent the set of e-mails with common random effects $\boldsymbol{\theta}_m^*$. After determining the cluster structure, the cluster-level parameters $\boldsymbol{\theta}_m^*$ are recomputed from the conditional density

$$(v) \quad p(\boldsymbol{\theta}_m^* | \cdot) = N(\mathbf{Q}_m \sum_{j \in C_m} \mathbf{W}'_j \tilde{\mathbf{u}}_j, \mathbf{Q}_m)$$

where $\mathbf{Q}_m^{-1} = \boldsymbol{\Theta}^{-1} + \sum_{j \in C_m} \mathbf{W}'_j \mathbf{W}_j$. Selecting a new value for the cluster parameter $\boldsymbol{\theta}_j^*$ changes the $\boldsymbol{\theta}_{jS}$ for the e-mails within that cluster.

- (f) The random effect γ_k for link k can be generated from a univariate mixture distribution.

$$(vi) \quad p(\gamma_k | \{\gamma_t, t \neq k\}, \{u_{ijk}\}, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_i\}, \{\boldsymbol{\theta}_j\}, \boldsymbol{\tau}) \sim q_{l0} F_b(\gamma_k | \cdot) + \sum_{t \neq k} q_{lt} \delta_{\gamma_t}$$

where

- $F_b(\gamma_k | \cdot)$ is the baseline posterior distribution given by the univariate normal distribution $p(\gamma_k | \{u_{ijk}\}, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_i\}, \{\boldsymbol{\theta}_j\}, \boldsymbol{\tau}) \sim N(\hat{\gamma}_k, v_k)$ where $v_k^{-1} = \boldsymbol{\tau}^{-1} + n_k$, n_k is the total number of observations pertaining to link k and $\hat{\gamma}_k = v_k \sum \tilde{u}_{ijk\gamma}$, where the summation is over all the observations pertaining to link k and

$$\tilde{u}_{ijk\gamma} = u_{ijk} - \mathbf{x}'_{ijk} \boldsymbol{\mu} - \mathbf{z}'_{jk} \boldsymbol{\lambda}_i - \mathbf{w}'_{ik} \boldsymbol{\theta}_j.$$

- $q_{t_0} \propto \alpha_3 \int f(\mathbf{u}_k | \gamma_k, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_i\}, \{\boldsymbol{\theta}_j\}) N(0, \tau) d\tau_k$, just α_3 times the marginal density of the utilities for link k under the normal baseline prior density $N(0, \tau)$. This marginal density is obtained by integrating out the link-level random effects and is the density at $\tilde{\mathbf{u}}_k$ of the multivariate normal distribution $N(0, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma} = \mathbf{1}\mathbf{1}' + I$, where $\mathbf{1}$ is a vector of 1's.
- $q_{t_t} \propto f(\mathbf{u}_k | \gamma_t, \boldsymbol{\mu}, \{\boldsymbol{\lambda}_j\}, \{\boldsymbol{\theta}_j\})$, the normal density of the utilities for link k evaluated using link t 's parameters.

The weights $q_{t_t}, \forall t$ are standardized to sum to 1. δ_s is a degenerate distribution with point mass at s . Thus, in the above full conditional distribution, with probability proportional to q_{t_0} we sample $\gamma_k s$ from the full conditional under the baseline population distribution, and with probability proportional to q_{t_t} we select from the distribution δ_{γ_t} , which means that we set $\gamma_k = \gamma_t$.

- (g) As with the user and e-mail random effects, the use of the Dirichlet process also induces a cluster structure on the link-specific parameters. A remixing step can be used to recompute the cluster-specific parameters to facilitate better convergence. Specifically, the parameter γ_m^* associated with the cluster Cl_m can be sampled from the full conditional distribution $p(\gamma_m^* |) \sim N(\mathbf{v}_m \boldsymbol{\Sigma} \tilde{\mathbf{u}}_{ijk\gamma}, \mathbf{v}_m)$. The summation in the full conditional is over all the observations associated with the links belonging to cluster m . The precision, $\mathbf{v}_m^{-1} = \tau^{-1} + n_m$ and n_m is the number of observations associated with the m th cluster Cl_m .

(h) The Wishart prior $W(\iota, (\mathbf{L})^{-1})$ for the precision matrix $\mathbf{\Lambda}^{-1}$, is conjugate to the normal distribution of the random effects. Thus the full conditional distribution of $\mathbf{\Lambda}^{-1}$ is also Wishart, and can be computed from the cluster-specific parameters λ_m^* as

$$(vii) \quad p(\mathbf{\Lambda}^{-1} | \{\lambda_m^*\}) \sim W\left(\left[\sum_{m=1}^{I^*} \lambda_m^* \lambda_m^{*'} + \mathbf{L}\right]^{-1}, \iota + I^*\right)$$

(i) As the inverse gamma prior $IG(a, b)$ is conjugate to the normal distribution of the e-mail random effects, the full conditional distribution of each variance of the e-mail-level baseline distribution, $\Theta_{ll}, l = 1$ to p , (the diagonal elements of Θ) is inverse gamma, and is given by

$$(viii) \quad p(\Theta_{ll} | \{\theta_{ml}^*\}) \sim IG\left(\frac{J^*}{2} + a, \left[\frac{1}{2} \sum_{m=1}^{J^*} \theta_{ml}^{*2} + b^{-1}\right]^{-1}\right)$$

where J^* represents the total number of distinct e-mail clusters in the dataset.

(j) The full conditional distribution of the across link variance τ is inverse gamma, and is given by

$$(ix) \quad p(\tau | \{\gamma_m^*\}) \sim IG\left(\frac{K^*}{2} + a, \left[\frac{1}{2} \sum_{m=1}^{K^*} \gamma_m^{*2} + b^{-1}\right]^{-1}\right)$$

where K^* is the total number of distinct link clusters.

(k) The precision parameters α_1, α_2 and α_3 can be sampled using data-augmentation (West, 1992). Here we show this step for α_1 . At the n th iteration of the Gibbs sampler, we first sample a latent variable ζ from the beta distribution $(\zeta | \alpha_1^{(n-1)}, I) \sim Be(\alpha_1^{(n-1)} + 1, I)$, which has mean $(\alpha_1^{(n-1)} + 1)/(\alpha_1^{(n-1)} + 1 + I)$. Then $\alpha_1^{(n)}$ is sampled from the mixture of gamma distributions, i.e.,

$$(\alpha_1^{(n)} | \zeta, I^{*(n-1)}) \sim \pi_\zeta Ga(a + I^{*(n-1)}, b - \log(\zeta)) + (1 - \pi_\zeta) Ga(a + I^{*(n-1)} - 1, b - \log(\zeta))$$

where the weights π_ζ are defined in odds forms by

$$\frac{\pi_\zeta}{1 - \pi_\zeta} = \frac{a + I^{*(n-1)} - 1}{I(b - \log(\zeta))}$$

The parameters α_2 and α_3 can be drawn similarly.

(1) Finally, as part of a data-augmentation step, the utilities can be sampled. As

$u_{ijk} \sim N(m, 1)$ a priori, where the mean $m = \mathbf{x}'_{ijk}\boldsymbol{\mu} + \mathbf{z}'_{jk}\boldsymbol{\lambda}_i + \mathbf{w}'_{ik}\boldsymbol{\theta}_j + \gamma_k$, the full

conditional distribution for the utility u_{ijk} is a truncated normal distribution

$u_{ijk} \sim tn(m, 1)$. This distribution is truncated from above at 0 if the dependent variable

$y_{ijk} = 0$ and is truncated from below at 0 if $y_{ijk} = 1$ (Albert and Chib, 1993).

Table 1
VARIABLE DESCRIPTORS

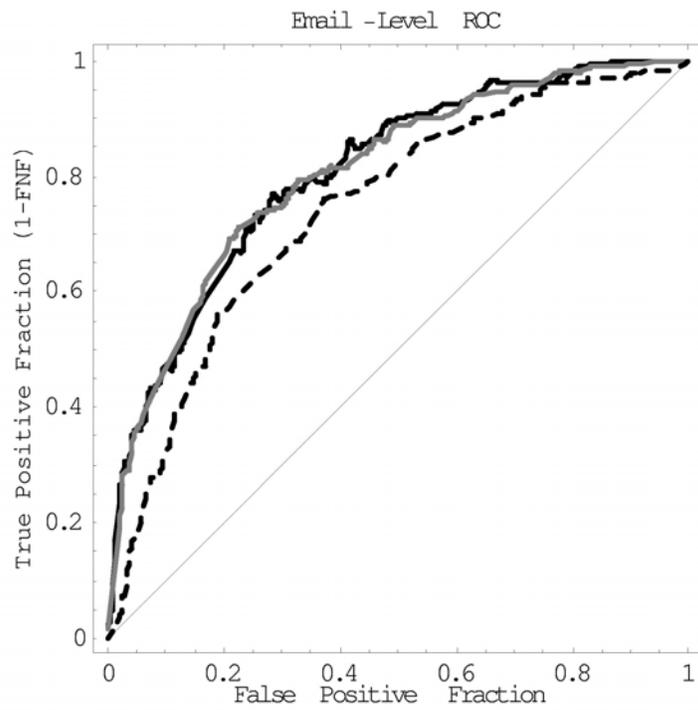
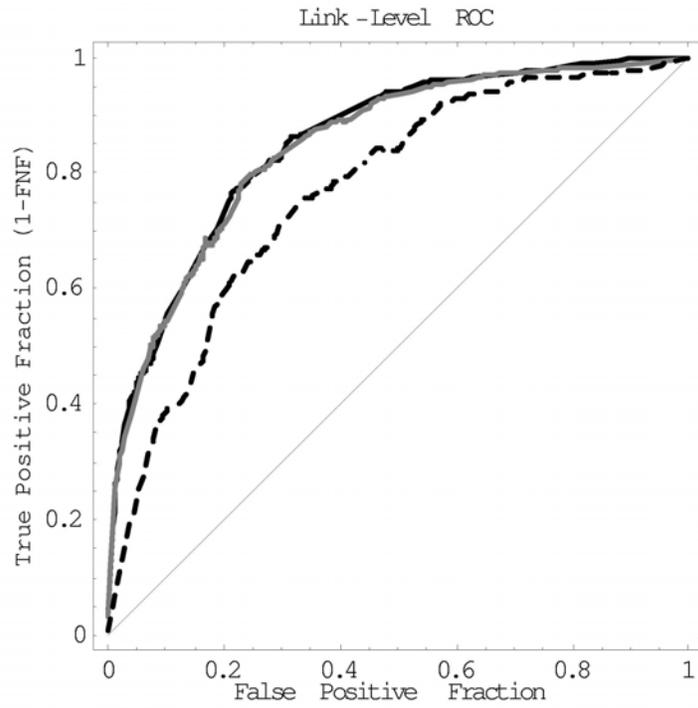
| Variable | Calibration Data | | Validation Data EXTANT | | Validation Data NOVEL | |
|--------------|------------------|-----------|---------------------------|-----------|--------------------------|-----------|
| | Mean | Std. Dev. | Mean | Std. Dev. | Mean | Std. Dev. |
| Link Click | 0.075 | 0.263 | 0.077 | 0.250 | 0.087 | 0.282 |
| E-mail Click | 0.307 | 0.461 | 0.278 | 0.448 | 0.350 | 0.477 |
| Text | 0.635 | 0.481 | 0.623 | 0.485 | 0.621 | 0.485 |
| Num-Items | 5.609 | 0.784 | 5.620 | 0.803 | 5.784 | 1.021 |
| Position | 3.305 | 1.656 | 3.310 | 1.663 | 3.392 | 1.747 |
| Since | 7.197 | 9.372 | 8.111 | 10.204 | 6.160 | 8.011 |
| Content1 | 0.181 | 0.385 | 0.180 | 0.384 | 0.179 | 0.384 |
| Content2 | 0.182 | 0.386 | 0.182 | 0.386 | 0.179 | 0.384 |
| Content3 | 0.182 | 0.386 | 0.182 | 0.386 | 0.179 | 0.384 |
| Content4 | 0.069 | 0.255 | 0.065 | 0.246 | 0.068 | 0.252 |
| Content5 | 0.062 | 0.242 | 0.064 | 0.244 | 0.065 | 0.246 |
| Content6 | 0.061 | 0.239 | 0.064 | 0.245 | 0.052 | 0.221 |
| Content7 | 0.036 | 0.187 | 0.038 | 0.191 | 0.035 | 0.184 |
| Content8 | 0.036 | 0.187 | 0.037 | 0.189 | 0.033 | 0.179 |
| Content9 | 0.022 | 0.148 | 0.023 | 0.150 | 0.033 | 0.178 |
| Content10 | 0.037 | 0.188 | 0.036 | 0.187 | 0.036 | 0.185 |
| Content11 | 0.032 | 0.177 | 0.035 | 0.184 | 0.031 | 0.173 |

Table 2
PARAMETER ESTIMATES FOR MODEL DP

| Variables | Fixed Effects μ | 95% Probability Interval | Std. of Base distribution Across users | Std. of Base distribution across e-mails |
|------------|------------------------|--------------------------|---|---|
| Intercept | -1.47 | (-2.97-0.59) | 0.51 | 0.45 |
| Content1 | 0.25 | (-0.48, 0.79) | 0.97 | 0.44 |
| Content2 | 1.07 | (0.57, 1.47) | 0.34 | 0.52 |
| Content3 | 0.21 | (-0.54, 0.71) | 0.65 | 0.47 |
| Content4 | 0.29 | (-0.56, 0.72) | 0.35 | 0.71 |
| Content5 | 1.21 | (0.21, 1.62) | 0.93 | 0.52 |
| Content6 | 0.79 | (0.13, 1.25) | 0.38 | 0.57 |
| Content7 | -0.67 | (-2.45, 0.19) | 0.49 | 0.72 |
| Content8 | 0.80 | (-0.11, 1.43) | 0.38 | 0.52 |
| Content9 | 0.28 | (-0.95, 0.93) | 0.45 | 0.56 |
| Content10 | 0.32 | (-0.59, 0.89) | 0.49 | 0.62 |
| Content11 | -1.20 | (-3.38, 0.27) | 0.54 | 1.45 |
| Position | -0.37 | (-0.59, -0.19) | 0.23 | 0.27 |
| Since | -0.24 | (-0.59, -0.19) | 0.17 | 0.28 |
| Num-Items | 0.01 | (-0.17, 0.13) | 0.18 | |
| Text | 0.29 | (-0.33, 0.65) | | |
| α_1 | 103.32 | (69.37, 130.87) | | |
| α_2 | 114.25 | (77.26, 144.88) | | |
| α_3 | 383.07 | (319.93, 431.89) | | |

1. The parentheses contain the 2.5th and the 97.5th percentiles.

Figure 1: Receiver Operating Characteristic Curves for Link-Level and e-mail-Level Probabilities.



References

- Aberdeen Group (2001, "e-Mail Marketing: Relevancy, Retention, and ROI," Working Paper, <http://www.aberdeen.com/abcompany/hottopics/emailmarketing/default.htm>
- Alba, Joseph, John Lynch, Barton Weitz, Chris Janiszewski, Richard Lutz, Alan Sawyer, and Stacy Wood, (1997) "Interactive Home Shopping: Consumer, Retailer, and Manufacturer Incentives to Participate in Electronic Marketplaces," *Journal of Marketing*, 61, 3 (July), 38-53.
- Allenby, G. M., N. Arora and J. L. Ginter (1998), "On the Heterogeneity of Demand," *Journal of Marketing Research*, 35, 3, 384-389.
- , and P. Rossi, (1999), "Marketing Models of Consumer Heterogeneity," *Annals of Econometrics, Marketing and Econometrics*. Eds. Wansbeek, T. and M. Wedel, 57-79.
- Ansari, Asim, Skander Essegaier and Rajeev Kohli (2000), "Internet Recommender Systems," *Journal of Marketing Research*, 37, 3, 363-375.
- , K. Jedidi, and S. Jagpal (2000) "A Hierarchical Approach for Modeling Heterogeneity in Structural Equation Models," *Marketing Science*, 328-347.
- Blackwell, D. and J. B. MacQueen (1973), "Ferguson Distribution Via Polya Urn Schemes," *The Annals of Statistics*, 1, 353-355.
- Bush, C. A. and S. N. MacEachern (1996) "A Semi-Parametric Bayesian Model for Randomized Block Designs," *Biometrika*, 83, 275-285.
- Doss, H. (1994) "Bayesian Nonparametric Estimation for Incomplete Data Via Successive Substitution Sampling," *The Annals of Statistics*, 22, 1763-1786.
- Egan, J. P. (1975), *Signal Detection Theory and ROC Analysis*, Academic Press, New York.
- Escobar, M. D. (1994), "Estimating Normal Means with a Dirichlet Process Prior," *Journal of the American Statistical Association*, 89, 268-277.
- , M. D., and M. West (1997), "Computing Bayesian Nonparametric Hierarchical Models," in *Practical Nonparametric and Semiparametric Bayesian Statistics*, (eds: P. Mueller et al.), New York: Springer Verlag.
- Ferguson, T. S. (1973), "A Bayesian Analysis of Some Nonparametric Problems," *The Annals of Statistics*, 1, 209-230.
- , T. S. (1974), "Prior Distributions on Spaces of Probability Measures," *The Annals of Statistics*, 2, 615-629.

- Foulds, L. R. (1984) *Combinational Optimization for Undergraduates*, Springer-Verlag, New York.
- Gershoff, Andrew and Patricia M. West (1998), "Using a Community of Knowledge to Build Intelligent Agents," *Marketing Letters*, 9, 79-91.
- Geisser, S., and Eddy, W. (1979), "A Predictive Approach to Model Selection," *Journal of the American Statistical Association*, 74, 153-160.
- Gelfand, A. E. (1996), "Model Determination Using Sampling-based Models," in *Monte Carlo Markov Chain in Practice*, W. R. Gilks, S. Richardson, and D. J. Spiegelhalter, ed. Chapman and Hall: London.
- Hanssens, Dominique M., and Barton A. Weitz (1980), "The Effectiveness of Industrial Print Advertisements Across Product Categories," *Journal of Marketing Research*, 17, 3 (August), 294-306.
- Hoque, Abeer Y., and Gerald L. Lohse (1999), "An Information Search Cost Perspective for Designing Interfaces for Electronic Commerce," *Journal of Marketing Research*, 36, 3 (August), 387-394.
- Houston, F. S. and Carol Scott (1984), "The Determinants of Advertising Page Exposure," *Journal of Advertising*, 13, 2, 27-33.
- MacEachern, S. N. (1994) "Estimating Normal Means with a Conjugate Style Dirichlet Process Prior," *Communications in Statistics: Simulation and Computation*, 23, 727-741.
- Manchanda, P., A. Ansari, and S. Gupta (1999), "The Shopping Basket: A Model for Multicategory Purchase Incidence Decisions," *Marketing Science*, 18, 2, 95-114.
- Mela, C. F., K. Jedidi and D. Bowman (1998), "The Long-Term Impact of Promotions on Consumer Stockpiling Behavior," *Journal of Marketing Research*, 35, 2, 250-262.
- Metz, C. E. (1986) "ROC Methodology in Radiologic Imaging," *Investigative Radiology*, Vol. 21, 720-733.
- Perkowitz, M. and Oren Etzioni (1997), "Adaptive Web Sites: Automatically Learning from User Access Patterns," in Proceedings of 6th International WWW Conference, Santa Clara, CA., 1997.
- Press, J. S. (1971), *Applied Multivariate Analysis*, Holt, Rinehart and Winston, Inc., New York.
- Rossi, Peter E., Robert E. McCulloch and Greg M. Allenby (1996), "The Value of Purchase History Data in Target Marketing," *Marketing Science*, 15, 4, 321-340.

- Sarukkai, Ramesh (2000), "Link Prediction and Path Analysis Using Markov Chains," *Computer Networks*, 33, 377-386.
- Schmittlein, D. C., D. G. Morrison, and R. Colombo (1987), "Counting Your Customers: Who Are They and What Will They Do Next?" *Management Science*, 33, 1, 1-24.
- Shaffer, G. and Z. John Zhang (1995), "Competitive Coupon Targeting," *Marketing Science*, 14, 395-416.
- Swets, J. A. (1979) "ROC Analysis Applied to the Evaluation of Medical Imaging Techniques," *Investigative Radiology*, Vol. 14, 109-121.
- Tellis, Gerard J. and Fred S. Zufryden (1995), "Tackling the Retailer Decision Maze: Which Brands to Discount, How Much, and When?" *Marketing Science*, 14, 3, 271-299.
- Wedel, M. and W. Kamakura (2000), *Market Segmentation: Conceptual and Methodological Foundations*, Second Edition. Kluwer, Boston.
- West, M. (1992), Hyperparameter Estimation in Dirichlet Process Mixture Models," *ISDS Discussion Paper #92-A03*, Duke University.
- , P. Muller, and M. D. Escobar (1994), "Hierarchical Priors and Mixture Models with Applications in Regression and Density Estimation," *Aspects of Uncertainty: A Tribute to D. V. Lindley* (eds. A. F. M. Smith and P. R. Freeman), London: John Wiley and Sons, 363-386.