



## Marketing Science

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

### Predicting Customer Value Using Clumpiness: From RFM to RFMC

Yao Zhang, Eric T. Bradlow, Dylan S. Small

To cite this article:

Yao Zhang, Eric T. Bradlow, Dylan S. Small (2015) Predicting Customer Value Using Clumpiness: From RFM to RFMC. Marketing Science 34(2):195-208. <https://doi.org/10.1287/mksc.2014.0873>

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

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

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

Copyright © 2015, INFORMS

Please scroll down for article—it is on subsequent pages



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

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

# Predicting Customer Value Using Clumpiness: From RFM to RFMC

Yao Zhang

Fixed Income Research Division, Credit Suisse, New York, New York 10010,  
yao.a.zhang@gmail.com

Eric T. Bradlow, Dylan S. Small

The Wharton School, University of Pennsylvania, Philadelphia, Pennsylvania 19104  
{ebradlow@wharton.upenn.edu, dsmall@wharton.upenn.edu}

In recent years, customer lifetime value (CLV) has gained increasing importance in both academia and practice. Although many advanced techniques have been proposed, the recency/frequency/monetary value (RFM) segmentation framework, and its related probability models, remain a CLV mainstay. In this article, we demonstrate the deficiency in RFM as a basis for summarizing customer history (data compression), and extend the framework to include clumpiness (C) by a metric-based approach. Our main empirical finding is that C adds to the predictive power, above and beyond RFM and firm marketing action, of both the churn, incidence, and monetary value parts of CLV. Hence, we recommend a significant implementation change: from RFM to RFMC.

This work is also motivated by noting that although statistical models based on RFM summaries can fit well in aggregate, their use can lead to significant micro-level (e.g., ranking of customers) prediction errors unless C is captured. A set of detailed empirical studies using data from a large North American retailer, in addition to six companies that vary in their business model: two traditional (e.g., CDNow.com) and four Internet (e.g., Hulu.com), demonstrate that the “clumpiness phenomena” is widely prevalent, and that companies with “bingeable content” have both high potential and high risk segments, previously unseen, but now uncovered because of the new framework: RFM to RFMC.

**Keywords:** customer lifetime value; RFM; clumpiness

**History:** Received: April 10, 2013; accepted: June 22, 2014; Preyas Desai served as the editor-in-chief and Scott Neslin served as associate editor for this article. Published online in *Articles in Advance* September 2, 2014, updated October 1, 2014.

## 1. Introduction

In marketing, customer lifetime value (CLV) is defined as the net present value of the cash flows attributed to the relationship with a customer, and therefore reflects a customer's future profitability (Gupta et al. 2004). By making good use of CLV as a marketing metric, managers tend to place greater emphasis on customer service and long-term customer satisfaction, rather than on maximizing short-term sales; and recent studies have shown CLV centrality as a good long-run firm strategy (Kumar 2008, Fader 2013). Thus, coupled with the easy availability of transaction data, growing empirical evidence on the impact of CLV on firm revenues and profitability have made CLV an increasingly important concept in both academia and practice. Companies such as Harrah's, IBM, Capital One, LL Bean, ING, and others routinely use CLV as a tool to manage and measure the success of their business (Davenport et al. 2010). Academics have written scores of articles and dozens of books on this topic in the past decade (e.g., Gupta et al. 2006, Kumar et al. 2008), because of its central importance.

Although many advanced statistical models and data mining techniques have been proposed to estimate CLV, recency, frequency, and monetary value (RFM) segmentation as a method, remains a mainstay of the industry because of its ease of implementation in practice. Recent research that shows its statistical sufficiency under certain assumptions (Fader et al. 2005) also provides a theoretical justification. More specifically, from the direct marketing literature historically, RFM is based on the practice/foundation that there exist three important variables to summarize a customer's purchase history: recency (R), frequency (F), and monetary value (M). The simplest model classifies customers into groups based on each of these variables, and generates a score for each group. Mailing or other marketing communication programs are then prioritized based on the scores of different RFM groups. Many researchers and consultants have developed more formal scoring models (i.e., regression-type models) that attempt to predict customers' future behavior (see, e.g., Baesens et al. 2002, Berry and Linoff 2004, Bolton 1998, Malthouse 2003, Malthouse and Blattberg

2005, Parr Rud 2001, Blattberg et al. 2008), but CLV via RFM still remains widely used.

Whereas RFM or related simple empirically driven scoring models attempt to predict customers' behavior in the future, and are therefore implicitly linked to CLV, they have several limitations (Fader et al. 2005, Kumar 2006). As a result, a series of more formal probability models have been developed that incorporate unobserved attrition, which links the RFM paradigm directly to CLV (Schmittlein et al. 1987; Fader et al. 2005, 2010). These models have resulted in a well-grounded behavioral story, and therefore can be used to value customers in terms of their future number of purchases, lifetime duration, or CLV (and separately understand the subcomponents of it, i.e., purchasing while alive versus churn). The research also (as mentioned above) shows that RFM variables are sufficient statistics for this class of CLV models (Fader et al. 2005), which provides a theoretical justification for their use, but also implies that R, F, and M provide a complete customer summary for CLV prediction. It is this tenet that we question and expand on here.

In particular, we utilize the recent research of Zhang et al. (2013) that proposes a new individual-level measure of a customer's history, (C) clumpiness—defined as the degree of nonconformity to equal spacing—and demonstrate its properties as an important component in the better understanding of CLV and its subcomponents. Our primary research questions therefore are the following: Although RFM is very attractive for practitioners because it only requires the monitoring of *three* easy-to-compute variables, do they fully explain key aspects of customer behavior? Are there key features missing? Is it really true that there are no other *systematic factors*, i.e., C, that can help in identifying valuable customers and predicting their future purchase behavior? We note that whereas extant research has developed more complicated underlying statistical models that allow for clumpiness (i.e., hidden Markov models (HMMs): Netzer et al. 2008, Schweidel et al. 2009), and hence the concept of clumpy data and the need to model it is not new, our approach is complementary (but very different) in that it is a purely metric-based approach that can be computed, stored, and used as easily as R, F, and M. In addition, as we demonstrate, by including C as an additional component of CLV prediction, researchers can determine *prior* to running more complicated models (i.e., HMMs versus RF-based models) whether there is a likely payoff in increasing statistical complexity; hence a value to researchers. Related to this idea of a priori complexity needed, and to clearly delineate our paper's contribution, it is equally important for us to state what this paper does not do. Although our paper does provide a new metric-based approach to improve the understanding of CLV and its subcomponents, we do

not provide a new class of probability models for CLV. This is an important task, beyond the scope of this paper, and is being addressed by many of the current aforementioned HMM papers. In addition, we view Zhang et al. (2013), and our work, as complementary and in the right sequence where that paper establishes the legitimacy of computing C in the way we do, and this paper discusses its application in customer valuation.

To demonstrate that C is broadly applicable as a metric, we have obtained seven data sets that vary in terms of their industry (e.g., retailing versus video consumption) and regularity in consumption to demonstrate the properties of C as both a predictor (X variable) of CLV, and an outcome (Y variable) to be correlated with demographics and marketing action. In particular, a data set from a large North American retailer (who wishes to remain anonymous) that includes data on purchase amount, what product categories the purchases were in (and for what dollar amounts), demographics for the consumers (allowing us to see if certain individuals are more prone to clumpy behavior), and marketing activity at the individual-level (across three different advertising channels) and six other data sets are utilized to lay out the business problems from a company's point of view that clumpiness suggests. In particular, a better capturing of clumpiness might enable firms to better estimate customer future value, thus enabling superior targeting, customer acquisition, and segmentation. In addition, when a clumpy customer returns, he or she might make a lot of purchases and becomes very profitable. As a result, a good measure of clumpiness enables companies to deliver special promotions to reactivate the right customers. From the investor perspective, clumpiness means high growth potential but large risk, thus a clumpiness measure might serve as an important factor that should be taken into account when firms make investment decisions.

Based on this discussion, we posit that clumpiness (C) should become the next building block to profile customers and an additional important driver/predictor of CLV and its subcomponents (churn, customer frequency, and customer monetary value). Thus, as an overarching goal, we intend to bring greater attention to clumpiness, which should be considered as an important feature besides RFM to both researchers and practitioners as a measure of key strategic importance. Also, by showing the pervasiveness of clumpiness across data sets, and its predictive power for customer value even after controlling for RFM and marketing activity, we are able to show that it is not an esoteric or isolated problem, and that it contains potential monetary value to the firm. To summarize, the empirical contributions of our paper include the following:

1. C is a significant predictor of out-of-sample individual-level customer value, even conditional on

R, F, M, and marketing activity; hence our title from RFM to RFMC.

2. Computing C utilizing purchase data (purchase clumpiness, purchase-C) and visitation data (visit clumpiness, visit-C) are nonredundant and differentially predictive statistics in their own right.

3. Segmentation based on RFMC provides additional insights above and beyond RFM alone in identifying future customer activity.

4. Models that assume stationary behavior misestimate the potential  $P(\text{alive})$  for clumpy customers as the distribution conditional on clumpiness is bimodal.

5. C is related to demographics and varies by purchase category.

6. Marketing activity is (modestly) correlated with C; hence an opportunity exists for future research to assess firms' ability to maximize profitability using the link between marketing activity and C as a strategic tool.

Although our empirical findings are limited to the degree that they are based on the data sets we obtained and analyzed, they clearly indicate the need to study clumpiness from a much deeper perspective.

The remainder of this article is organized as follows. In §2, we begin with a brief review of clumpiness measures that allows us to extend RFM to RFMC. In §3, we present our empirical analysis framework that describes the many ways in which we utilize C as both an independent and dependent variable. In §4, we conduct a detailed empirical study of our data that provides strong evidence in support of the importance of clumpiness when valuing customers. Furthermore, we address a subtle issue (portended above), which is that clumpiness can be either visitation based (visit-C) or purchase based (purchase-C), and demonstrate their nonredundancy and differential prediction power. In addition, §4 contains a descriptive analysis that looks at the relationship between clumpiness, demographics, and marketing variables. Finally in §5 we conclude with a discussion of several additional open issues that arise from this work. To summarize, for many current digital consumables (where consumption may be clumpy), we believe that this research raises questions about the viability of RFM as a sufficient form of data compression.

## 2. Clumpiness Measures

To begin to explicate our RFMC framework, we need to describe a procedure, as simple as RFM, to compute C (clumpiness). Although growing attention has been placed on the increasing pattern of clumpy data in many empirical areas, there is very little research on constructing a well-defined and careful measure that could act as a foundation for our metric-based approach. The related hot hand effect has long been a widespread belief in sports, and has triggered a

branch of interesting research that could shed some light on this domain (Gilovich et al. 1985; Tversky and Gilovich 1989, 2005; Bar-Eli et al. 2006). However, many concerns have been raised about the low power of existing hot hand significance tests (Dorsey-Palmateer and Smith 2004, Miyoshi 2000, Wardrop 1999, Frame et al. 2003). To better capture clumpiness, Zhang et al. (2013) provides a rich discussion of this issue and proposes a new class of clumpiness measures that are shown to have higher statistical power in extensive simulations under a wide variety of data generating mechanisms (alternative hypotheses as compared to the null of stationarity). Although details are provided in Zhang et al. (2013), we note a few desirable properties<sup>1</sup> about these measures.

- *Minimum*. The measure should be the minimum if the events are equally spaced.
- *Maximum*. The measure should be the maximum if all of the events are gathered together.
- *Continuity*. Shifting event times by a very small amount should only change the measure by a small amount.
- *Convergence*. As events move closer (further apart), the measure should increase (decrease).

The four properties above make intuitive sense, and also provide a comprehensive description of the measure's dynamics. As shown in Zhang et al. (2013), any convex and symmetric function of intervisit times (IETs) for visit-C, or interpurchase times for purchase-C, satisfies all of the desired properties. We describe the method using intervisit times for ease of explication, and hence visit-C, except when otherwise mentioned; but an identical method is used for interpurchase times to compute purchase-C.

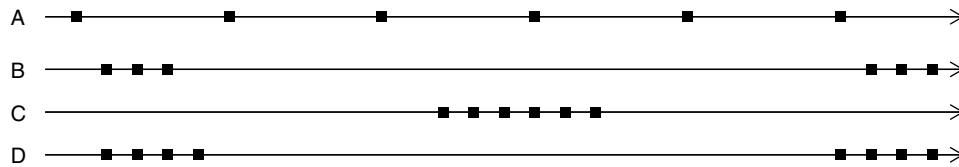
For a chosen clumpiness measure, a formal statistical test is then needed to determine the significance of an observed clumpiness value. The null hypothesis is random sampling without replacement, where  $n$  (the number of events) and  $N$  (the number of trials) are known, and Monte Carlo simulation is applied to compute the Z-table, the table of clumpiness critical values. Once the Z-table is generated, a test of clumpiness can be implemented: i.e., the null hypothesis of randomness is rejected and the individual customer is judged to be clumpy at the  $\alpha$ -level when the clumpiness measure is larger than the corresponding  $\alpha$ -level clumpiness critical value. For the case that  $n$  is much smaller than  $N$ , a useful Poisson approximation is available for those critical values.

For all analyses in this research, we chose  $H_p$ , an individual-level entropy measure (metric) from the class of convex and symmetric IET statistics, as the

<sup>1</sup> These four properties hold when C is computed for a fixed number of time periods ( $N$ ) and a fixed number of events ( $n$ ). When  $n$  or  $N$  change, we discuss later the properties of C.



Figure 1 A Few Examples: How Clumpiness Is Calculated



measure of clumpiness, due to its slightly improved and robust performance shown in Zhang et al. (2013).<sup>2</sup> The expression for  $H_p$  is given by

$$H_p = 1 + \frac{\sum_{i=1}^{n+1} \log(x_i) \cdot x_i}{\log(n+1)}, \quad (1)$$

where, as mentioned,  $n$  and  $N$  denote the number of events and the number of trials (potential events) for one sequence of incidence data in a given observation period, and  $t_i$  and  $x_i$  are the  $i_{th}$  occurrence of event time and IETs, respectively.<sup>3</sup> Note that in Equation (1), we divide each IET,  $x_i$ , by  $N+1$  to control for the length of the observation period.

To better illustrate the computation and features of our clumpiness measure, let us consider four hypothetical customers who make purchases among  $N=30$  time periods. Their transaction histories are displayed in Figure 1, where the black boxes indicate the occurrence of purchases. Take customer A for example who made purchases in time periods 1, 6, 11, 16, 21, and 26, so  $x_1 = 1$ , and  $x_2, \dots, x_7 = 5$  and  $n=6$ . After scaling,  $H_p = 1 + [(1/31) \cdot \log(1/31) + (6 \cdot 5/31) \cdot \log(5/31)] / \log(6+1) = 0.036$ . Since it is less than the corresponding critical value  $z(6, 30) = 0.257$ , customer A is classified as nonclumpy at the 5% significance level.

This measure not only allows us to categorize customers in terms of clumpiness but also enables us to compare the degree of clumpiness across individuals. For example, although customers B and C are both judged as clumpy, B's measure  $H_p = 0.48$  is greater than C's of  $H_p = 0.34$ , and since  $n$  and  $N$  are identical, B is "unambiguously" more clumpy than C. We put unambiguously in quotes because although it is statistically unambiguous when  $n$  and  $N$  are the same to compare C (larger C is more clumpy), to the eye (unlike  $R$  and  $F$ ), it may not be as trivial. When C is applied to two customers with different  $n$  and  $N$ , one could compare  $p$ -values (smaller  $p$ -value if more clumpy) or use the empirical distribution of clumpiness within an  $n$  and  $N$  bucket (larger percentile C within your group) to norm/order the values. Further details of the computation of clumpiness, and its associated statistical test are included in Appendices A and B for

those looking to fully replicate our findings and apply it to their own data sets.

Finally, we note that since the clumpiness measure is computed by scaled interevent times (i.e., dividing by  $N+1$ ), it is conceptually invariant to the scaling of the time units. Empirically, however, it is not always true (i.e., the computed value will not be identical) because of the discreteness and nonlinearity in  $H_p$ . In particular, if the selected time unit is too long, data points are aggregated and information about the IETs would be lost. On the other hand, a too short time unit would give rise to possibly over-fit patterns. Hence, the choice of  $t$  in practice should be matched to the scaling of a business's decision time periodicity (e.g., weekly, monthly, or quarterly). Furthermore, since computing C is trivial and fast, one could do this for multiple scalings and check the robustness. With this measure in hand, we can now extend RFM to RFMC.

### 3. RFMC Empirical Analysis Framework for Studying Clumpiness

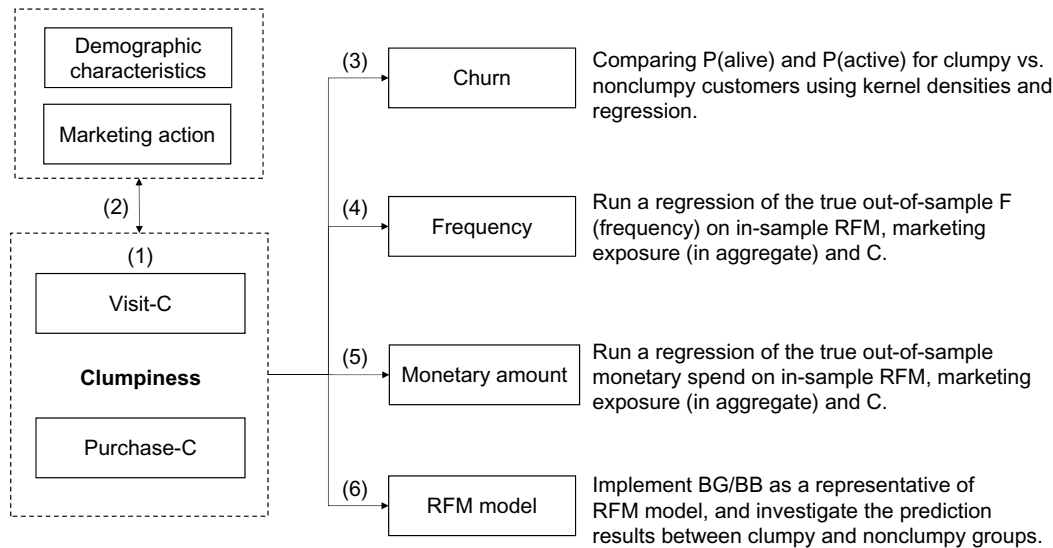
Before applying our clumpiness measures to the data, we provide Figure 2 that lays out how we have empirically studied clumpiness in a number of different ways. First, as labeled (1) in Figure 2, we consider visit-C and purchase-C as conceptually distinct but computationally identical (see §2). That is, from a firm's point of view, clumpy visitation may indicate heightened periods of search, whereas clumpy purchasing may indicate heightened periods of consumer affinity. Second, labeled (2) and described more fully in §4, we explore the degree to which clumpiness varies systematically across people (as described by their demographic characteristics) and marketing action (described by firm expenditures). Although certainly not causal, to the degree that this allows a firm to understand their customer base better or to influence clumpiness through differential expenditure can be managerially important. Finally, Figure 2 highlights the four major empirical contributions of this paper, which we next describe in more detail.

Labeled (3) in Figure 2, we study the use of C as an X variable in understanding customer churn. This is done by first computing  $P(\text{alive})$  for customers under a stationary transaction model (the BG/BB model) and comparing the values for those people deemed to be clumpy versus not. In this manner, we will demonstrate the bimodality for clumpy people (i.e., they

<sup>2</sup> We also applied the other clumpiness measures as in Zhang et al. (2013), and there is no significant difference.

<sup>3</sup> For  $i=2, \dots, n$ ,  $x_i = t_i - t_{i-1}$ ;  $x_1 = t_1$  and  $x_{n+1} = N+1 - t_n$ . It follows that  $\sum_{i=1}^{n+1} x_i = N+1$ .

**Figure 2** Analysis Framework Utilizing C as a Predictor of CLV and Its Subcomponents, and Its Relation to Marketing Action and Demographics Utilizing C as a Dependent Variable



may be dead or alive, and it is uncertain) versus the unimodal nature for the nonclumpy customers. To further our understanding of the relationship between  $C$  and churn, we also compare observed out-of-sample activity versus predicted  $P(\text{active})$  for clumpy versus nonclumpy customers and demonstrate the increased activity out-of-sample for the clumpy ones. Analysis component (4) highlights the use of  $C$  as a predictor of frequency ( $F$ ) by regressing out-of-sample  $F$  on in-sample  $C$ ,  $R$ ,  $F$ ,  $M$ , and marketing expenditure and demonstrating the statistically significant predictive power of  $C$ . Component (5) continues the exploration of  $C$  as a predictor of CLV components by using out-of-sample monetary value ( $M$ ) as the  $Y$  variable in a regression with (as in (4)) in-sample  $C$ ,  $R$ ,  $F$ ,  $M$ , and marketing action as predictors. Last, we compare in component (6) the prediction errors for a stationary probability model across clumpy and nonclumpy customers. This analysis demonstrates the systematic prediction errors for certain types of customers when they exist in the data and provides researchers and managers some insights on when stationary models are most likely to be sufficient for analyzing their data set *before* doing any formal analyses.

With this framework in hand, we describe next our application of this analysis framework systematically through the six analysis components.

## 4. Empirical Application of $C$ to Multiple Data Sets

We illustrate our framework from Figure 2 by its application to seven different data sets that span a number of different industry verticals and data complexity levels commonly seen in current marketing applications.

Specifically, for only one of our data sets do we have marketing, demographics, and purchasing information (as compared to visitation), and hence we provide in this research partial empirical evidence for Figure 2 for the distinctness of visit and purchase  $C$  (component (1)), the role of marketing activity (component (2)), and the role of  $C$  in predicting monetary value (component (5)). For all other analyses that are based on incidence data, we can complete them for all seven data sets. Next, we provide a set of analyses on our most complete data set and follow that with a number of analyses on all of our obtained data sets.

### 4.1. Distinctness of Visit and Purchase Based Clumpiness

To explore the distinctness between visit and purchase-based  $C$ , we analyzed a comprehensive data set from a large North American retailer that prefers to remain anonymous. This data set contains a cohort of 42,000 randomly selected customers that tracked both their online and offline interaction with the company over two years. The data includes their online website visitation on a daily basis as well as the purchases made by each of these customers, and what product categories the purchases were in (and for what dollar amounts). Thus, for our purposes, this allows us to compute both a  $42,000 \times 730$  (daily) visitation and purchase incidence matrix on which visit- $C$  and purchase- $C$  will be built.

To begin our application of Figure 2, we applied our  $H_p$  measure (for both visit- $C$  and purchase- $C$ ) to each of the 42,000 customers. Approximately 70% of the population based on their online visitation are categorized as visit-clumpy,<sup>4</sup> suggesting that this is

<sup>4</sup> Here  $\alpha = 0.05$  was chosen as the probability threshold to determine clumpiness.

**Table 1** Visit-C versus Purchase-C for 42,000 Customers of a Large North American Retailer

Visit/purchase	Nonclumpy (%)	Clumpy (%)	Total (%)
Nonclumpy	26	4	30
Clumpy	54	16	70
Total	80	20	100

a widespread phenomenon for this data, which we later explore more generally for other data sets. If we use the interpurchase time to calculate purchase-C, the clumpy population only constitutes 20% of the cohort. Thus, at least for this data set, purchase-based clumpiness is distinct but correlated with ( $p < 0.001$ ) visitation-based clumpiness as displayed in Table 1 and Appendix D. The relationship between purchase-based clumpiness and visitation-based clumpiness in different types of data sets is an interesting area for future study.

#### 4.2. Clumpiness and Its Relationship to Demographics

For the same large North American retailer data set as in §4.1, demographics for the consumers (allowing us to see if certain individuals are more prone to clumpy behavior) are included. A summary of the data is shown in Table 2. As we can see, the customers have an average age of 39, spend an average \$125 when they purchase, and come roughly once every other week.

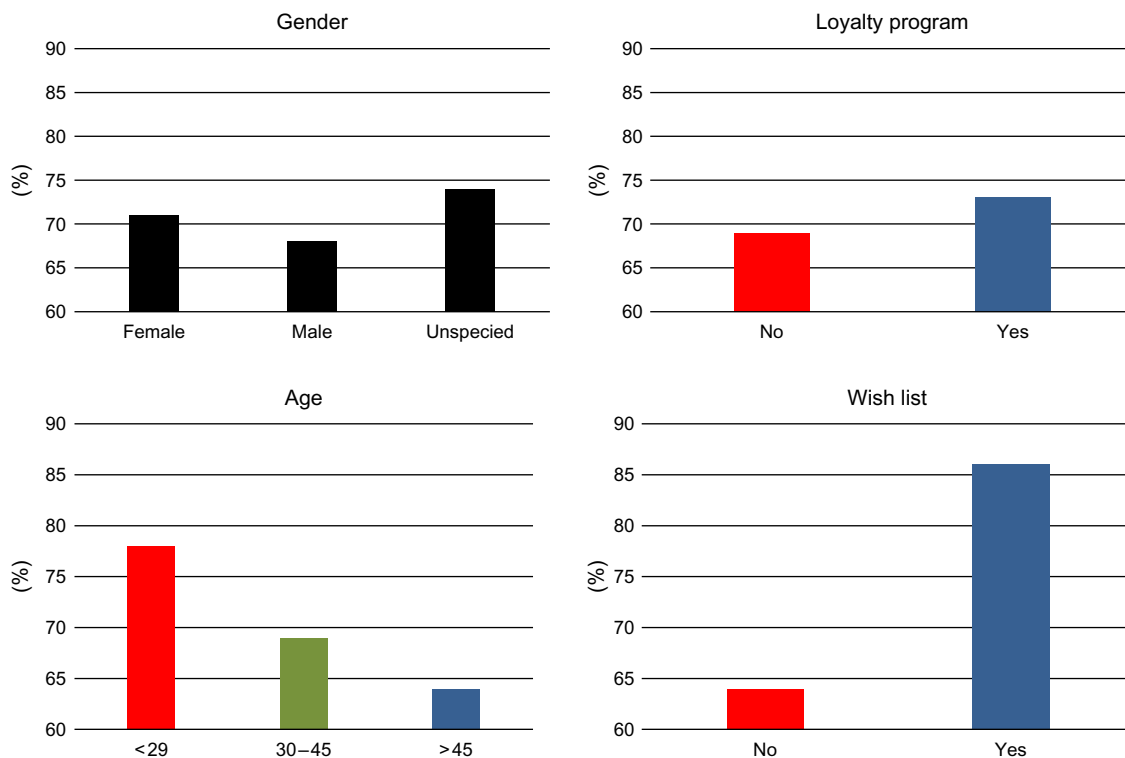
**Table 2** Demographics and RFM Summary of 42,000 Customers from a Large North American Retailer

Variable	Summary
Gender	F: 76%, M: 11%, U: 13%
Age	Mean: 39, SD: 14
Loyalty program	Y: 36%, N: 64%
Wish list	Y: 17%, N: 83%
Recency	Mean: 432, SD: 146
Frequency	Mean: 50, SD: 83
Monetary amount	Mean: 125, SD: 102

We first look at how visit-C is related to personal characteristics (purchase-C results are available upon request). As shown in Figure 3, we find weakly significant gender differences in clumpy behavior ( $p < 0.1$ , women are more clumpy than men), customers who are older are less clumpy ( $p < 0.01$ ), and customers who are loyalty program customers of the retailer or who have a wish list with the firm (i.e., a wedding registry list) are more clumpy. Although these findings may not be immediately monetizable by the firm, they do provide insights into observable characteristics that the firm can use for C-based targeting.

#### 4.3. Clumpiness and Marketing Exposure

To descriptively explore the relationship between C and marketing exposure in a bidirectional way (Figure 2, point (2)), we report a correlation matrix between marketing activity and visit and purchase-based C (Table 3).

**Figure 3** (Color online) Percentage of Visit-Based Clumpy Customers at  $\alpha = 0.05$  Level Broken Down by Demographics

**Table 3** A Correlation Matrix Between Marketing Activity and Visit and Purchase-Based C

	Visit-C	Purchase-C	Catalog	Direct mail	Email
Visit-C	1.00	0.14	0.01	0.02	0.08
Purchase-C	0.14	1.00	−0.01	0.01	0.01
Catalog	0.01	−0.01	1.00	0.25	0.49
Direct mail	0.02	0.01	0.25	1.00	0.28
Email	0.08	0.01	0.49	0.28	1.00

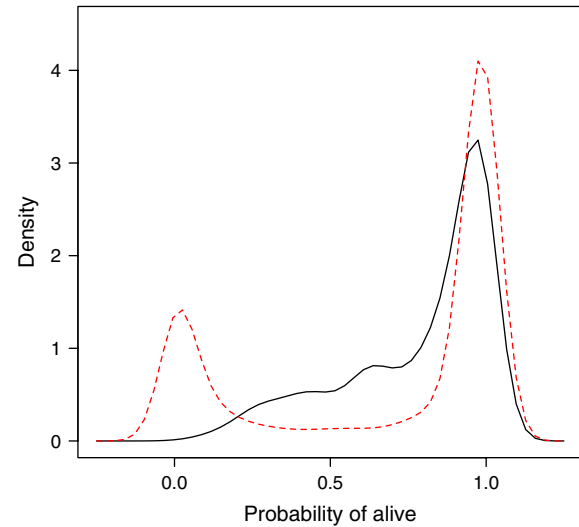
As shown, only quantity of email sent to an individual has a significant impact ( $p < 0.05$ ) on visit-C,<sup>5</sup> which suggests a few areas for future study. First, email (among email, catalog, and direct mail) is known to be a more short-term marketing tool (Chandon et al. 2000) and hence its increased correlation with C has some empirical validity based on past studies. Second, since we observe that email is correlated with visit-C, but not purchase-C, this suggests an opportunity for the firm to translate visit patterns to purchase.

#### 4.4. Clumpiness and Its Relationship to $P(\text{alive})$ and $P(0)$

One main tenet of this paper, as discussed, is to explore the relationship between clumpiness and components of CLV (as per Figure 2, components (3)–(5)). In this subsection, we begin that exploration by assessing the differences (if they exist) between the  $P(\text{alive})$  of customers who are designated as clumpy versus those who are not. To this end, we utilized the first 18 months of data from the large North American retailer as our in-sample calibration period, and held out the last six months for validation. We operationalized our churn benchmark (albeit there are numerous choices including  $P(0)$  below) by computing  $P(\text{alive})$  at the end of the observation period assuming a BG/BB model (Fader et al. 2010) of behavior (in which a customer buys at a constant rate until he churns, see computational details in Appendix C). We plot the density of  $P(\text{alive})$  (using a kernel density) segmented by visit-C clumpy and visit-C nonclumpy customers at the  $\alpha = 0.05$  level in Figure 4.

The density of  $P(\text{alive})$  among the clumpy group is clearly bimodal, which demonstrates the large variation within this group and the verbiage that “clumpy customers may be dead or they may be dormant” waiting to return—i.e., two-humped with respect to churn. This demonstrates one dimension by which RF-based models may not provide a complete characterization of a customer’s history; they do not buy at constant propensity while alive, hence periods of inactivity are erroneously attributed to death. Content-based reasons

**Figure 4** (Color online) Distribution of  $P(\text{alive})$  Between Nonclumpy and Clumpy Groups (Kernel Density)



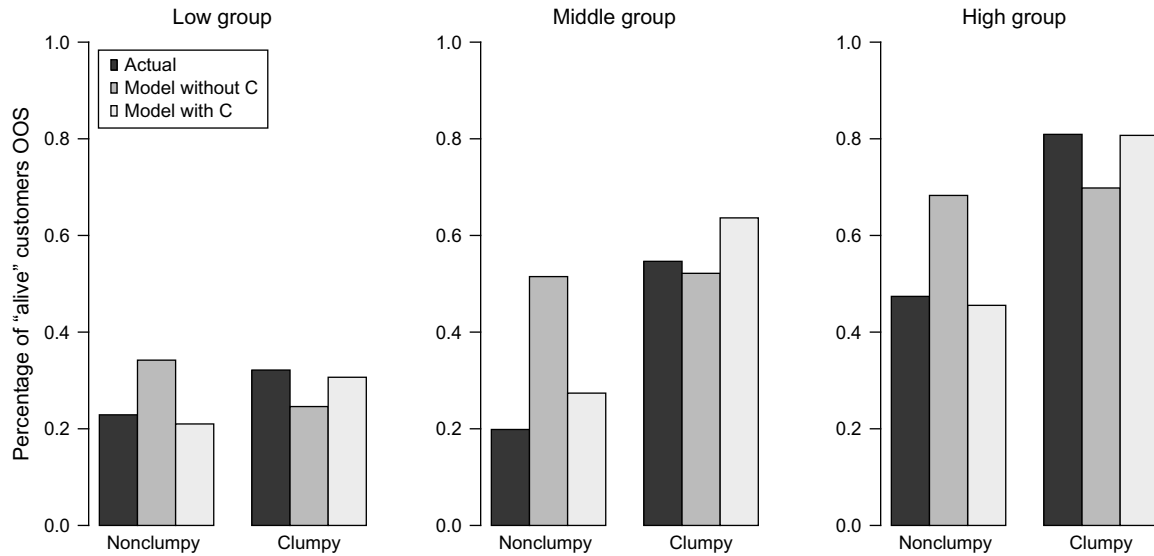
for such bimodality (e.g., a new season of shows is available at a digital content provider) is a rich area for future study but beyond the scope of this paper. Yet, as a prospectus for future researchers we might imagine that goal completion (Kivetz et al. 2006) and unintended consumption (Stilley et al. 2010) would be two process-level stories worth pursuing.

Although the previous analyses provide some evidence about the impact of C on churn, we wanted to provide an additional set of exploratory results. In this vein, we now demonstrate that including C (as a predictor variable) helps one reduce the type I and II errors for classifying a customer as alive or not, a crucial component of CLV models. Specifically, we ran a logistic regression using TRUE out-of-sample incidence (1 if any observation occurred during the six-month validation period and hence the customer is alive, 0 otherwise) as the dependent variable and BG/BB-based  $P(0)$ , probability of no observations in the out-of-sample period, and observed C as regressors, and compare that to only using  $P(0)$ , to assess C’s extra predictive power for out-of-sample activity (if it exists). See Appendix C for computational details on computing  $P(0)$ . The results are displayed in Figures 5 and 6.

Figures 5 and 6 show three important findings. First, as we can see, out-of-sample those customers with higher in-sample clumpiness return significantly more often than nonclumpy customers. Second, this pattern widens as the amount of in-sample clumpiness increases from low to medium to high. Last, a model that includes in-sample C in its regressors to predict  $P(\text{activity})$  out-of-sample significantly outperforms a model that excludes it. Each of these findings support the heavier return pattern of clumpy customers beyond that predicted by a stationary (nonclumpy) model.

<sup>5</sup> This correlation analysis was supplemented by a logistic regression of clumpiness status as a dependent variable against marketing activity by channel. As with Table 3, email volume is a significant predictor of visit-C ( $p < 0.05$ ).



**Figure 5** A Comparison of Actual versus Predicted Active Customers Out-of-Sample from a Model With and Without C for Low, Medium, and High Clumpy Customers

#### 4.5. Clumpiness and Out-of-Sample Frequency

To further illustrate that clumpiness can be valuable as a component of CLV predictions, we also show that C adds extra predictive power for out-of-sample frequency and purchasing (component (4), Figure 2) even after conditioning on existing RFM and marketing exposure. In particular, we ran two regressions of the TRUE out-of-sample visit and purchase-based F (frequency) on in-sample RFM, marketing exposure (in aggregate) and C for the large North American retailer data set. All of the variables have a positive effect, but importantly visit-C based clumpy customers are predicted to make more purchases ( $p < 0.01$ ) in terms of frequency. The result is similar if using purchase-based C (except that

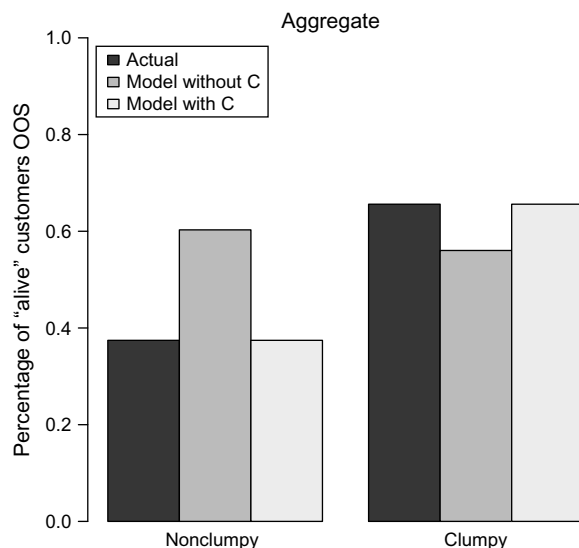
**Table 4** In-Sample RFMC and MKTG on Out-of-Sample Frequency

	Visit-C	Purchase-C
R	0.002 (0.00)	0.002 (0.00)
F	0.017 (0.00)	0.017 (0.00)
M	0.001 (0.00)	0.001 (0.00)
C	0.235 (0.00)	0.998 (0.00)
MKTG	0.008 (0.00)	0.008 (0.00)

*Notes.* Reported are two sets of regression coefficients and  $p$ -values. The first column uses visit-C as the dependent variable; the second column is the one with purchase-C.

the effect is much stronger,  $p < 0.01$  when comparing the two regression coefficients<sup>6</sup>), and is shown in Table 4. It is important for us to note the interesting (and central to this research) finding that purchasing behavior (purchase-C) predicts out-of-sample purchasing more significantly than visit-based C and that both visit and purchase-based C are significant predictors of out-of-sample frequency even after controlling for RFM and marketing.

To illustrate that our findings are not unique to this retailer, we extended this analysis for frequency to six other companies where we were able to obtain data: two traditional online businesses (CDNow, Mecoxlane), and four large Internet companies (Hulu, YouTube, Amazon, and eBay). The data for the two traditional businesses and Hulu.com were obtained from the Wharton Customer Analytics Initiative (WCAI).<sup>7</sup> CDNow is one of

**Figure 6** Aggregate Comparison of Actual versus Predicted Active Customers Out-of-Sample Using a Model With and Without C

<sup>6</sup> The increase in  $R$ -squared is 7% by adding purchase-C to the regression and 2% by adding visit-C.

<sup>7</sup> The Wharton Customer Analytics Initiative (WCAI) is the pre-eminent academic research center based in the Wharton School's Marketing Department, focusing on the development and application of customer analytics methods.

**Table 5** Data Summary for Additional Applications

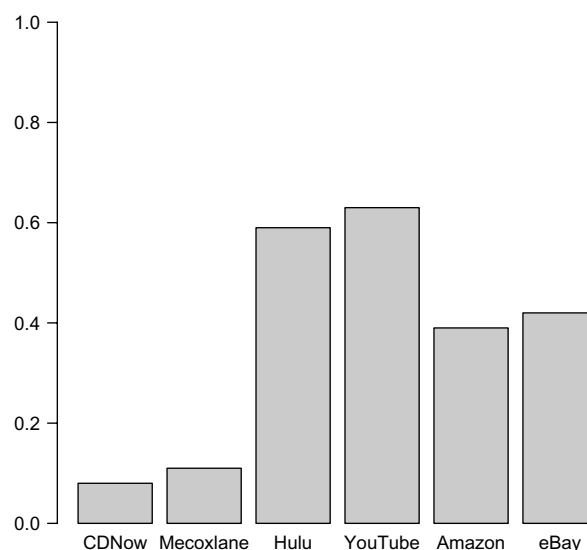
Company	Start	End	Type	Sample size
CDNow	January 1997	June 1998	Daily	500
Mecoxlane	March 2009	April 2010	Weekly	180
Hulu	January 2010	April 2010	Daily	1,000
YouTube	January 2011	December 2011	Daily	5,000
Amazon	January 2011	December 2011	Daily	5,000
eBay	January 2011	December 2011	Daily	5,000

the oldest and largest online retailers, having sold different forms of music (and related products) on the World Wide Web since 1994, and Mecoxlane ([www.m18.com](http://www.m18.com)) is one of China's leading online platforms for apparel and accessories, which has experienced substantial growth since its inception in 1996, and was listed on the NASDAQ exchange in 2010. Hulu.com is an online video content company that is one of the top two online video content properties in the United States (the other is YouTube). For the other three popular Internet companies, we used the data from the comScore Web Behavior Database on WRDS (Wharton Research Data Services; <https://wrds-web.wharton.upenn.edu/wrds/>), which offers web-wide visitation and transaction behavior based on a random sample from a cross-section of more than two million Internet users in the United States. We randomly selected 5,000 customers for each company, tracked their daily browsing history in the 2011 calendar year, and generated a daily incidence matrix of individuals-by-days, where the entries are 1s or 0s indicating whether that individual visited the website on that day. A summary of the data sets is provided in Table 5. Similar to the retailer data, we used the first three quarters of data for calibration to compute our clumpiness measure  $H_p$ , and held out the last quarter of data to compute out-of-sample measures noting again for these six data sets that the analyses are based on visit-C only.

Before reporting the regression results as per Table 4, we report the percentage of clumpy customers for each company in Figure 7 and the out-of-sample frequency for clumpy and nonclumpy customers in Figure 8 to portend the likely regression findings. It is easy to see that within the two traditional companies, there are roughly only 10% of customers identified to be clumpy. In other words, customers' behavior is relatively stable/stationary in these traditional businesses. Conversely, the clumpiness phenomena is widely prevalent in the four Internet companies.

Furthermore, when comparing the results across firms in Figure 8, we observe that for the two traditional firms, the two clumpy and nonclumpy groups do not differ significantly. On the other hand, for the four Internet companies, it is clear, that the clumpy group has a larger group average and standard deviation,

**Figure 7** Percentage of Clumpiness Customers



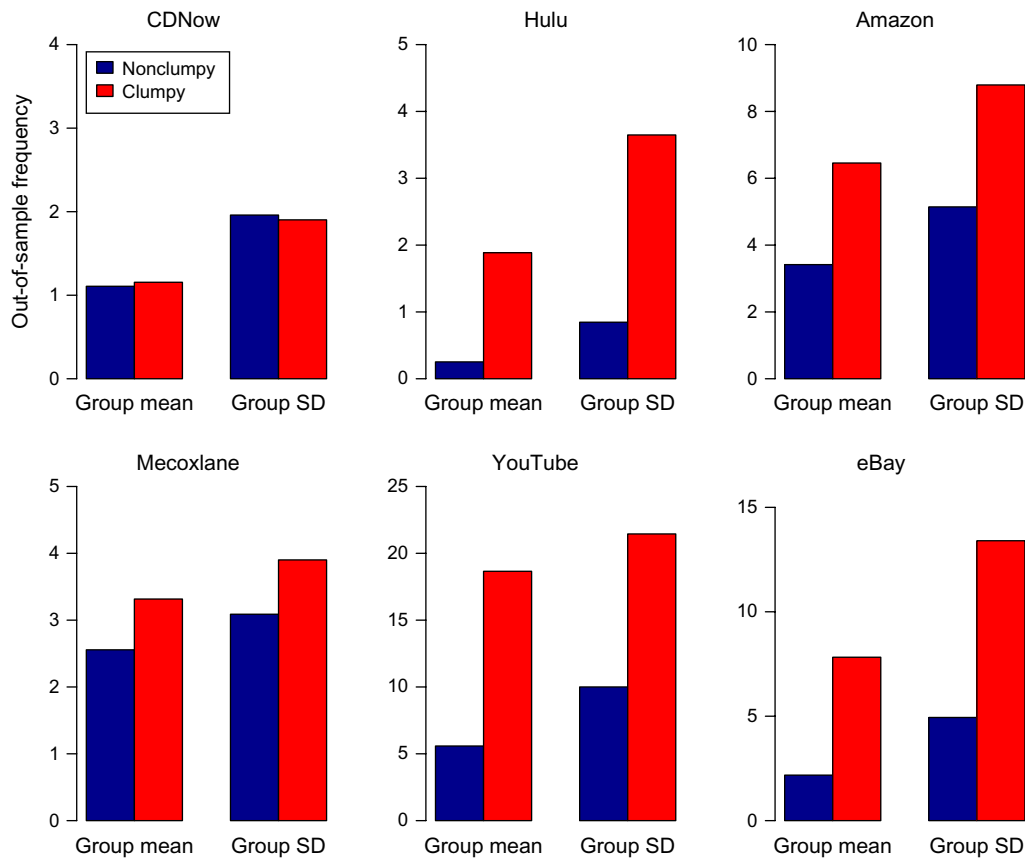
which is consistent with our findings from the large North American retailer. The two results combined suggest that many companies may have high potential and high-risk customers in a way that is different from traditional companies with regular buying patterns. Whether this pattern holds in general is an excellent empirically based research project for the future. However, it does portend why the BG/BB and RFM models that were historically built on traditional data sets may have room for improvement in the digital/Internet space.

We now extend the key regression analyses of the retailer data set to all six data sets to demonstrate the broad predictive power of  $C$  for out-of-sample frequency. In particular, since these six data sets are visitation only without purchase, demographics, and marketing action (as described), we ran a regression of the TRUE out-of-sample  $F$  (frequency) on in-sample  $RF$  and  $C$ . The coefficients and corresponding  $p$ -values are displayed in Table 6. Except for the two traditional firms, visit- $C$  has a significantly positive impact and predictive value for out-of-sample frequency supporting the central tenet of this paper and confirming Table 4's results more broadly.

Although the increase in  $R$ -squared by adding  $C$  is modest (consistently between 5%–7% for Amazon, YouTube, Hulu, and eBay), the results of Figure 8 and the last row of Table 6 suggest that  $C$ 's additional predictive relationship with frequency and what moderates it is an important area for future study.

#### 4.6. Clumpiness and Monetary Value

In the previous section, we explored the role of  $C$  in predicting frequency; we now return to the large North American retailer data set to explore how  $C$  impacts dollar amount and customer value (component (5) of

**Figure 8** (Color online) Comparison of Out-of-Sample Frequency Between Nonclumpy and Clumpy Groups

analysis plan, Figure 2). We ran a regression of the TRUE out-of-sample total purchase dollar amount on in-sample RFM, aggregate marketing exposure, and in-sample C, using both visit-C and purchase-C. Not unexpectedly, as per marketing practice, all of in-sample R, F, M, and marketing activity have significantly positive effects on out-of-sample value. However, the findings around visit-C and purchase-C are more subtle. In particular, whereas visit-C is not significant in predicting out-of-sample dollar volume, purchase-C, even after controlling for R, F, M, and MKTG, is highly significant. This separation between C on F (Table 4 in the previous section) and C on dollar volume (Table 7), and differences between visitation-based and purchase-based C, is a potentially important managerially relevant finding, suggesting visit-C drives frequency (Table 4), not dollars (Table 7), but purchase-C really drives dollars.

Last, we looked at how the impact of purchase-based C on the dollar amount varies across product categories. That is, we have looked at clumpiness across seven different product subcategories (intimates, home and furniture, kids, women's clothing, women's accessories, men's accessories, men's clothing) provided by the retailer. Our findings suggest over a 2:1 ratio of variation of clumpiness across categories with the highest being women's clothing and home goods (15%) and the lowest being kids (7%). For six (out of seven) categories, purchase-C is predictive of purchase dollar amount out-of-sample, indicated by a Y in Table 8, even after controlling for R, F, M, and MKTG. This suggests that an understanding of the context specific nature of clumpiness is a future vital area of research and could be done in a more systematic manner by taking an extensive set of product categories (e.g., the IRI data set, Bronnenberg et al. 2008) and a set of

**Table 6** In-Sample RFC on Out-of-Sample Frequency

	CDNow	Mecoxlane	Amazon	YouTube	Hulu	eBay
R	0.002 (0.00)	0.040 (0.00)	0.014 (0.00)	0.037 (0.00)	0.012 (0.00)	0.028 (0.00)
F	0.286 (0.00)	0.237 (0.00)	0.413 (0.00)	0.382 (0.00)	0.116 (0.00)	0.397 (0.00)
C	-0.188 (0.14)	0.763 (0.17)	0.243 (0.00)	2.521 (0.00)	0.561 (0.00)	0.622 (0.00)

Note. Reported are regression coefficients and *p*-values.

**Table 7** In-Sample RFMC and MKTG on Out-of-Sample Dollar Amounts

	Visit-C	Purchase-C
R	0.32 (0.00)	0.30 (0.00)
F	2.51 (0.00)	2.48 (0.00)
M	1.87 (0.00)	1.87 (0.00)
C	−5.08 (0.74)	101.96 (0.00)
MKTG	1.40 (0.00)	1.40 (0.00)

Notes. Reported are two sets of regression coefficients and  $p$ -values. The first column is the one with visit-C; the second column is the one with purchase-C.

**Table 8** Purchase Clumpiness for Seven Different Product Subcategories, and Whether It Is Predictive ( $Y/N$ ) of Out-of-Sample Customer Value at the  $\alpha = 0.05$  Level

Home	Intimates	Kids	Men	Men's accessories	Women	Women's accessories
15%	8%	7%	13%	12%	15%	12%
Y	Y	N	Y	Y	Y	Y

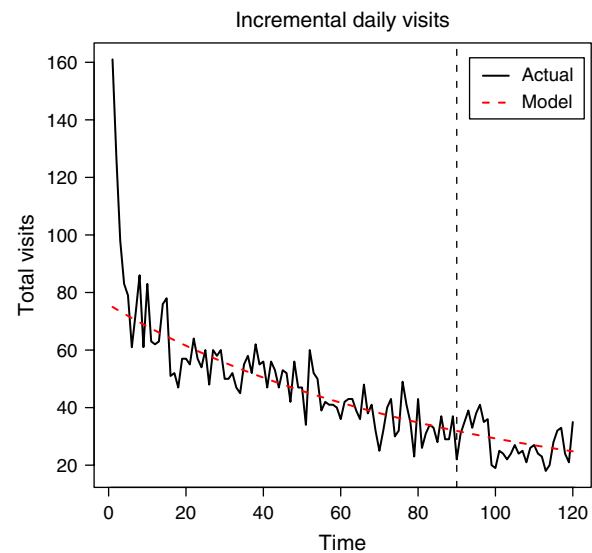
product category characteristics (e.g., Macé and Neslin 2004) and doing a systematic study.

Although Table 8 demonstrates a significant predictive power of purchase-C on out-of-sample purchase dollars, in Table 9 we show that the amounts (in real dollar terms) are quite significant. We report model-free evidence reporting both the purchase amounts and the percentage difference between clumpy and nonclumpy buyers.

#### 4.7. RFM Models: Aggregate versus Disaggregate Fit

We conclude this section by taking a deep look at the impact of deviation from a simple RFM story on customer valuation as per component (6) of Figure 2. We use our Hulu data set as an example (its application to the other six data sets are available upon request), and utilized the BG/BB (as we did with the retailer data) as a benchmark model (Fader et al. 2010). In other words, we use it as a representative of RF-based models in the following analysis; but we note that our results would be robust to other choices within the RF-class. In Figure 9, we compute the expected number of visits for the cohort of customers for each of the 120 days, with the vertical dotted line indicating the end of our chosen in-sample period at 90 days. The BG/BB model predictions track the actual (incremental) trajectory very well in both the calibration period (except the very beginning) and the validation period; a common (but possibly not complete) way to assess fit.

**Figure 9** (Color online) Aggregate Fit of BG/BB Model to Hulu.com Data



Despite the strong indication of good overall fit, we want to provide here a more nuanced view of model fit. It is true that the BG/BB provides accurate predictions overall, but what about at the segment level? That is, the BG/BB or other RFM models are commonly used to identify the most profitable cohort of customers to the firm, but can it equally be used to look at individual-level customer value that is also of interest. In particular, even if the fit is good overall, if it gets the tail of the distribution wrong, it will make severe CLV-component errors. Thus we compute the average of out-of-sample frequency under a BG/BB model (actual and model prediction) within nine  $R$  and  $F$  groups, created via equal one-third splits. The results are displayed in Table 10.

The last line of the table shows that the model only underestimates out-of-sample frequency by 4% in total, consistent with extant research. However, the prediction errors are significant within many subgroups. Take the cohort of customers with high  $F$  and low  $R$ , for example, for which there is significant prediction error. Why is this the case? As discussed before, the BG/BB tells a buy-till-you-die story, so a low  $R$  coupled with high  $F$  means high predicted probability of attrition. However, it might not be true; those customers might just be clumpy and mistakenly seen as inactive. Instead, they might be dormant but will return and continue generating visits. As per Figure 9, can we find a way to

**Table 9** Purchase Amounts and Percentage Difference Between Nonclumpy and Clumpy Groups in Out-of-Sample Purchase Amounts by Product Subcategory

	Home	Intimates	Kids	Men	Men's accessories	Women	Women's accessories
Nonclumpy (\$)	38	32	0.17	23	9	235	56
Clumpy (\$)	73	95	0	36	18	335	96
Percent change (%)	92	197	−100	57	100	43	71



**Table 10** A Summary of Prediction Errors of BG/BB Model, by RF Segmentation for Hulu.com Data

Nodes		Number of users	Out-of-sample			
Frequency	Recency		Actual	Estimated	Error	Percentage (%)
L	L	116	18	5.4	12.6	70
L	M	76	14	33.4	−19.4	−139
L	H	32	19	29.2	−10.2	−54
M	L	45	15	0.5	14.5	97
M	M	62	28	31.8	−3.8	−13
M	H	38	67	64.8	2.2	3
H	L	8	18	0	−18.0	100
H	M	34	36	12.1	23.9	66
H	H	97	400	464.7	−64.7	−16
Total		508	615	642	−27	−4

help differentiate people between clumpy and inactive? Our clumpiness measure exactly evaluates and informs about this separation. To investigate this, an aggregate set of box-and-whisker plots is employed where we look at prediction errors. The results are shown in Figure 10.

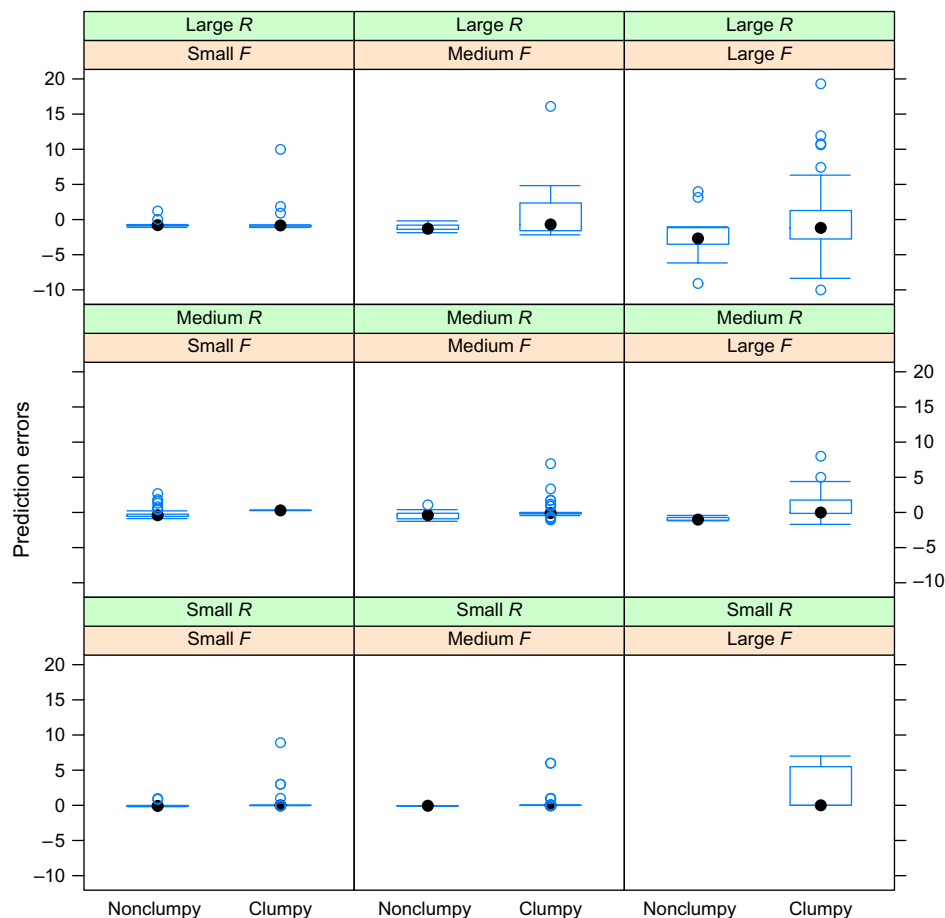
It is clear that clumpy groups have larger prediction errors using the BG/BB in general. In other words, a buy-till-you-die story performs well for nonclumpy

customers, but not for clumpy ones. Looking at people with low  $R$  specifically, the major underestimation of this cohort results from clumpy people who actually return. Summarizing, since clumpiness can be computed ahead of any formal model fitting, researchers may find  $C$  a useful tool to determine needed model complexity in advance.

## 5. Conclusion and Future Research

The move toward a customer-centric approach to marketing, coupled with the increasing availability of customer transaction data, has led to an interest in both the notion and the calculation of customer lifetime value. RFM is a mainstay of estimating CLV for many businesses.

The key result of our research is that besides RFM, clumpiness is also an important driver of profiling customers and estimating CLV. We illustrate it by conducting a detailed empirical study of a data set from a North American retailer. From a methodological perspective, we show that even though RFM-based models are able to provide good overall estimates of in-sample and out-of-sample frequency in some cases, it will lead to significant prediction errors at the individual level if

**Figure 10** (Color online) Comparison of Prediction Errors Between Nonclumpy and Clumpy Groups Broken Down by  $R$  and  $F$ 

clumpiness is not captured and it exists. We also extend the analysis to a variety of data from six companies: two traditional and four Internet. This demonstrates that the clumpiness phenomena is widely prevalent on the Internet or at least is worth exploring. From the standpoint of managerial implementation, clumpy customers mean high potential, high risk. Through properly managing and delivering the right promotion, they can lead to huge value to the company in a way that might be missed otherwise.

Future marketing research, based on our findings here, would benefit from (1) better or more statistically powerful measures to capture clumpiness; (2) better statistical models that capture the degree of clumpiness; (3) studies that combine a complete set of demographic and marketing-mix variables to study the drivers of clumpiness including field experiments that will allow firms to assess the link between marketing and C, and then marketing and firm outcomes; and (4) applying the clumpiness measures to more data sets to understand empirical generalizations. Hopefully, this research is a good first step.

#### Appendix A. Step-by-Step Instruction to Compute $H_p$

1. Convert the individual-level transaction data into incidence/binary data, if necessary.
2. Compute the Inter-event times (IETs).

$$x_i = \begin{cases} t_1, & \text{if } i = 1, \\ t_i - t_{i-1}, & \text{if } i = 2, \dots, n, \\ N + 1 - t_n, & \text{if } i = n + 1, \end{cases}$$

where  $t_i$  and  $x_i$  are the  $i_{th}$  occurrence of event time and IETs, respectively, and  $n$  and  $N$  represent the number of visits and total time intervals.

3. Rescale the IETs. Divide IETs by  $N + 1$ .
4. Compute the normalized entropy-like

$$H_p: 1 + \frac{\sum_{i=1}^{n+1} \log(x_i) \cdot x_i}{\log(n+1)}.$$

#### Appendix B. Test of Clumpiness

For a chosen clumpiness measure ( $H_p$  is used in our paper), a formal statistical test is needed to determine the significance of clumpiness. Larger values than would be expected under a model of randomness then provides an indication of clumpiness, with smaller values than expected indicating stability. Although the null hypothesis is random sampling without replacement, where the sample size ( $n$ ) and population size ( $N$ ) are known, the null distribution for the new clumpiness measures cannot be generally derived in closed form. To facilitate our analysis, Monte Carlo simulation is applied to compute the Z-table, the table of clumpiness critical values. The pseudo-code is presented as follows for the general case, where  $n$  and  $N$  denote the number of events and the number of trials for one sequence of incidence data in a given observation period:

1. Initialize the iteration number  $M$  and level of significance  $\alpha$ .
2. Given  $n$  and  $N$ , for  $m$  in 1 to  $M$ :
  - (a) Take a sample of size  $n$  from  $N$  days without replacement.

- (b) Calculate the clumpiness measure of the random sample.

3. Find the  $\alpha$ -percentile of the computed measures in the random sample as the critical value.

Once the Z-table is generated, a test of clumpiness can be implemented for each sequence of data. The null hypothesis of randomness is rejected and this sequence is judged to be clumpy when the clumpiness measure is larger than the corresponding critical value.

For the case that  $n$  is much smaller than  $N$ , a useful approximation is available for those critical values. Using the well-known law of rare events, a Poisson distribution can be used as a good approximation of the binomial distribution if  $n$  is sufficiently large and  $p$  is sufficiently small. Similarly, the Dirichlet distribution is the limiting case of the distribution of IETs in the random draws without replacement (under the Bernoulli distribution conditional on the total number of events). There is a rule of thumb stating that the approximation result is excellent if  $N \geq 20$  and  $n/N \leq 0.1$ .

#### Appendix C. Key Results of BG/BB Model

BG/BB is a benchmark model for buyer behavior in a discrete-time, noncontractual setting. We present some key results below. (The associated derivations can be found in Fader et al. 2010.)

Let the random variable  $X(n) = \sum_{i=1}^n Y_i$  denote the number of transactions occurring across the first  $n$  transaction opportunities

$$E(X(n) | \alpha, \beta, \gamma, \delta) = \frac{\alpha}{\alpha + \beta} \frac{1}{\delta - 1} \left\{ \delta - \frac{\Gamma(\gamma + \delta)\Gamma(n + \delta + 1)}{\Gamma(\delta)\Gamma(\gamma + \delta + n)} \right\}.$$

More generally, let the random variable  $X(n, n + n^*) = \sum_{i=n+1}^{n+n^*} Y_i$  denote the number of transactions in the interval  $(n, n + n^*)$ .  $P(0)$  can be generated by setting  $x^*$  equal to 0

$$E(X(n, n + n^*) | \alpha, \beta, \gamma, \delta) = \frac{\alpha}{\alpha + \beta} \frac{1}{\delta - 1} \left\{ \frac{\Gamma(\gamma + \delta)\Gamma(n + \delta + 1)}{\Gamma(\delta)\Gamma(\gamma + \delta + n)} - \frac{\Gamma(\gamma + \delta)\Gamma(n + n^* + \delta + 1)}{\Gamma(\delta)\Gamma(\gamma + \delta + n + n^*)} \right\},$$

$$P(X(n, n + n^*)) = x^* | \alpha, \beta, \gamma, \delta = \delta_{x^*=0} \left\{ 1 - \frac{B(\gamma, \delta + n)}{B(\gamma, \delta)} \right\} + \binom{n^*}{x^*} \frac{B(\alpha + x^*, \beta + n^* - x^*)}{B(\alpha, \beta)} \frac{B(\gamma, \delta + n + n^*)}{B(\gamma, \delta)} + \sum_{i=x^*}^{n^*-1} \binom{i}{x^*} \frac{B(\alpha + x^*, \beta + i - x^*)}{B(\alpha, \beta)} \frac{B(\gamma + 1, \delta + n + 1)}{B(\gamma, \delta)}.$$

In most customer-based analysis settings we are interested in making statements about customers conditional on their observed purchase history  $(x, t_x, n)$

$$P(\text{alive at } n+1 | \alpha, \beta, \gamma, \delta, x, t_x, n) = \frac{B(\alpha + x, \beta + n - x)}{B(\alpha, \beta)} \frac{B(\gamma, \delta + n + 1)}{B(\gamma, \delta)} \cdot \frac{1}{L(\alpha, \beta, \gamma, \delta | x, t_x, n)},$$

$$\begin{aligned}
& E(X(n, n+n^*) | \alpha, \beta, \gamma, \delta, x, t_x, n) \\
&= \frac{1}{L(\alpha, \beta, \gamma, \delta | x, t_x, n)} \frac{B(\alpha+x+1, \beta+n-x)}{B(\alpha, \beta)} \\
&\quad \cdot \frac{\Gamma(\gamma+\delta)}{(\gamma-1)\Gamma(\delta)} \left\{ \frac{\Gamma(n+\delta+1)}{\Gamma(\gamma+\delta+n)} - \frac{\Gamma(n+n^*+\delta+1)}{\Gamma(\gamma+\delta+n+n^*)} \right\}, \\
& \text{DERT}(d | \alpha, \beta, \gamma, \delta, x, t_x, n) \\
&= \frac{B(\alpha+x+1, \beta+n-x)}{B(\alpha, \beta)} \frac{B(\gamma, \delta+n+1)}{B(\gamma, \delta)(1+d)} \\
&\quad \cdot \frac{{}_2F_1(1, \delta+n+1; \gamma+\delta+n+1; 1/(1+d))}{L(\alpha, \beta, \gamma, \delta | x, t_x, n)}.
\end{aligned}$$

where  ${}_2F_1()$  is the Gaussian hypergeometric function and  $d$  is the discount rate.

#### Appendix D. Correlation Matrix of RFMC

	R	F	M	Visit-C	Purchase-C
R	1.00	0.33	0.05	0.24	0.08
F	0.33	1.00	0.02	0.28	0.06
M	0.05	0.02	1.00	0.01	0.01
Visit-C	0.24	0.28	0.01	1.00	0.14
Purchase-C	0.08	0.06	0.01	0.14	1.00

#### References

- Baesens B, Viaene S, Van den Poel D, Vanthienen J, Dedene G (2002) Bayesian neural network learning for repeat purchase modelling in direct marketing. *Eur. J. Oper. Res.* 138(1):191–211.
- Bar-Eli M, Avugos S, Raab M (2006) Twenty years of hot hand research: Review and critique. *Psych. Sport Exercise* 7(6): 525–553.
- Berry MJA, Linoff GS (2004) *Data Mining Techniques: For Marketing, Sales, and Customer Relationship Management* (Wiley, Indianapolis).
- Blattberg R, Kim B, Neslin S (2008) *Database Marketing: Analyzing and Managing Customers*, Vol. 18 (Springer Verlag, New York).
- Bolton R (1998) A dynamic model of the duration of the customer's relationship with a continuous service provider: The role of satisfaction. *Marketing Sci.* 17(1):45–65.
- Bronnenberg BJ, Kruger MW, Mela CF (2008) Database paper the IRI marketing data set. *Marketing Sci.* 27(4):745–748.
- Chandon P, Wansink B, Laurent G (2000) A benefit congruency framework of sales promotion effectiveness. *J. Marketing* 64(4): 65–81.
- Davenport TH, Harris J, Shapiro J (2010) Competing on talent analytics. *Harvard Bus. Rev.* 88(10):52–58.
- Dorsey-Palmateer R, Smith G (2004) Bowlers hot hands. *Amer. Statistician* 58(1):38–45.
- Fader P (2013) *Customer Centricity: Focus on the Right Customers for Strategic Advantage* (Wharton Executive Essentials, Philadelphia).
- Fader PS, Bruce GSH, Ka LL (2005) RFM and CLV: Using iso-value curves for customer base analysis. *J. Marketing Res.* 42(4): 415–430.
- Fader PS, Hardie BGS, Shang J (2010) Customer-base analysis in a discrete-time noncontractual setting. *Marketing Sci.* 29(6): 1086–1108.
- Frame D, Hughson E, Leach JC (2003) Runs, regimes, and rationality: The hot hand strikes back. Working paper, Leeds School of Business, University of Colorado, Boulder.
- Gilovich T, Vallone R, Tversky A (1985) The hot hand in basketball: On the misperception of random sequences. *Cognitive Psych.* 17(3):295–314.
- Gupta S, Lehmann D, Stuart J (2004) Valuing customers. *J. Marketing Res.* 41(11):7–18.
- Gupta S, Hanssens D, Hardie B, Kahn W, Kumar V, Lin N, Ravishanker N, Sriram S (2006) Modeling customer lifetime value. *J. Service Res.* 9(2):139–155.
- Kivetz R, Urminsky O, Zheng Y (2006) The goal-gradient hypothesis resurrected: Purchase acceleration, illusionary goal progress, and customer retention. *J. Marketing Res.* 43(1):39–58.
- Kumar V (2006) CLV: A path to higher profitability. Technical report, Working paper, University of Connecticut, Storrs.
- Kumar V (2008) *Customer Lifetime Value: The Path to Profitability* (Now Publishers, Boston).
- Kumar V, Venkatesan R, Bohling T, Beckmann D (2008) The power of CLV: Managing customer lifetime value at IBM. *Marketing Sci.* 27(4):585–599.
- Macé S, Neslin SA (2004) The determinants of pre- and postpromotion dips in sales of frequently purchased goods. *J. Marketing Res.* 41(3):339–350.
- Malthouse E (2003) Scoring models. *Kellogg on Integrated Marketing* (John Wiley & Sons, Hoboken, NJ), 227–249.
- Malthouse E, Blattberg R (2005) Can we predict customer lifetime value? *J. Interactive Marketing* 19(1):2–16.
- Miyoshi H (2000) Is the hot-hands phenomenon a misperception of random events? *Japanese Psych. Res.* 42(2):128–133.
- Netzer O, Lattin J, Srinivasan V (2008) A hidden Markov model of customer relationship dynamics. *Marketing Sci.* 27(2):185–204.
- Parr Rud O (2001) *Data Mining Cookbook: Modeling Data for Marketing, Risk, and Customer Relationship Management* (Wiley, New York).
- Schmittlein D, Morrison D, Colombo R (1987) Counting your customers: Who are they and what will they do next? *Management Sci.* 33(1):1–24.
- Schweidel D, Bradlow E, Fader P (2009) Portfolio dynamics for customers of a multiservice provider. *Management Sci.* 57(3): 471–486.
- Stilley KM, Inman JJ, Wakefield KL (2010) Planning to make unplanned purchases? The role of in-store slack in budget deviation. *J. Consumer Res.* 37(2):264–278.
- Tversky A, Gilovich T (1989) The “hot hand”: Statistical reality or cognitive illusion? *Chance* 2(4):31–34.
- Tversky A, Gilovich T (2005) The cold facts about the “hot hand” in basketball. *Anthology Statist. Sports* 16:169.
- Wardrop R (1999) Statistical tests for the hot-hand in basketball in a controlled setting. *Amer. Statistician* 1(1):1–20.
- Zhang Y, Bradlow ET, Small DS (2013) New measures of clumpiness for incidence data. *J. Appl. Statist.* 40(11):2533–2548.

#### CORRECTION

In this article, “Predicting Customer Value Using Clumpiness: From RFM to RFMC” by Yao Zhang, Eric T. Bradlow, and Dylan S. Small (first published in *Articles in Advance*, September 2, 2014, *Marketing Science*, DOI:10.1287/mksc.2014.0873), the equation in point 4 of Appendix A was corrected as follows:

$$H_p: 1 + \frac{\sum_{i=1}^{n+1} \log(x_i) \cdot x_i}{\log(n+1)}.$$