

Design and Evaluation of Optimal Free Trials

Hema Yoganarasimhan *
University of Washington

Ebrahim Barzegary*
University of Washington

Abhishek Pani*
Bright Machines

November 15, 2021

Abstract

Free trial promotions are a commonly used customer acquisition strategy in the Software as a Service (SaaS) industry. We use data from a large-scale field experiment to study the effect of trial length on customer-level outcomes. We find that the 7-days trial is the best average treatment that maximizes customer acquisition, retention, and profitability. In terms of mechanism, we rule out the demand cannibalization theory, find support for the consumer learning hypothesis, and show that long stretches of inactivity at the end of the trial are associated with lower conversions. We then develop a framework for personalized targeting policy design and evaluation. We first learn a lasso model of outcomes as a function of users' pre-treatment variables and treatment. Next, we use individual-level predictions of the outcome to assign the optimal treatment to each user. We then evaluate the personalized policy using the inverse propensity score reward estimator. We find that a personalization based on lasso leads to 6.8% improvement in subscription compared to a uniform 30-days for all policy. It also performs well on long-term customer retention and revenues in our setting. Segmentation analysis suggests that skilled and experienced users are more likely to benefit from longer trials. Finally, we show that personalized policies do not always outperform uniform policies, and one should be careful when designing and evaluating personalized policies. In our setting, personalized policies based on other outcomes and heterogeneous-treatment effects, estimators (e.g., causal forests, random forests) perform worse than a simple 7-days for all policy.

Keywords: free trials, targeting, personalization, policy evaluation, field experiment, machine learning, digital marketing, Software as a Service

*We are grateful to an anonymous firm for providing the data and to UW-Foster High-Performance Computing Lab for providing us with computing resources. We thank the participants of the 2018 Marketing Science conference and the Triennial Choice Symposium, and 2021 UT Dallas Bass Conference. Thanks are also due to seminar audiences at the Emory University, Johns Hopkins University, University of California Berkeley, University of Houston, University of Maryland, University of Rochester, University of South Carolina, and University of Southern California, for their feedback. Please address all correspondence to: hemay@uw.edu.

1 Introduction

1.1 SaaS Business Model and Free Trials

Over the last few years, one of the big trends in the software industry has been the migration of software firms from the perpetual licensing business model to the “Software as a Service” (SaaS) model. In the SaaS model, the software is sold as a service, i.e., consumers can subscribe to the software based on monthly or annual contracts. Global revenues for the SaaS industry now exceed 200 billion USD (Gartner, 2019). This shift in the business model has fundamentally changed the marketing and promotional activities of software firms. In particular, it has allowed firms to leverage a new type of customer acquisition strategy: free trial promotions, where new users get a limited time to try the software for free.

Free trials are now almost universal in the SaaS industry because software is inherently *experience good*, and free trials allow consumers to try the software product without risk. However, we do not have a good understanding of how long these trials should be or the exact mechanism through which they work. In the industry, we observe trial lengths ranging anywhere from one week to three months; e.g., Microsoft 365 offers a 30 days free trial, whereas Google’s G Suite offers a 14 days free trial. There are pros and cons associated with both long and short trials. A short trial period is less likely to lead to free-riding or demand cannibalization and is associated with lower acquisition costs. On the other hand, an extended trial period can enhance consumer learning by giving consumers more time to learn about product features and functionalities. Longer trials can also create stickiness/engagement and increase switching-back costs. That said, if users do not use the product more with a longer trial, they are more likely to conclude that the product is not useful or forget about it. Thus, longer trials lack the deadline or urgency effect (Zhu et al., 2018).

While the above arguments make a global case for shorter/longer trials, the exact mechanism at work and the magnitude of its effect can be heterogeneous across consumers. In principle, if there is significant heterogeneity in consumers’ response to the length of free trials, SaaS firms may benefit from assigning each consumer a different trial length depending on her/his demographics and skills. The idea of personalizing the length of free trial promotions is akin to third-degree price discrimination because we effectively offer different prices to different consumers over a fixed period. Indeed, SaaS free trials are particularly well-suited to personalization because of a few reasons. First, software services have zero marginal costs, and there are no direct cost implications of offering different trial lengths to different consumers. Second, it is easy to implement a personalized free trial policy at scale for digital services, unlike physical products. Finally, consumers are less likely to react adversely to receiving different trial lengths (unlike prices). However, it is not clear whether personalizing the length of free trials improves customer acquisition and firm revenues, and if yes, what is the best approach to design and evaluate personalized free trials.

1.2 Research Agenda and Challenges

In this paper, we are interested in understanding the role of trial length on customer acquisition and profitability for digital experience goods. We focus on the following research questions. First, does the length of a free trial promotion affect customer acquisition, and if so, what is the ideal trial length? Second, what is the mechanism

through which trial length affects conversions? Third, is there heterogeneity in users' responsiveness to trial lengths? If yes, how can we personalize the assignment of trial lengths based on users' demographics and skills, and what are the gains from doing so? Further, what types of customers benefit from shorter trials vs. longer trials? Finally, how do personalized targeting policies that maximize short-run outcomes (i.e., customer acquisition) perform on long-run metrics such as consumer retention and revenue?

We face three main challenges in answering these questions. First, from a data perspective, we need a setting where trial length assignment is exogenous to user attributes. Further, to understand the mechanism through which trial length affects conversions, we need to observe usage and activity during the trial period. Second, from a methodological perspective, we need a framework to design and evaluate personalized targeting policies in high-dimensional settings and identify the optimal policy. A series of recent papers at the intersection of machine learning and causal inference provide heterogeneous treatment effects estimators, which we can use to personalize treatment assignment (Athey and Imbens, 2016; Wager and Athey, 2018). Similarly, a series of papers in marketing have combined powerful predictive machine learning models with experimental (or quasi-experimental) data to develop personalized targeting policies (Rafieian and Yoganarasimhan, 2021; Simester et al., 2020). However, the optimal policy that each of these papers/methods arrive at in a given empirical context can differ. Thus far, we have little to no understanding of how these methods compare in their ability to design effective targeting policies. This brings us to the third challenge. We need to be able to evaluate the performance of each policy *offline* (without deploying it in the field). Evaluation is essential because deploying a policy in the field to estimate its effectiveness is costly in time and money. Moreover, given the size of the policy space, it is simply not feasible to test each policy in the field.

1.3 Our Approach and Findings

To overcome these challenges and answer our research questions, we combine a three-pronged framework to design and evaluate personalized targeting policies with data from a large-scale free trial experiment conducted by a major SaaS firm. The firm sells a suite of related software products (e.g., Microsoft 365, Google G Suite) and is the leading player in its category, with close to monopolistic market power. At the time of this study, the firm used to give users a 30-days free trial for each of its software products, during which they had unlimited access to the software suite. Then, the firm conducted a large-scale field experiment, where new users who started a free trial for one of the firm's products were randomly assigned to one of 7, 14, or 30-days trial length conditions. It also monitored the subscription and retention decisions of the users in the experiment for two years. The firm also collected data on users' pre-treatment characteristics (e.g., skill level and job) and post-treatment product usage during the trial period.

First, we quantify the average treatment effect of trial length on subscription. We find that the firm can do significantly better by simply assigning the 7-days trial to all consumers (which is the best uniform policy). This leads to a 5.59% gain in subscriptions over the baseline of 30 days for all policy in the test data. In contrast, the 14-days for all policy does not significantly increase subscriptions. This finding suggests that simply shortening the trial length to 7 days will lead to higher subscriptions. At the time of the experiment, the firm offered a standard 30-day free trial to all its consumers. So better performance of the much shorter

7-day trial was surprising, especially since the reasons proposed in the analytical literature for the efficacy of free trials mostly support longer trials, e.g., switching costs, consumer learning, software complexity, and signaling. (See $\chi 2$ for a detailed discussion of the analytical literature on free trials). Therefore, we next examine the mechanism through which trial length affects conversion and present some evidence for why a shorter trial works better in this setting and examine the generalizability of these results. To that end, we leverage the usage data during the trial period to understand the mechanism through which trial length affects subscriptions. We show that there are two opposing effects of trial length. On the one hand, as trial length increases, product usage and consumer learning about the software increases. This increase in usage can have a positive effect on subscriptions. On the other hand, as trial length increases, the gap between the last active day and the end of the trial increases, while the average number of active days and usage per day reduces. These factors are associated with lower subscriptions. In our case, the latter effect dominates the former, and shorter trials are better.

Our analysis presents three key findings relevant to the theories on the role of free trials for experience goods. First, we rule out the demand cannibalization or free riding hypothesis advocated by many theoretical papers by showing that users who use the product more during the trial are more likely to subscribe (Cheng and Liu, 2012). Second, we provide empirical support for the consumer learning hypothesis, since we show that longer trials lead to more usage, which in turn is associated with higher subscriptions (Dey et al., 2013). Third, we identify a novel mechanism that plays a significant role in the effectiveness of free trials – the negative effect of long stretches of inactivity at the end of the trial on subscription.

Next, we develop a two-step approach to personalized policy design since an unstructured search for the optimal policy is not feasible in our high-dimensional setting. In the first stage, we learn a lasso model of outcomes (subscription) as a function of the users’ pre-treatment demographic variables and their trial length. Then in the second stage, we use the individual-level predictions of the outcome to assign the optimal treatment for each user. Then, we use the Inverse Propensity Score (IPS) reward estimator, popular in the counterfactual policy evaluation literature in computer science, for offline policy evaluation (Horvitz and Thompson, 1952; Dudík et al., 2011).

Based on this approach, we show that the personalized free trial policy leads to over 6.8% improvement in subscription compared to the baseline uniform policy of giving a 30-day trial for all. That said, the magnitude of gains from personalization (over the best uniform policy of 7 days for all) are modest (which is in line with the recent findings on personalization of marketing interventions in digital settings; e.g., Rafeian and Yoganarasimhan (2021)). Further, we find that customers’ experience and skill level affect their usage, which affects their subscription patterns. Beginners and inexperienced users show only a small increase in usage with longer trial periods. Further, when given longer trials, they end up with long periods of inactivity at the end of the trial period, which negatively affects their likelihood of subscribing. Thus, it is better to give them short trials. In contrast, long trials are better for experienced users because it allows them to use the software more and they are not as negatively influenced by periods of inactivity later in the trial period. Overall, our findings suggest that simpler products and experienced users are more likely to benefit from longer trials.

Next, we find that the personalized policy, designed to optimize subscriptions, also performs well on long-term metrics, with a 7.96% increase in customer retention (as measured by subscription length) and 11.61% increase in revenues. We also consider two alternative personalized policies designed to maximize subscription length and revenues and compare their performance with that of the subscription-optimal policy. Interestingly, we find that the subscription-optimal policy always performs the best, even on long-run outcomes. While this finding is specific to this context, it nevertheless shows that optimizing low-variance intermediate outcomes (i.e., statistical surrogates) can be revenue- or loyalty-optimal in some settings.

Finally, we consider counterfactual policies based on four other outcome estimators: (1) linear regression, (2) CART, (3) random forests, and (4) XGBoost, and two heterogeneous treatment effect estimators: (1) causal tree, and (2) generalized random forests. We find our lasso-based personalized policy continues to perform the best, followed by the policy based on XGBoost (6.17% improvement). However, policies based on other outcome estimators (e.g., random forests, regressions) perform poorly. Interestingly, policies based on the recently developed heterogeneous treatment effects estimators (causal tree and causal forest) also perform poorly. Causal tree is unable to personalize the policy at all. Causal forest personalizes policy by a small amount, but the gains from doing so are marginal. While our findings are specific to this context, it nevertheless suggests that naively using these methods to develop personalized targeting policies can lead to sub-optimal outcomes. This is particularly important since these methods are gaining traction in the marketing literature and are being used without evaluation using off-policy methods; see for example Guo et al. (2017) and Fong et al. (2019).

Our research makes three main contributions to the literature. First, from a substantive perspective, we present the first empirical study that establishes the causal effect of trial length on conversions and provides insight into the mechanisms at play. Second, from a methodological perspective, we present a framework that managers and researchers can use to design and evaluate personalized targeting strategies applicable to a broad range of marketing interventions. Finally, from a managerial perspective, we show that the policies designed to optimize short-run conversions also perform well on long-run outcomes in our setting, and may be worth considering in other similar settings. Importantly, managers should recognize that many popular estimators can give rise to poorly designed personalized policies, which are no better than simple uniform policies. Offline policy evaluation is thus a critical step before implementing any policy.

2 Related Literature

Our paper relates to the research that examines the effectiveness of free trials on the purchase of experience goods, especially digital and software products. Analytical papers in this area have proposed a multitude of theories capturing the pros and cons of offering free trials. Mechanisms such as switching costs, network effects, quality signaling, and consumer learning are often proposed as reasons for offering free trials. In contrast, free-riding and demand cannibalization are offered as reasons against offering free trials. See Cheng and Liu (2012), Dey et al. (2013), and Wang and Özkan-Seely (2018) for further details. In spite of this rich theory literature, very few empirical papers have examined whether and how free trials work in practice. In an early paper, Scott (1976) uses a small field experiment to examine if users given a two-week free

trial are more likely to purchase a newspaper subscription. Interestingly, she finds that free trials do not lead to more subscriptions compared to the control condition. While the number of participants may not have been sufficient to detect small effects and the context was very different from digital SaaS products, it nevertheless raises the question of whether free trials can be an effective marketing strategy. More recently, two empirical papers study free trials using observational data. Foubert and Gijsbrechts (2016) build a model of consumer learning and show that while free trials can enhance adoption, ill-timed free trials can also suppress adoption. Using a bayesian learning approach, Sunada (2018) compares the profitability of different free trial configurations. However, neither of these papers examines how trial length affects subscriptions/revenues because they lack variation in the length of the free trials in their data. In contrast, we use data from a large-scale field experiment with exogenous variation in the length of free trials to identify the optimal trial length for each user. In addition, we contribute to this literature by leveraging the individual-level software usage data during the trial period to rule out some of the earlier theories proposed in this context, e.g., free riding. To the best of our knowledge, our paper provides the first comprehensive empirical analysis of how trial length affects the purchase of digital experience goods.

Second, our paper relates to the marketing literature on real-time customization and personalization of digital products and promotions using machine learning methods. This literature has used a wide range of methods for the personalization tasks in a variety of contexts: website design using dynamic programming and adaptive experiments (Hauser et al., 2009), display ads using multi-arm bandits (Schwartz et al., 2017), ranking of search engine results using feature engineering, and boosted trees (Yoganarasimhan, 2020), mobile ads using behavioral and contextual features (Rafieian and Yoganarasimhan, 2021), and the sequence of ads shown in mobile apps using batch reinforcement learning and optimal dynamic auctions (Rafieian, 2019a,b). We add to this literature in two ways. First, we document the gains from personalizing the duration of a new type of promotional strategy: the length of time-limited free trials for digital experience goods using a targeting framework based on data from a large-scale field experiment. Second, we show that while personalization can help, it may not always be the case. Indeed, in our setting, many commonly used methods for personalization often perform worse than a robust uniform policy based on average treatment effects. While these findings are specific to our context, it nevertheless suggests that managers should be careful in designing and evaluating personalized targeting policies.

Our paper also relates to the theoretical and empirical research on personalized policy design and evaluation in computer science and economics. In an early theoretical paper, Manski (2004) presents a method that finds the optimal treatment for each observation by minimizing a regret function. Recent theoretical papers in this area include Swaminathan and Joachims (2015) and Swaminathan et al. (2017), Kitagawa and Tetenov (2018), and Athey and Wager (2017). There is also a small but growing list of marketing papers in this area. Hitsch and Misra (2018) propose a heterogeneous treatment effects estimator based on kNN, develop targeting policies based on it, and evaluate the performance of their policies using the IPS estimator on a test data. However, their estimator does not work in our setting because it requires all the covariates to be continuous since it based on Euclidean distance. Simester et al. (2020) examine how

managers can evaluate targeting policies efficiently. They compare two types of randomization approaches: (a) randomization by action and (b) randomization by policy. They provide two valuable insights. First, they note that randomization by action is preferable to randomization by policy because it allows us use off policy evaluation to evaluate any policy. Second, they note that when comparing two policies we should recognize that if the policies recommend the same action for some customers then the difference in the performance of the policy for these customers is exactly zero. In another particularly relevant paper, Simester et al. (2019) investigate how data from field experiments can be used to design targeting policies for new customers or new regimes, and also use the IPS estimator to evaluate the performance of a series of personalized policies. They present comparisons for a broad range of methods and show that model-based methods in general (and lasso in particular) offers the best performance, though this advantage vanishes if the setting and/or consumers change significantly. Our paper also echoes this finding: the lasso-based personalized policy performs the best in our setting too. Further, we also provide comparisons to personalized policies based on the newly proposed heterogeneous treatment effects estimators (e.g., causal forest), and show that the lasso-based policy continues to perform the best.

Our paper is relevant to the literature on statistical surrogates (Prentice, 1989; VanderWeele, 2013). In our setting, subscription can be interpreted as an intermediate outcome or surrogate for long-run retention and revenue. Interestingly, we find that personalized policies optimized on the short-term outcome or surrogate do well (or better than) policies optimized directly on the long-term outcomes. We attribute this to the fact that long-term outcomes have higher variance and fewer observations in our setting. In a recent paper, Yang et al. (2020) use surrogates to impute the missing long-term outcomes and then use the imputed long-term outcomes to develop targeting policies. Their results confirm our broader finding that short-term outcomes can be sufficient to derive targeting policies that are optimal from a long-run perspective.

More broadly, our work relates to the large stream of marketing literature that has examined and contrasted the short vs. long run effects of promotions (Mela et al., 1997; Pauwels et al., 2002). The main takeaway from this literature is that consumers who are exposed to frequent price promotions become price sensitive and engage in forward buying over time. While these early papers focused on consumer packaged goods, Anderson and Simester (2004) conduct a field experiment on price promotions in the context of durable goods sold through catalogs. They find evidence in support of both forward-buying and increased deal sensitivity. Our paper adds to this literature by examining the long-run effect of free-trial promotions on long-run subscription and revenue for digital SaaS products. While free-trial promotions can be viewed as a price discount (i.e., zero price for a fixed period), forward-buying is not feasible in our setting and consumers are exposed to the promotion only during the sign-up period (i.e., no expectation of future free trials). In this case, we find that targeted free-trial promotions that maximize short-run revenue (or subscriptions) also perform well on long-run outcomes (two-year revenue).

3 Setting and Data

In this section, we describe our application setting and data.

3.1 Setting

Our data come from a major SaaS firm that sells a suite of software products. The suite includes a set of related software products (similar to Excel, Word, PowerPoint in Microsoft’s MS-Office). The firm is the leading player in its category, with close to monopolistic market power. Users can either subscribe to single-product plans that allow them access to one software product or to bundled plans that allow them to use several products at the same time. Bundles are designed to target specific segments of consumers and consist of a set of complementary products. The prices of the plans vary significantly and depend on the bundle, the type of subscription (regular or educational), and the length of commitment (monthly or annual). Standard subscriptions run from \$30 to \$140 per month depending on the products in the bundle and come with a monthly renewal option. (To preserve the firm’s anonymity, we have multiplied all the dollar values in the paper by a constant number.) If the user is willing to commit to an annual subscription, they receive over 30% discount in price. However, users in annual contracts have to pay a sizable penalty to unsubscribe before the end of their commitment. The firm also offers educational licenses at a discounted rate to students and educational institutions, and these constitute 20.8% of the subscriptions in our data.

3.2 Field Experiment

At the time of this study, the firm used to give users a 30-day free trial for each of its software products, during which they had unlimited access to the product.¹ In order to access the product after the trial period, users need a subscription to a plan or bundle that includes that product.

To evaluate the effectiveness of different trial lengths, the firm conducted a large-scale field experiment that ran from December 1st 2015 to January 6th 2016 and spanned six major geographic markets – Australia and New Zealand, France, Germany, Japan, United Kingdom, and United States of America. During the experiment period, users who started a free trial for any of the firm’s four most popular products were randomly assigned to one of 7, 14 or 30 days free trial length buckets. These three trial lengths were chosen because they are the most commonly used ones in the industry and represent a vast majority of the SaaS free trials. Treatment assignment was at user level, i.e., once a user was assigned to a trial length, her/his trial length for the other three popular products was also set at the same length. The length of the free trial for other products during this period remained unchanged at 30 days. The summary statistics for the treatment assignment and subscriptions are shown in Table 1.

The experiment was carefully designed and implemented to rule out the possibility of self-selection into treatments, a common problem in field experiments. In our setting, if users can see which treatment (or free trial length) they are assigned to prior to starting their trial, then users who find their treatment undesirable may choose to not start the trial. In that case, the observed sample of users in each treatment condition would no longer be random, and this in turn would bias the estimated treatment effects. Moreover, since the experimenter cannot obtain data on those who choose to not to start their free trials, there is no way to address this problem econometrically. To avoid these types of self-selection problems, the firm designed

¹This free trial is at the software product level, i.e., users start a separate trial for each software product, and their trial for a given product expires 30 days from the point at which they started the free trial for it.

	7 Days Trial	14 Days Trial	30 Days Trial	Total
Number of observations (N)	51,040	51,017	235,667	337,724
Percent of total observations	15.11	15.11	69.78	100
Number of subscriptions	7,835	7,635	34,564	50,034
Percent of total subscriptions	15.66	15.26	69.08	100
Subscription rate within group (in %)	15.36	14.96	14.67	14.81

Table 1: Summary statistics of treatment assignment and subscription rates.

the experiment so that users were informed of their trial-length only after starting their trial. In order to try a software product, users had to take the following steps: (1) sign up with the firm by creating an ID, (2) download an app manager that manages the download and installation of all the firm’s products, and (3) click on an embedded *start trial* button to start the trial for a given product. Only at this point in time, they are shown the length of their free trial as the time left before their trial expires (e.g., “Your free trial expires in 7 days”). While users can simply quit or choose to not use the product at this point, their identities and actions are nevertheless captured in our data and incorporated in our analysis.

Finally, it is important to note that treatment assignment was unconfounded with other marketing mix variables. In this context, it is useful to discuss prices since they can vary across products and users. The price that a user gets for a product/bundle depends only on two user-level observables – the geographic location of the user and her/his job (students get a discount). Both of these variables are observed in the data, and treatment assignment is orthogonal to these variables. Thus, price is not confounded with treatment.²

In sum, the design satisfies the two main conditions necessary for the experiment to be deemed clean – (1) unconfoundedness and (2) compliance (Mutz et al., 2019).

3.3 Data

We have data on 337,724 users who were part of the experiment. For each user i , we observe the following information – (1) Treatment assignment (W_i), (2) Pre-treatment demographic data (X_i), and (3) Post-treatment behavioral data (Z_i). The treatment variable denotes the trial length that the user was assigned to – 7, 14, or 30 days. The variables under the latter two categories are described in detail below.

3.3.1 Pre-treatment demographic data

1. Geographic region: The geographic region/country that the user belongs to (one of the six described in $\times 3.1$). It is automatically inferred from the user’s IP address.
2. Operating system: The OS installed on the user’s computer. It can take eight possible values, e.g., Windows 7, Mac OS Yosemite. It is inferred by the firm based on the compatibility of the products downloaded with the user’s OS.

²While the distribution of prices shown to users is the same across all treatment arms, the distribution of prices paid by the subscribers can be different under each treatment arm. This is because the length of the free trial can influence which types of consumers subscribe and the products/bundles that they subscribe to. These differences can lead to downstream differences in the revenues under treatments and targeting policies. We discuss these issues in detail in $\times??$.

Variable	Number of sub-categories	Share of top sub-categories			
		1 st	2 nd	3 rd	4 th
Geographic region	6	55.02%	13.66%	9.12%	8.83%
Operating system	8	28.97%	21.4%	14.04%	13.98%
Signup channel	42	81.56%	8.14%	3.47%	0.81%
Job	14	28.20%	21.90%	20.34%	8.46%
Skill	5	69.05%	12.75%	10.77%	7.38%
Business segment	7	35.41%	32.74%	18.40%	7.81%

Table 2: Summary statistics for the pre-treatment categorical variables.

3. Sign-up channel: The channel through which users came to sign-up for the free trial. In total, there are 42 possible sign-up channels, e.g., from the legacy version of the software, from the firm’s website, through third-parties, and so on.
4. Skill: A self-reported measure of the user’s fluency in using the firm’s software suite. This can take four possible values – beginner, intermediate, experienced, and mixed.
5. Job: The user’s job-title (self-reported). The firm gives users 13 job-titles to pick from, e.g., student, business professional, hobbyist.
6. Business segment: The self-reported business segment that the user belongs to. Users can choose from six options here, e.g., educational institution, individual, enterprise.

The last three variables are self-reported though not open-ended, i.e., users are required to pick one option from a list provided by the firm. However, users may choose not to report these values, in which case, the missing values are recorded as “unknown”. We treat this as an additional category for each of these three variables in our analysis.³ The six pre-treatment variables and their summary statistics are shown in Table 2.

3.3.2 Post-treatment behavioral data

For all the users in our data, we observe their subscription and renewal decisions for approximately 24 months (from December 2015 till November 2017). We have data on:

1. Subscription information: We have data on whether a user subscribes or not, and the date and type of subscription (product or bundle of products) if she does subscribe.
2. Subscription length: Number of months that the user is a subscriber of one or more products/bundles during the 24-month observation period. If a user does not subscribe to any of the firm’s products during the observation period, then this number is zero by default.⁴

³Only a small fraction of people choose to not report these data. For example, the percentage of users with “unknown” Skill and Job is 7.4% and 21.9%, respectively.

⁴If a user unsubscribes for a period of time and then comes back, her subscription length is the total number of months that she was a paying customer of the firm. If a user subscribes to two or more plans, we aggregate the length of subscription all plans and report the total. So the subscription length can be greater than 24 months for such users.

Variable	Mean	Standard Deviation	Min	25%	50%	75%	Max	Number of Observations
Subscription	0.148	0.355	0	0	0	0	1	337,724
Subscription length (all)	2.23	6.37	0	0	0	0	108	334,223
Subscription length (subscribers)	16.02	8.43	0	10	17	22	108	46,533
Revenue (all)	79.13	285	0	0	0	0	20,208	334,223
Revenue (subscribers)	568	552	0	242	420	666	20,208	46,533

Table 3: Summary statistics of subscription, subscription length, and revenue outcomes. All the revenue numbers are scaled by a constant factor to preserve the firm’s anonymity.

3. Revenue: The total revenue (in scaled dollars) generated by the user over the 2-year observation period. This is a function of the user’s subscription date, the products and/or bundles that she subscribes to, the price that she pays for her subscription, and subscription length.

The summary statistics of these outcome variables are presented in Table 3. Both subscription length and revenue are shown for: (a) all users and (b) the subset of users who subscribed. There are a couple of points to note here. First, we do not have access to the subscription length and revenue data for team subscriptions and government subscriptions (which constitute a total of 3501 subscriptions). Hence, the number of observations used to calculate the summary statistics for subscription length and revenue for subscribers is lower. Second, the minimum subscription length observed in the data for subscribers is zero because we have a few users (58 users) who immediately unsubscribed after subscribing (within one month), in which case the firm returns their money and records their subscription length and revenue as zero. Based on Table 3, we see that approximately 14.8% users who start a free trial subscribe, and the average subscription length of subscribers is about 16 months (which is a little over a year).

We also observe the following product download and usage data for the duration of a user’s trial period.

1. Products downloaded: The date and time-stamp of each product downloaded by the user.
2. Indicator for software use: An indicator for whether the user used the software at all.
3. Number of active days: Total number of days in which the user used the software during the trial period. For example, if a user with a 7-day trial uses the software on the first and third day, this variable is two.
4. Usage during trial: Each product in the software suite has thousands of functionalities. Functionalities can be thought of as micro-tasks and are defined at the click and key-stroke level; e.g., save a file, click undo, and create a table. The firm captures all this information and we have data on the total count of the functionalities used by the user during her trial period.
5. Dormancy length: The number of days between the last active day and the last day of trial, as shown in Figure 1. For example, if a user with a 30-days trial last used the software on day 20, then her dormancy length is 10.



Figure 1: Dormancy length: Number of days between the last active day and the end of the trial period.

Variable	Mean	Std	Min	25%	50%	75%	Max	N
Total downloaded packages	1.17	0.41	1.0	1.00	1.00	1.00	4.00	337,724
Indicator for software use	0.83	0.37	0.0	1.00	1.00	1.00	1.00	303,514
Number of active days	3.03	3.94	0.0	1.00	2.00	4.00	30.00	303,514
Usage during trial	1,733	7,220	0	47	257	1,086	488,666	303,514
Log usage during trial	5.09	2.74	0.0	3.87	5.55	6.99	13.10	303,514
Dormancy length	16.87	11.23	0.0	6.00	15.00	29.00	30.00	303,514

Table 4: Summary statistics for usage features.

We present the summary statistics of these usage variables in Table 4. The usage data are also missing (at random) for a subset of users and we report the summary statistics for non-missing observations. As we can see, most users download only one software product; only 13.6% of people download more than one product. Further, 83% of users try the software at least once. However, the number of active days is relatively small; the average user uses the software for only three days during the trial period. Next, we see that an average user uses 1,733 functionalities during her trial. However, notice that this variable is very skewed with the variance much higher than the mean. So we use the natural log of this variable in all our analyses going forward. Finally, we see that the average dormancy length is close to 17 days, which means that many users stop using the software much before the end of trial period.

Finally, we refer interested readers to Tables A1 and A2 in Appendix A for the summary statistics of outcome and usage variables by trial length, respectively.

3.3.3 Training and Test Data

To design and test counterfactual free trial policies, we partition the data into two independent samples.

Training Data: This is the data that is used for both learning the model parameters as well as model selection (or hyper-parameter optimization through cross-validation).

Test Data: This is a hold-out data on which we can evaluate the performance of the policies designed based on the models built on training data.

We use 70% of the data for training (and validation) and 30% for test. See Table A3 in the Appendix A for a detailed breakdown of how the data are split across the two data-sets. Note that while the joint distributions of the variables in the two samples should be the same theoretically, there will be some minor differences between the two data-sets due to the randomness in splitting in a finite sample. It is important to keep this in mind when comparing results *across* the two data-sets.

Data	Treatment	Subscription rate	Subscription rate difference	t-statistics	Percentage gain over baseline
Training data	7 days	0.1532	0.0064	3.08	4.34
	14 days	0.1490	0.0021	1.03	1.45
	30 days	0.1468	—	—	—
Test data	7 days	0.1544	0.0082	2.58	5.59
	14 days	0.1511	0.0048	1.51	3.28
	30 days	0.1463	—	—	—

Table 5: Average effect of the 7- and 14-day treatments on subscription; compared to the control condition of 30-day free-trial. Baseline subscription rate (for 30-day case): 14.68 in training data and 14.63 in test data.

4 Main Effect of Trial Length on Subscription

We now document the main effect of trial length on subscription and present some evidence for the mechanism behind this effect. For expositional simplicity, we focus on subscription here and present a detailed analysis long-run outcomes such as revenue and retention in $\chi 7$.

4.1 Average Treatment Effect

In a fully randomized experiment (such as ours), the average effect of a treatment can be estimated by simply comparing the average of the outcome of interest across treatments. We set the 30-day condition as the control and estimate the average effects of the 14- and 7-days trials on subscriptions for training and test data. The results from this analysis are shown in Table 5.

The 7-day trial increases the subscription rate by 4.34% over the baseline of the 30-day condition in the training data and by 5.59% in the test data. However, in both data sets, the effect of the 14-day trial is not significantly different from that of the 30-day trial. These results suggest that a uniform targeting policy that gives the 7-day treatment to all users can significantly increase subscriptions.⁵ We also see that the average treatment effect is fairly small compared to the outcome, which is either zero or one. This is understandable since the effect of trial length is likely to be small compared to other factors that affect customer acquisition. Finally, note that the gains and subscription rates in the training and test data are slightly different. As discussed earlier, this is due to the randomness in the splitting procedure.

Next, to ensure that these results are not driven by any problems with randomization, we conduct a series of randomization checks. We present the details of these tests in Appendix $\chi B.2$ and discuss them briefly here. First, we conduct a joint test of orthogonality of pre-treatment variables and treatment assignment (McKenzie, 2017). This is done by regressing the treatment variable on the entire set of pre-treatment variables (with dummies for each sub-category shown in Table 2). We find that the pre-treatment variables have no predictive power when it comes to predicting treatment, which suggests that randomization was done correctly. Note that this approach to checking for potential issues with randomization is preferable to the old practice of showing tables of means for pre-treatment variables across treatment arms and running a battery of t-tests for

⁵In general, it is a better practice to obtain ATEs directly based on mean comparisons without using regression-based approaches (Imbens and Rubin, 2015). Nevertheless, in Appendix B.1, we present the ATEs based on regressions (with and without controls) and they are statistically indistinguishable from those shown in the main text.

Outcome variable	Intercept	14 days trial	30 days trial	R^2	N
Total downloaded packages	1.137 (0.002)	0.01 (0.002)	0.017 (0.002)	0.000	337724
Indicator for software use	0.828 (0.002)	0.004 (0.002)	0.009 (0.002)	0.000	303514
Number of active days	1.747 (0.018)	0.625 (0.026)	1.711 (0.02)	0.028	303514
Number of active days/trial length	0.25 (0.001)	-0.08 (0.001)	-0.134 (0.001)	0.078	303514
Log usage during trial	4.77 (0.013)	0.196 (0.018)	0.411 (0.014)	0.003	303514
Log average daily usage during trial	3.197 (0.009)	-0.357 (0.012)	-0.737 (0.009)	0.022	303514
Dormancy length	4.631 (0.043)	5.135 (0.06)	16.432 (0.047)	0.337	303514

Table 6: Regression of usage features on trial length. Standard errors in parentheses.

a variety of reasons; see Bruhn and McKenzie (2009) and Mutz et al. (2019) for detailed discussions.⁶ Next, we regress the outcome variable (subscription outcome) on the treatment variable and all the pre-treatment variables. We find that the treatment effects are very similar to those in Table 5, which again suggests that there are no issues with randomization.⁷

4.2 Mechanism

At the time of the experiment, the firm offered a standard 30-day free trial to all its consumers. The better performance of the much shorter 7-day trial was both surprising and inexplicable for many reasons. First, the firm sells a complicated suite of software with multiple products and functionalities. So we would have expected that giving consumers more time to familiarize themselves with it and learn the software would produce better outcomes. Second, the reasons proposed in the theory literature for the efficacy of free trials largely support longer free trials, e.g., switching costs, consumer learning, software complexity, and signaling. Thus, it is not obvious why a shorter trial works better. Therefore, we now examine how trial length affects conversion and present some evidence for why a shorter trial works better in this setting. In the process, we also discuss the generalizability of our findings and the mechanisms proposed.

Intuitively, trial length can affect how consumers download, use, and interact with the software; and differences in these usage variables can lead to different subscription outcomes. So we first examine whether and how trial length affects usage. We regress each of the usage variables shown in Table 4 on trial length and present the results in Table 6. Since trial length is randomly assigned, we can interpret these results causally. First, we find that longer free trials lead to more product downloads and more usage. Further investigation suggests that this increase in downloads mainly comes from the higher downloads of products 1 and 3, which are complements (see Figure A1 in Appendix C). This suggests that giving longer trial lengths to users increases their probability of exploring other complementary products. Next, we see that a larger fraction of people try the software at least once with a longer trial, and the number of active days and log usage also increases with trial length. However, the rate of increase in the number of active days and usage is sub-linear compared to the increase in trial length. For instance, going from 7 to 14 days increases the

⁶Further, presenting tables of means for each pre-treatment variable is not feasible in our case since all our pre-treatment variables are categorical with a large number of sub-categories.

⁷Later in the paper, we use the empirical propensities to evaluate the gains from our models. So any minor discrepancies in the propensities of treatment allocation are taken care of; see Equation (4) and the discussion around it.

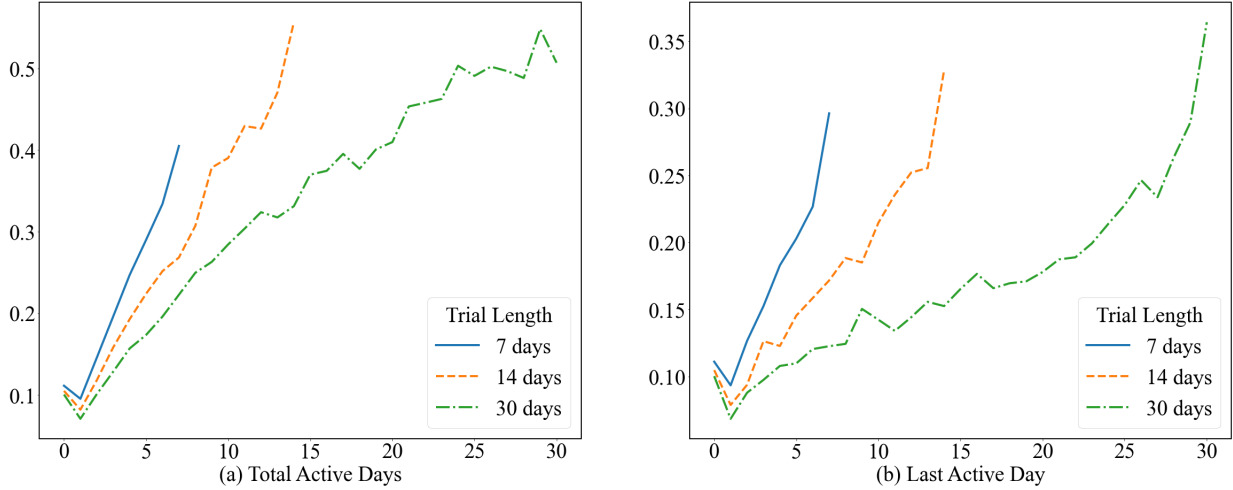


Figure 2: (a) The subscription rate based on the last day of trial use for different trial lengths. (b) The subscription rate based on the number of active days for different trial lengths.

number of active days by 0.625, which is much smaller than 7 days (the increase in trial length). Thus, when we normalize the number of active days by trial length, the average number of days during which a user is active during her trial reduces as trial length increases. The same pattern holds for log usage; while total usage increases as trial length increases, average daily usage falls. Finally, we find that the dormancy period increases as trial length increases. While the average dormancy length is 4.6 days for the 7-days trial, it is over 21 days for the 30-days trial.

Next, we examine whether and how usage is associated with subscription. The left panel of Figure 2 shows the probability of subscription as a function of the total number of active days for each trial length. As we can see, users who are active for more days are also more likely to subscribe, and this pattern is true for all three trial lengths. However, given the same level of activity, shorter trial lengths are associated with higher conversion. For example, users who were active for five days are more likely to subscribe when they are in the 7-day condition, compared to the 14 or 30-days condition. Next, in the right panel of Figure 2, we show the probability of subscription as a function of the last active day for all three trial lengths. We see that users whose last active day is earlier in the trial period are less likely to subscribe. Further, for the same last active day, users with shorter trials are more likely to subscribe. Recall that dormancy length is defined as trial length minus the last active day. So this suggests that users who have not used the product for long periods at the end of the trial period are less likely to subscribe.

We now check if the preliminary patterns shown in Figure 2 hold after we control for other usage and user-specific observables. In Table 7, we present the results from a regression with the user's subscription decision as the outcome variable and her trial length and usage variables as explanatory variables. We find that after controlling for everything else, users who log more active days and use the product more are more likely to subscribe. Further, users who have longer dormancy periods are less likely to subscribe. This is understandable because a user who has not used the software for a long time by the end of her trial period is

	coef	std err	z	$P > z $	[0.025	0.975]
Indicator for using the software	-0.5145	0.036	-14.252	0.000	-0.585	-0.444
Total downloaded packages	0.5632	0.013	43.789	0.000	0.538	0.588
Number of active days	0.0440	0.002	19.241	0.000	0.040	0.049
Log usage during trial	0.0620	0.005	11.267	0.000	0.051	0.073
Dormancy length	-0.0297	0.001	-30.141	0.000	-0.032	-0.028

Table 7: Regression of subscription on usage features and trial length, with all the pre-treatment variables included as controls (not shown in the table above).

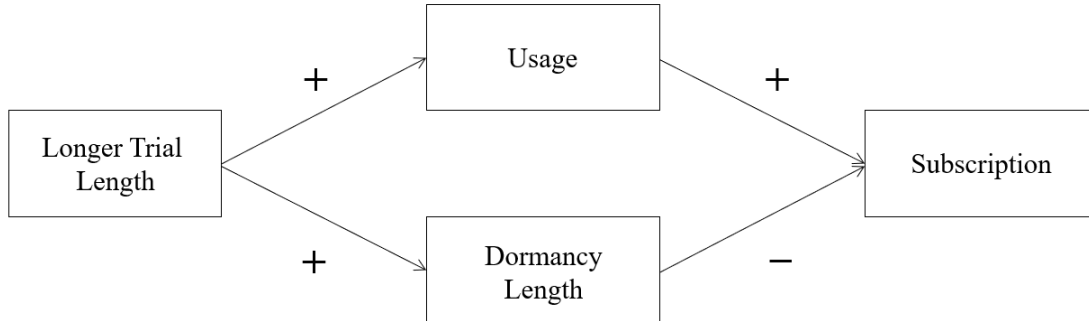


Figure 3: The effect of trial length on usage and dormancy length, and subsequently subscription.

likely to forget about it and/or conclude that the product is not useful (Zhu et al., 2018).

Together, the above findings suggest that two opposing effects of trial length on usage and subscription. We depict these effects in Figure 3. On the one hand, as trial length increases, product usage and consumer learning about the software increases. This increase in usage can have a positive effect on subscriptions. On the other hand, as trial length increases, the gap between the last active day and the end of the trial increases, while the average number of active days and usage per day reduces. These factors are associated with lower subscriptions. In our case, it seems that the latter effect dominates the former, and hence shorter trials are better.⁸

Our analysis presents three key findings relevant to the broader theories on the role of free trials for experience goods. First, we rule out the well-known demand cannibalization hypothesis advocated by many theoretical papers (Cheng and Liu, 2012; Dey et al., 2013). These papers argue that, with longer trials, free-riders can use the product extensively during the trial, get their project/job done, and avoid subscribing. However, the results in Figure 2 and Table 7 rule out the free-riding hypothesis because users who use the product heavily during the trial are also more likely to subscribe. However, this evidence is for the full population of users. Second, we provide empirical support for the consumer learning hypothesis proposed in analytical papers (e.g., (Dey et al., 2013)) since we find that longer trials lead to more usage, which in turn is associated with higher subscription. Third, we identify a novel mechanism that plays a significant role in the effectiveness of free trials – the negative effect of long dormancy periods on subscription. We provide more

⁸The results in Table 7 should only be interpreted as suggestive evidence for the second half of the mechanism shown in Figure 3 (and not causally). This is because the unobserved attributes of the user that drive usage may also drive subscription.

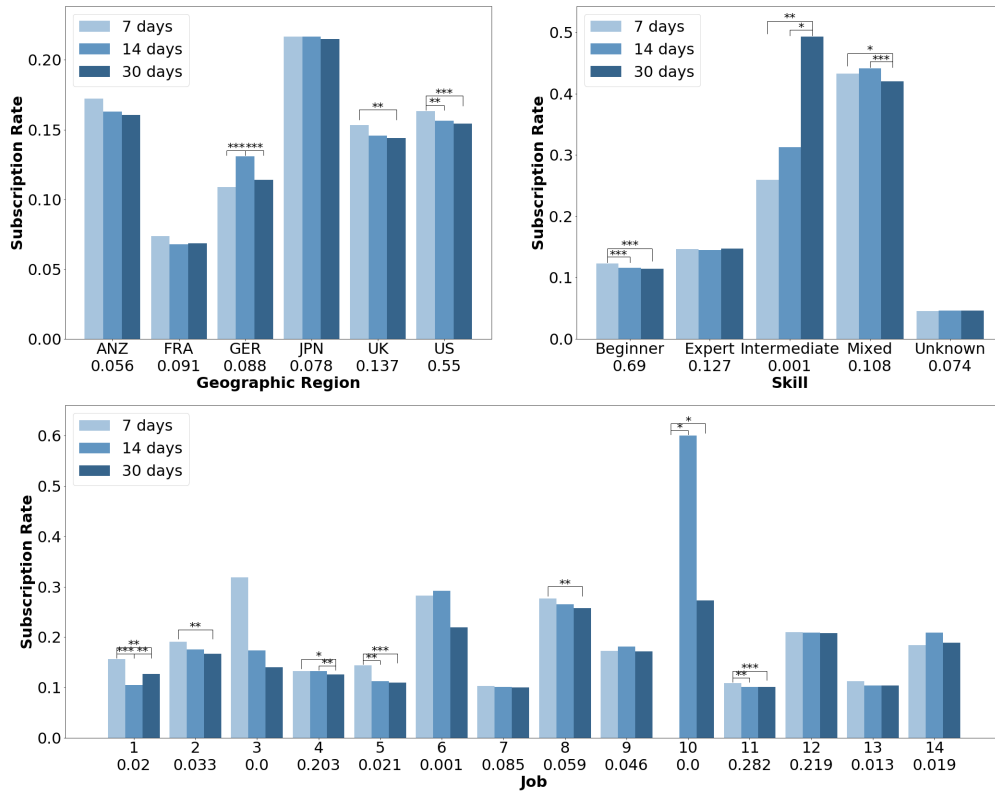


Figure 4: Heterogeneity in consumers’ response to the three trial lengths within three categories – Geographic region, Skill, and Job. The six geographic regions shown are: Australia and New Zealand, France, Germany, Japan, and United States of America (in that order). Under each sub-category, the fraction of users in that sub-category are shown. We do not include sub-category names for Job to preserve the firm’s anonymity. (* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$)

evidence in support of these ideas in $\times 6$, where we discuss the heterogeneous response of different types of users.

4.3 Heterogeneity in Users’ Responsiveness

So far, we have shown that the 7-days trial is the best average treatment and provided some intuition for why. However, the effect of trial length could be heterogeneous across users and the mechanisms discussed earlier could be differentially important for different types of users. We now examine if this is indeed the case.

In the top left panel of Figure 4, we partition the data into six sub-groups based on the user’s geographic region and present the average subscription rates for the three trial lengths for each region. The results suggest that there is some heterogeneity in response rates by region. For example, the 14-day trial is more effective in Germany while the 7-day trial is more effective in the United States of America. Next, we perform a

similar exercise on skill-level and job (see the top right and bottom panels in Figure 4). Again, we find that users’ responsiveness to the treatment is a function of their skill level and job. For instance, the 7-day trial is significantly better for Beginners, whereas the 14-day trial is more effective for Mixed-skill users.

These results suggest that users’ responsiveness to trial lengths is heterogeneous on many pre-treatment variables. If the firm can successfully exploit the different sources of heterogeneity and personalize its free trial assignment at the individual-level, then it may be able to further improve subscriptions.

5 Counterfactual Analysis: Personalized Policy Design and Evaluation

Given that the preliminary evidence above suggests that benefits the firm can benefit from personalizing free trial assignment. In §5.1, we describe the procedure we use to design the personalized policy. Next, in §5.2, we present the gains from the personalized policy in our setting. Next, in §5.3, we compare the performance of our approach to other personalized policies. Finally, in §5.4, we examine why some consumers respond better to shorter trials (vs. longer trials) and tie the policy-prescribed segmentation to mechanisms discussed in §4.2.

5.1 Optimal Policy Design

Let $i \in \{1, \dots, N\}$ denote the set of independent and identically distributed users, where each user is characterized by a pre-treatment covariate vector $X_i \in X$ of dimension D . Let $W_i \in W$ denote the treatment or intervention that i receives. $W = \{0, \dots, W-1\}$ refers to the set of treatments, and the total number of treatments is W . Finally, let $Y(X_i, W_i)$ denote the outcome for a user i with pre-treatment variables X_i when she is allocated treatment W_i .

A personalized treatment assignment policy, π , is defined as a mapping between users and treatments such that each user is allocated one treatment, $\pi : X \rightarrow W$. The firm’s goal is to choose a policy π such that it maximizes the expectation of outcomes, $\frac{1}{N} \mathbb{E} \left[\sum_{i=1}^N Y(X_i, W_i^\pi) \right]$. Thus, for policy π and outcome of interest Y , we can write our reward function as $R(\pi, Y) = \frac{1}{N} \sum_{i=1}^N \mathbb{E} [Y(X_i, \pi(X_i))]$. Thus, given a reward function $R(\pi, Y)$, the optimal personalized policy is given by:

$$\pi^* = \arg \max_{\pi \in \Pi} [R(\pi, Y)], \quad (1)$$

where Π is the set of all possible policies.

The problem of finding the optimal personalized policy is equivalent to one of finding the policy π^* that maximizes the reward function $R(\pi, Y)$. As discussed in §1, this is a non-trivial problem since the cardinality of the policy space can be quite large.⁹ So, a direct search over the policy space to find the optimal policy is infeasible. Therefore, we adopt a two-step approach to find the optimal policy π^* that avoids this problem. To do so, we make the three standard assumptions on: (1) unconfoundedness, (2) SUTVA, and (3) positivity. Given that our data comes from a fully randomized experiment assumptions (1) and (2) are automatically satisfied.

⁹The total number of possible policies is $W^{\prod_{d=1}^D c_d}$, when we have D pre-treatment variables and the d -th variable can take c_d different values. This number can be extremely high even in simple settings. In our application, the cardinality of the policy space is equal to $3^{987,840}$.

Further, assumption 2 is satisfied because we do not expect any network effects in our setting (since the experiment was run on unconnected users distributed all over the world).

With these assumptions in place, we can design the optimal personalized policy if we either have estimates of the outcome of interest or pairwise treatment effects. In the main analysis, we design our personalized policy based using outcome estimates based on lasso (Tibshirani, 1996; Friedman et al., 2010). That is, we model the subscription outcome using as $f(x, w) = \mathbb{E}[Y|X_i = x, W_i = w]$, where $f(\cdot)$ is a lasso model. Lasso estimates a linear regression that minimizes the MSE with an additional term to penalize model complexity as shown below:

$$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3) = \arg \min \sum_{i=1}^n (Y_i - X_i\beta_1 - W_i\beta_2 - X_iW_i\beta_3)^2 + \lambda(j|\beta_1|_1 + j|\beta_2|_1 + j|\beta_3|_1), \quad (2)$$

where $j|\beta_i|_1$ is the L1 norm of the vector β_i and is equal to the sum of the absolute value of the elements of vector β_i . Intuitively, if there are multiple weak (and correlated) predictors, lasso will pick a subset of them and force the coefficients of others to zero, thereby estimating a simpler model. Model-selection in lasso is data-driven, i.e., λ is a hyper-parameter that is learned from the data (and not assumed). Please see Appendix D for details of our lasso estimation.

Next, using estimate of the expected outcome, $\hat{y}(x = X_i, w)$, from the lasso model, we obtain the optimal personalized policy based on for observation i as:

$$\pi_{lasso}(X_i) = w, \quad \text{where } w = \arg \max_{w \in W} \hat{y}(x = X_i, w) \quad (3)$$

Our personalized free trial policy, π_{lasso} partitions the population into three segments: 7-, 14- and 30-days optimal segments, which constitute 68.9%, 23.2%, and 7.9% of the population, respectively.

5.2 Empirical Policy Evaluation: Gains from Personalization

We now empirically evaluate and quantify the gains from the personalized free trial policy π_{lasso} over non-personalized policies. To do so, we first define three uniform (one length for all) policies:

π_{30} – This policy prescribes the 30-days treatment for all users. It was used by the firm at the time of the experiment and we therefore use it as the baseline policy in all our comparisons.

π_{14} – This policy prescribes the 14-day treatment for all users.

π_7 – This policy prescribes the 7-day treatment for all users. Since we found that 7 days is the best average treatment in $\mathcal{X}_4.1$, this is the best uniform policy.

We evaluate the expected reward from the policies (both personalized and uniform) using the Inverse Propensity Score (IPS) estimator, that has been extensively used in the off-policy evaluation literature (Horvitz and Thompson, 1952; Dudík et al., 2011), and has recently been applied in the marketing too (Simester et al., 2019; Rafieian and Yoganarasimhan, 2021; Hitsch and Misra, 2018). For any given policy π , this estimator takes all the observations where the user received the policy-prescribed treatment and scales them up by

Policy category	Policy	Estimated Subscription (%)		Increase in subscription (%)	
		Training Set	Test Set	Training Set	Test Set
Personalized based on lasso	<i>lasso</i>	15.85	15.62	7.97	6.81
Uniform	7	15.32	15.44	4.34	5.59
	14	14.90	15.11	1.45	3.28
	30 (Baseline)	14.68	14.63	—	—
Alternative personalized policies	<i>reg</i>	15.89	15.33	8.21	4.83
	<i>cart</i>	15.32	15.44	4.34	5.59
	<i>r-forest</i>	17.42	14.82	18.67	1.32
	<i>xgboost</i>	16.00	15.53	8.98	6.17
	<i>c-forest</i>	15.58	15.46	6.09	5.71
	<i>c.tree</i>	15.32	15.44	4.34	5.59

Table 8: Gains in subscription from implementing different counterfactual free-trial policies. The results for policies π_{cart} , $\pi_{c.tree}$, and π_7 are the same since they prescribe the 7-days treatment to all users.

their propensity of receiving the treatment assigned to them. This scaling gives us a pseudo-population that received the policy-prescribed treatment. Thus, the average of the outcome for this pseudo-population gives us an unbiased estimate of the reward for the full population, if we were to implement the proposed policy in the field. Formally:

$$\hat{R}_{IPS}(\pi, Y) = \frac{1}{N} \sum_{i=1}^N \frac{1[W_i = \pi(X_i)]Y_i}{\hat{e}_{\pi(X_i)}(W_i)}, \quad (4)$$

where $\hat{e}_{\pi(X_i)}(W_i)$ is the probability that a user whom the policy prescribes treatment $\pi(X_i)$ is given W_i .¹⁰

We present the expected rewards (or subscription rates) from all the three uniform policies and π_{lasso} in the top panel of Table 8. The key finding is that personalization based on pre-treatment demographic variables leads to over 6.8% improvement in subscription compared to the baseline uniform policy of giving a 30-days trial for all. Further, we see that the personalized policy also does better than the best uniform policy of 7 days for all. To examine if this difference is significant, we conduct a paired t-test based on bootstrapping as follows. We repeatedly (20 rounds) split the entire data into training and test (in the same proportion used in the main analysis, i.e., 0.7/0.3). Then, in each round, we train a lasso model on the training set using a five-fold cross-validation and calculate the IPS-rewards (based on Equation 4) for both π_7 and π_{lasso} on the test data. Finally, we run a two-sided paired t-test to compare lasso’s performance with the uniform all 7-days policy. The t-statistic and p-value for the two-sided test are 3.123 and 0.0056, respectively, which confirms that the personalized policy π_{lasso} is better than the best uniform policy π_7 .

That said, notice the magnitude of gains from personalization (over the best uniform policy) is modest. This is understandable since the personalized policy assigns about 70% of users to the 7-day treatment, and the gains from personalization only accrue from the remaining 30% of users who are allocated the 14- or

¹⁰In theory, in a randomized experiment, the propensity of treatment assignment is orthogonal to the treatment prescribed by any policy π . Thus, $e(W_i = w; X_i) = e(W_i = w) \forall w \in \mathcal{W}$ is known and constant for all observations. However, in practice, within the set of users for whom policy π prescribes w , the empirical treatment propensities might not be the same as that in the full data. So we use the empirical propensity, defined as: $\hat{e}_{\pi(X_i)}(W_i) = \frac{\frac{1}{N} \sum_{j=1}^N 1[W_j = W_i, \pi(X_j) = \pi(X_i)]}{\frac{1}{N} \sum_{j=1}^N 1[\pi(X_j) = \pi(X_i)]}$.

30-days treatments. As Simester et al. (2020) point out, this is because the difference in the the performance of the two policies for users assigned to the 7-day trial is exactly zero. Further, our treatment effect is small compared to the outcome – a common occurrence for marketing interventions such as advertising or promotions (Lewis and Rao, 2015). These findings are consistent with the recent literature on personalization (Yoganarasimhan, 2020; Rafieian and Yoganarasimhan, 2021), which demonstrate positive but moderate gains from personalization digital interventions.

5.3 Comparisons and Robustness Checks

So far, we used lasso as the outcome estimator to design our personalized policy. We now examine whether counterfactual personalized policies based on other outcome and heterogeneous treatment effects estimators perform better. Specifically, we consider policies based on four additional outcome estimators: (1) linear regression, (2) CART, (3) random forests, and (4) XGBoost, and two heterogeneous treatment effect estimators: (1) causal tree, and (2) generalized random forests. The technical details of these models and their tuning details are shown in Appendices E and F.

First, we find that each of these policies behaves quite differently when it comes to treatment assignment (see Table A8 in Appendix G for details). Interestingly, we find that policies based on CART and causal tree do not personalize treatment assignment and end up giving the 7-day treatment to all users, that is: π_{cart} and $\pi_{c.tree}$. We also find that $\pi_{xgboost}$ is quite similar to π_{lasso} . Both prescribe the 7-days trial to 70% of users, the 14-days trial to 20% of users, and the 30-days trial to 10% of users. In contrast, π_{reg} and $\pi_{r-forest}$ prescribe the 7-day treatment to the least number of users while $\pi_{c-forest}$ prescribes the 7-day treatment to 91% of users (and the 30-day treatment to no one).

Next, in the bottom panel of Table 8, we present the performance of these policies. We find that π_{lasso} continues to be the best, and the second-best policy is $\pi_{xgboost}$. There are two main takeaways here. First, poorly designed personalized policies (e.g., those based on regression and random forest) can actually do worse than the best uniform policy on the test data. Second, we do not find much correlation between an outcome estimator’s predictive ability and its efficacy in policy design. For instance, random forest has a lower mean squared error on the test data compared to lasso, but $\pi_{r-forest}$ is much worse than π_{lasso} (see Table A9 in Appendix G). This is likely because the objective function in outcome estimation methods is predictive ability, which is different from policy design or performance. In sum, our findings suggest that managers should be careful in both designing and evaluating personalized policies. It is critical to: (1) not conflate a model’s predictive ability with its ability to form policy, and (2) evaluate the performance of each policy on an independent test data with appropriate policy evaluation metrics.

Next, we find that the recently proposed heterogeneous treatment effects estimators, causal tree and causal forest, perform poorly when it comes to personalized policy design. Our results suggest that managers may be better off adopting the best uniform policy instead of investing resources in personalizing policies based on these methods. This is an important finding since these methods are gaining traction in the marketing literature and researchers are starting to use them (e.g., Guo et al. (2017), Fong et al. (2019)). Our findings suggest that relying on heterogeneous treatment effects estimators can be sub-optimal.

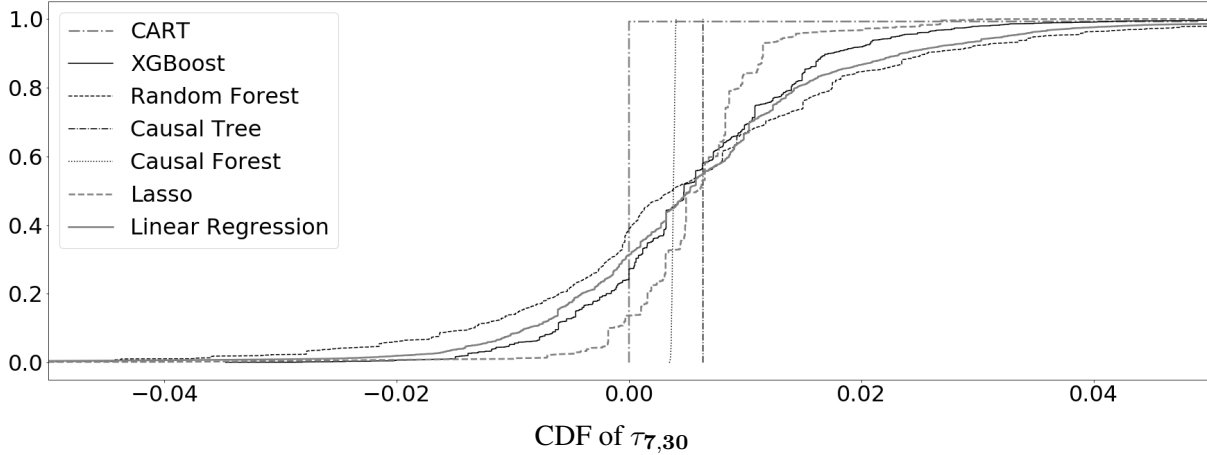


Figure 5: The CDF of estimated CATEs for 7 vs 30 days of free trial from using different methods (for test data).

Finally, we examine if there is any relationship between the estimates of treatment effects based on a specific method and the performance of the policy based on it. Figure 5 shows the CDF of $\tau_{7,30}$ for all the methods used for policy design.^{11 12} The first pattern that stands out is that treatment effects estimates based on CART, causal tree, and casual forest show very little heterogeneity (see the three vertical lines to the right of zero). This explains why policies based on these methods perform poorly – they are unable to personalize the policy sufficiently to optimize users’ response at the individual level. In contrast, the treatment effect estimates based on linear regression and random forest show the maximum amount of heterogeneity (see the two rightmost curves). This pattern, in combination with the poor performance of π_{reg} and π_{r_forest} on test data (and their extremely good performance on training data) hints at overfitting problems. That is, these models seem to infer much more heterogeneity than is true in the data. Interestingly, we see that the CDFs of treatment effects based on lasso and XGBoost lie somewhere in between the above two groups. They show sufficient heterogeneity, but not too much. Hence, policies based on these methods are able to personalize the treatment sufficiently without overfitting. Recall that the dispersion in treatment assignment for these two policies is higher than that in π_{cart} , π_{c_tree} , and π_{c_forest} , but lower than that in π_{reg} and π_{r_forest} . Thus, the ideal estimators for policy design are those that are able to capture sufficient heterogeneity to personalize effectively without overfitting (i.e., capture spurious heterogeneity).

6 Segmentation Analysis and Additional Evidence for Mechanism

So far, we focused on the question of “Who (should get a treatment)”. We now examine the question of “Why (should s/he get a specific treatment)”. Understanding why some users respond well to longer trials while others respond better to shorter trials can give us insight into consumers’ preferences and decision-making

¹¹The estimated distributions of $\tau_{14,30}$, and $\tau_{7,14}$ are shown in Figure A2 in Appendix G and their interpretations are largely similar to that presented here for $\tau_{7,30}$.

¹²For the outcome estimators, we can estimate treatment effects from outcome estimates as $E[Y|X_i = x; W_i = 7 \text{ days}] - E[Y|X_i = x; W_i = 30 \text{ days}]$. For heterogeneous treatment methods, these estimates are directly available (see Equation (A.4)).

process. These insights are valuable for two reasons. First, from the firm's perspective, they can be leveraged to improve other marketing interventions such as advertising and pricing. Second, from a research perspective, this gives us a better understanding of the sources of heterogeneity in the effectiveness of trial length on conversion and mechanisms at play, which can be generalized to other settings.

We now correlate a user's optimal treatment with her pre-treatment demographic and post-treatment behavioral variables. In the process, we present additional evidence for the mechanism through which trial length affects conversion, as discussed in 4.2. We conduct three sets of analyses to understand the mechanism and characterize the three segments. First, we quantify the differential effect of trial length on the download and usage behavior of the three segments. Second, we characterize the heterogeneity in the effect of usage on subscription across the three segments. Finally, we correlate a user's optimal treatment with her/his pre-treatment demographics and post-treatment outcomes to characterize the three segments. We refer readers to Appendix H for the details of these analyses and provide a summary of the three segments below.

7-day optimal segment: A vast majority of these users are beginners or students, and they are the least likely to subscribe. These users use the product more when given longer trials but don't scale up their usage as much as the 14-day optimal segment. This is understandable because most of them lack the skills to use the product extensively, even if given the opportunity to do so. Further, the negative effect of long dormancy periods is the most severe for this segment. This is understandable since these users are less familiar with the product, so when they go for a long period without using the software, they are more likely to conclude that the software is not useful and choose not to subscribe.

Overall, we find that short trials are more effective for beginners and new users because even though there are some positive effects of learning and usage, extended periods of inactivity at the end of long trials can have a strong negative effect on their subscription. One might wonder if this result simply stems from the fact that beginners have short tasks that require more than a week to complete (but still less time than 14/30 days), which leads them to have lower subscription rates when assigned the 14 and 30-day trial (i.e., a more complex version of the demand cannibalization hypothesis). However, if this explanation is true, then we should find that beginners/7-day optimal users who are assigned to the 14- and 30-day trial should be less likely to subscribe if they use the product more. However, we find the opposite; see Appendix H.4.

14-day optimal segment: These users are more likely to be mixed-skill, and they have the highest usage and subscription rates. This segment takes the most advantage of longer trials, i.e., they use the product the most and have the shortest periods of dormancy when given longer trials. It seems like these users actually try the product's features and evaluate the product carefully before deciding whether to subscribe or not. However, the effect of usage on subscriptions is lower for these users (compared to the other two segments). This is likely because they are figuring out whether the software is a good fit or not, and more usage may lead some users to learn that it is not a good fit. Further, the magnitude of the negative effect of dormancy length on subscription is also high for them. That is why the 30-days trial is not optimal for these users: the higher usage that comes with a more extended trial does not translate to big differences in

subscription, but they still get hit by the increase in dormancy length with the 30-days trial. On the other hand, when they are given only 7-days, they cannot use the product much, and the benefit from higher usage is not realized. Thus, 14 days is ideal for these users.

30-day optimal segment: These users are more experienced than average are less likely to be students and hobbyists, and more likely to sign up through the app manager instead of the website. These factors suggest that these are more likely to be experienced/legacy users who are already familiar with the software. Long dormancy periods have the least negative effect on these users. This is understandable because these users are likely to be already aware of and experienced with the software. Thus, they are unlikely to infer that the product is not useful if they don't use it for a few days at the end of the trial. Further, longer trials lead to more usage for these users, and the effect of usage on subscription is also high. Thus, giving them 30-days for trial is good.

One interesting pattern in the above findings is the non-monotonicity of usage across the three segments. We find that 7-day optimal users use the product the least, followed by the 30-day optimal users, while the 14-day optimal users use the product the most. This can be explained by the relative expertise-levels of the three groups. The extent to which a user uses the product depends on two factors: (1) how much do they need to evaluate the product?, and (2) how much can they evaluate product? The 7-day optimal users, who are pre-dominantly beginners have the least ability to explore the product features, and hence use it the least. In contrast, the 30-day optimal users, who are more likely to be experts and legacy users, have the highest ability to evaluate the product. However, given their expertise and familiarity with the software, they can do this without extensive usage. Finally, the 14-day optimal users, who are more likely to be mixed-skill users, have both high need to evaluate the product and sufficient ability to explore it. Hence, they have the highest usage.

It worth mentioning that our findings provide partial support to the theories proposed in the literature on the relationship between users' skill-level/experience and the effectiveness of free trials. For example, Dey et al. (2013) argue that longer trials are beneficial only when the learning rate is sufficiently large. We find that this is true in our case as well. However, this prior analytical research does not consider the negative effect of dormancy length on subscription, especially for beginners and new users. They argue that beginners should be given longer free trials because longer trials allow them to learn about the product, which increases their likelihood of subscription. In contrast, we find that short trials are optimal for beginners. While longer trials have a positive impact on the usage and subscription of this group, they are also the group whose subscription is most negatively affected by longer dormancy periods. Thus, ignoring the negative impact of dormancy length can lead us to make sub-optimal allocations of trial lengths for different segments.

Our findings suggest that firms and managers should take into account the heterogeneity in the evolution of usage and inactivity (as trial length increases) for different consumer types and customize trial lengths based on these patterns. In our setting, users require some skill and need to invest the effort to learn and use the software effectively. In particular, beginners and inexperienced users are unable to scale up their usage with longer trials, and therefore have longer periods of inactivity later in the trial period (which has a detrimental effect on subscription). However, if the software is simple and easy to use, we would not see

such periods of inactivity. Interestingly, this suggests that simpler products may benefit from longer trials (especially for beginners), whereas more complex products may benefit from shorter trials. In sum, both the complexity of the product and the skill of the user jointly determine usage and activity (or inactivity), which then affects subscription. Our results provide some general guidelines to firms on how to pick the right trial length for different products and segments.

7 Long-term Outcomes: Consumer Loyalty and Profitability

So far we have focused on short-run outcomes in our policy design and evaluation. However, a policy that maximizes subscriptions (or short-run conversions) may not be the best long-run policy if it brings in users who are less profitable or less loyal. For example, a policy that increases subscriptions among students (who get a significant educational discount and hence pay lower prices) and/or users who subscribe to lower-end products/bundles (that are priced much lower than the all-inclusive software suite) at the expense of high-end users can lead to lower revenues. Similarly, a policy designed to maximize subscriptions may do so at the expense of long-term retention, i.e., it may bring in the less loyal consumers who churn within a short period. Thus, a subscription-optimal policy may in fact be sub-optimal from the perspective of long-run outcomes (Gupta et al., 2006; Fader and Hardie, 2009; McCarthy et al., 2017). In this section, we therefore examine two important post-subscription outcomes of interest for the firm.

Consumer loyalty, as measured by subscription length or the number of months a user subscribes to the service over the two year period after the experiment.

Consumer profitability, as measured by the revenue generated by the user over the two years after the experiment. (In SaaS settings, revenues and profits can be treated as equivalent since the marginal cost of serving an additional user is close to zero.)

7.1 Gains in Retention and Revenue from Counterfactual Policies

We first show the average treatment effect of the three trial lengths on retention and revenue in Table 9.¹³ We find that the 7-days trial continues to be the best. In the test data, it increases retention by 6.4% and revenue by 7.91%. The average effect of the 14-days trial is both smaller in magnitude and not significant in the training data.¹⁴ These results largely mirror our findings on the average treatment effect of subscription, i.e., the 7-days trial is the best treatment.

Next, we examine how the uniform and personalized targeting policies described in §5.2 perform on

¹³A minor point is that we do not have access to subscription length and revenue data for all subscribers (recall the discussion in §3.3.2). So we treat the missing observations as zero in calculations. The results remain qualitatively unchanged if we instead work with the subset of users for whom we have non-missing revenue data.

¹⁴Note that 14-days trial outperforms the 7-days trial in the test data (on revenue) even though the 7-days trial is the best policy in the training data. We present a brief explanation for this discrepancy now. In general, estimates from one data set are valid in another data set only when the joint distribution of outcomes and covariates are similar in both data sets. However, in finite samples, there are usually some minor differences in training and test data due to the randomness in the splitting procedure. In our case, the main difference is this – the distributions of subscription length for the 14-day condition in the training and test data are different. This is however not the case for the 7- or 30-day conditions; see Table A16 in Appendix I. Thus, the estimate of subscription length from the training data does not translate well to test data, and this leads to the large difference in the subscription length and revenue estimates across the training and test data sets.

Data	Treatment	Average Subscription Length	Retention difference	t-statistics	Percentage gain over baseline
Training data	7 days	2.32	0.16	4.27	7.27
	14 days	2.22	0.06	1.54	2.59
	30 days	2.17	—	—	—
Test data	7 days	2.33	0.14	2.42	6.28
	14 days	2.32	0.13	2.27	5.91
	30 days	2.19	—	—	—

Data	Treatment	Average Revenue	Revenue difference	t-statistics	Percentage gain over baseline
Training data	7 days	82.65	6.17	3.75	8.06
	14 days	79.08	2.59	1.58	3.38
	30 days	76.49	—	—	—
Test data	7 days	83.72	6.14	2.42	7.92
	14 days	84.02	6.44	2.53	8.30
	30 days	77.58	—	—	—

Table 9: Average effect of the 7- and 14-day treatments on long-term variables (subscription length and revenue) compared to the control condition of 30-day free-trial.

Category	Policy	Subscription Length				Revenue			
		Estimate (Months)		Increase (%)		Estimate (\$)		Increase (%)	
		Training	Test	Training	Test	Training	Test	Training	Test
Personalized	<i>lasso</i>	2.39	2.36	10.42	7.96	85.96	86.67	12.38	11.72
Uniform	7	2.32	2.33	7.27	6.28	82.65	83.72	8.06	7.92
	14	2.22	2.32	2.59	5.91	79.08	84.02	3.38	8.30
	30 (Baseline)	2.17	2.19	—	—	76.49	77.58	—	—

Table 10: IPS estimates of the average subscription length and revenue under counterfactual policies (three uniform and one personalized).

the two long-term outcomes of interest. To derive the the IPS estimates of average subscription length and revenue under policy π , we first segment users into three groups based on the policy-prescribed treatment: (1) $\pi(X_i) = 7$ days, (2) $\pi(X_i) = 14$ days, and (3) $\pi(X_i) = 30$ days. Then, we use the observed subscription lengths and revenues as the outcome variables (Y_i) in Equation (4) to estimate the IPS rewards for these outcomes. Table 10 shows the results from this analysis. The main takeaway is that the personalized policy, π_{lasso} , which was designed to maximize subscriptions also does well on consumer loyalty and revenue compared to the other uniform policies. This is valuable from the the firms' perspective because it suggests that policies optimized for short-run outcomes are largely aligned with long-run outcomes as well.

7.2 Gains in Short-run vs. Long-run Outcomes

An interesting empirical pattern here is that the gains in subscription length and revenues are quantitatively different from the gain in subscription (compare the percentage increases in Tables ?? and 10). We now discuss the source of this difference.

We can write down the expected subscription length (denoted by Y_i^l) conditional on treatment W_i as:

$$E(Y_i^l | W_i) = \Pr(Y_i^s | W_i) E[T_{end} - T_{start} | W_i, Y_i^s = 1], \quad (5)$$

where $\Pr(Y_i^s | W_i)$ is the probability that user i will subscribe conditional on receiving treatment W_i and $E[T_{end} - T_{start} | W_i, Y_i^s = 1]$ is i 's expected length of subscription conditional on receiving treatment W_i and subscribing ($Y_i^s = 1$). The reason for the discrepancy in the gains on the two outcomes – subscription and subscription length – becomes apparent from Equation (5). If trial length affects not just subscription, but also how long a subscriber will remain loyal to the firm, then the gains in Y_i^l will be naturally different from the gains in subscription. To examine if this is true in our data, we present the summary statistics for $E[T_{end} - T_{start} | W_i, Y_i^s = 1]$ for the three trial lengths in Table A16 in Appendix I. We see that there are some small differences in this metric across the three trial lengths, which account for the differences between the gains in subscription and gains in subscription length.

Similarly, we can write the expected revenue (denoted by Y_i^r) conditional on treatment W_i as:

$$E(Y_i^r | W_i) = \Pr(Y_i^s | W_i) E[T_{end} - T_{start} | W_i, Y_i^s = 1] E[\text{Price}_i | W_i, X_i, Y_i^s = 1]. \quad (6)$$

This is similar to Equation (5), with the additional $E[\text{Price}_i | W_i, X_i, Y_i^s = 1]$ term. It suggests that trial length can influence revenues through three channels – (1) subscriptions, (2) length of subscription, and (3) price of the product subscribed. The first two were already discussed in the paragraph above. We now examine whether the products that consumers subscribe to and the prices that they pay are also a function of trial length. That is, we examine whether $E[\text{Price}_i | W_i, X_i, Y_i^s = 1]$ is indeed a function of W_i in our data. Note that the price that a subscriber pays is a function of both the product that s/he subscribes to (e.g., single product, all-inclusive bundle) as well her/his demographics (e.g., students pay lower prices for the same product). In Table A17 in Appendix I, we present the distribution of products and subscription type by trial length for all the subscribers in our data. Again, we see that there are some minor differences in product and subscription types across trial lengths, which explain the differences in revenue gains.

7.3 Optimizing on Long-run Outcomes

So far we saw that a personalized policy designed to optimize short-run conversions also does well on long-run outcomes. However, this still begs the question of how it compares to policies directly optimized to maximize long-run outcomes. In practice, the problem with using retention/revenues until some period T (e.g., two years) is that we have to wait till T to identify the best policy and then implement it. This is both sub-optimal and impractical from a firm's perspective. In contrast, using a short-term outcome such as subscriptions to design policy and then projecting the policy gains on long-term objectives (e.g., revenues) is both practical and feasible. However, a policy optimized on short-run outcomes may still perform worse than one directly optimized on long-term outcomes. Therefore, we examine and compare the performance of the policy designed to maximize subscriptions with policies designed to maximize customer loyalty or profitability, and see which performs better in our context.

Dataset	Policy optimized on	Subscription	Total Revenue	Subscription Length
Training data	Subscription	15.85	85.96	2.39
	Total Revenue	15.60	86.41	2.35
	Subscription Length	15.71	84.77	2.40
Test data	Subscription	15.62	86.67	2.36
	Total Revenue	15.45	84.28	2.33
	Subscription Length	15.53	84.86	2.35

Table 11: Expected mean of the three outcomes of interest under policies optimizing each outcome.

To that end, we now design two other personalized policies designed to maximize: (1) subscription length and (2) revenue. The policy design follows the same procedure described in 5.1, but with revenue (Y_i^r) and subscription length (Y_i^l) as our outcome variables. That is, we first estimate two separate lasso models with the above two variables as outcome variables, and then assign policy based on them.

Table 11 compares the performance of the three personalized policies on the three outcome variables of interest to the firm: subscriptions, subscription length, and revenue.¹⁵ Interestingly, we find that the policy optimized on short-run conversions (subscriptions) also performs the best on retention and revenue. There are three reasons for this. First, as we saw in the previous section, conditional on subscription, the differences in retention length and products purchased are relatively minor. Hence, optimizing on subscription is largely consistent with optimizing on retention/revenue. Second, recall that the long-run outcomes are missing for about 7% of the subscribers. So the policies based on these outcome have less information for training, which compromises their generalizability and hampers their performance on the test data. Third, subscription is a binary outcome and as such has no variance in the positive realizations. In contrast, the variance in the long-run outcome variables (subscription length and revenue) can be quite high. This variance makes it harder to generalize models based on these outcomes, which in turn adversely affects the performance of policies based on them.

In summary, our findings suggest that if there are no significant differences in customer loyalty and profitability as a function of the promotional channel through which the user converted (trial length in this case), then optimizing low-variance short-run conversions will also lead to more generalizable policies that will also perform well on long-run outcomes.

8 Conclusions

Free trials are now a commonly used promotional strategy for SaaS products and other digital experience goods. In this paper, we examine the effect of trial length on consumers' subscription and retention decisions using data from a large-scale field experiment run by a leading SaaS firm, where the firm randomly assigned new users to 7, 14, or 30 days of free trial. We leverage two unique features of the data in our study: (a) the exogenous assignment of trial length and (2) the user's post-treatment product download and usage data during the trial period.

¹⁵See Appendix J for a discussion of how the treatment allocation varies across the three policies.

We find that the shortest trial length (7-days) is the best average treatment and maximizes both short- and long-run outcomes, customer acquisition, retention, and profitability. While this result is likely to be specific to our setting, we examine the behavioral underpinning of these findings and provide some evidence on the mechanisms at play. We rule out the demand cannibalization or free riding theory, find support for the consumer learning hypothesis, and identify a novel mechanism that plays a significant role in the effectiveness of free trials – the negative effect of long stretches of inactivity at the end of the trial on subscription.

We then develop a personalized targeting policy based on lasso and show that it can lead to over 6.8% improvement in subscription compared to the baseline uniform policy of giving a 30-day trial for all. Further exploration of usage within different consumer segments in our personalization scheme suggests that simpler products and experienced users are more likely to benefit from longer trials. Finally, we find that the personalized policy designed to optimize subscriptions also performs well on long-term metrics such as customer retention and revenues in our setting. We also compare the performance of our benchmark personalized policy with alternative personalized policies developed based on other well-known outcome estimators (e.g., random forests) and the recently developed heterogeneous treatment effects estimators (e.g., generalized random forests). We find that many alternative personalized policies perform poorly, and are often worse than the simple uniform 7-days for all policy. A key managerial takeaway is that firms should not naively assume that personalization based on the most advanced estimators always helps. Instead, they should develop personalized policies based on a number of methods and carefully evaluate them using offline IPS estimators before investing resources in deploying personalized policies in the field.

Our paper opens many avenues for future research. First, while our analysis indicates that product usage during the trial period affects users' subscription decisions, we do not causally tie usage to subscriptions because usage is endogenous. Nevertheless, future research may be able to use treatment assignment (e.g., trial length) as an instrument that exogenously shifts usage and directly estimate the effect of usage on purchase. This can provide further insight into the question of whether encouraging usage (either through free trial or product improvements) can lead to better purchase outcomes. Second, our analysis suggests that personalized policies do not always perform better than a simple uniform policy. One interesting finding is that outcome estimators that have high predictive ability do not necessarily do well on personalized policy design (compare the performance of models in Table A9 in Web Appendix G with Table 8). Moreover, the finding that recently developed CATE estimators such as causal forest do not perform well in our setting is also surprising. Further investigation into the question of whether these results are generalizable would be a useful next step.

References

- E. T. Anderson and D. I. Simester. Long-run effects of promotion depth on new versus established customers: three field studies. *Marketing Science*, 23(1):4–20, 2004.
- S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- S. Athey and S. Wager. Efficient policy learning. *arXiv preprint arXiv:1702.02896*, 2017.
- M. Bruhn and D. McKenzie. In Pursuit of Balance: Randomization in Practice in Development Field Experiments. *American Economic Journal: Applied Economics*, 1(4):200–232, 2009.
- H. K. Cheng and Y. Liu. Optimal software free trial strategy: The impact of network externalities and consumer uncertainty. *Information Systems Research*, 23(2):488–504, 2012.
- D. Dey, A. Lahiri, and D. Liu. Consumer learning and time-locked trials of software products. *Journal of Management Information Systems*, 30(2):239–268, 2013.
- M. Dudík, J. Langford, and L. Li. Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, pages 1097–1104. Omnipress, 2011.
- P. S. Fader and B. G. Hardie. Probability Models for Customer-base Analysis. *Journal of Interactive Marketing*, 23(1): 61–69, 2009.
- N. Fong, Y. Zhang, X. Luo, and X. Wang. Targeted promotions on an e-book platform: Crowding out, heterogeneity, and opportunity costs. *Journal of Marketing Research*, 56(2):310–323, 2019.
- B. Foubert and E. Gijbrecchts. Try it, you’ll like it—or will you? the perils of early free-trial promotions for high-tech service adoption. *Marketing Science*, 35(5):810–826, 2016.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.
- Gartner. Forecast: Public cloud services, worldwide, 2016-2022, 4q18 update., 2019. URL <https://www.gartner.com/en/newsroom/press-releases/2019-04-02-gartner-forecasts-worldwide-public-cloud-revenue-to-g>.
- T. Guo, S. Sriram, and P. Manchanda. The effect of information disclosure on industry payments to physicians. *Available at SSRN 3064769*, 2017.
- S. Gupta, D. Hanssens, B. Hardie, W. Kahn, V. Kumar, N. Lin, N. Ravishanker, and S. Sriram. Modeling customer lifetime value. *Journal of Service Research*, 9(2):139–155, 2006.
- J. R. Hauser, G. L. Urban, G. Liberali, and M. Braun. Website morphing. *Marketing Science*, 28(2):202–223, 2009.
- G. J. Hitsch and S. Misra. Heterogeneous treatment effects and optimal targeting policy evaluation. *Available at SSRN 3111957*, 2018.
- D. G. Horvitz and D. J. Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260):663–685, 1952.
- G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- T. Kitagawa and A. Tetenov. Who should be treated? empirical welfare maximization methods for treatment choice. *Econometrica*, 86(2):591–616, 2018.
- R. A. Lewis and J. M. Rao. The unfavorable economics of measuring the returns to advertising. *The Quarterly Journal of Economics*, 130(4):1941–1973, 2015.
- C. F. Manski. Statistical treatment rules for heterogeneous populations. *Econometrica*, 72(4):1221–1246, 2004.
- D. M. McCarthy, P. S. Fader, and B. G. Hardie. Valuing subscription-based businesses using publicly disclosed customer data. *Journal of Marketing*, 81(1):17–35, 2017.
- D. McKenzie. Should we require balance t-tests of baseline observables in randomized experiments?, 2017. URL <https://blogs.worldbank.org/impactevaluations/should-we-require-balance-t-tests-baseline-observables-randomized-experiments>.
- C. F. Mela, S. Gupta, and D. R. Lehmann. The long-term impact of promotion and advertising on consumer brand choice. *Journal of Marketing research*, 34(2):248–261, 1997.
- D. C. Mutz, R. Pemantle, and P. Pham. The perils of balance testing in experimental design: Messy analyses of clean data. *The American Statistician*, 73(1):32–42, 2019.
- K. Pauwels, D. M. Hanssens, and S. Siddarth. The long-term effects of price promotions on category incidence, brand

- choice, and purchase quantity. *Journal of Marketing Research*, pages 421–439, 2002.
- R. L. Prentice. Surrogate endpoints in clinical trials: definition and operational criteria. *Statistics in medicine*, 8(4): 431–440, 1989.
- O. Rafieian. Optimizing user engagement through adaptive ad sequencing. Technical report, Working paper, 2019a.
- O. Rafieian. Revenue-optimal dynamic auctions for adaptive ad sequencing. Technical report, Working paper, 2019b.
- O. Rafieian and H. Yoganasimhan. Targeting and privacy in mobile advertising. *Marketing Science*, 2021.
- E. M. Schwartz, E. T. Bradlow, and P. S. Fader. Customer acquisition via display advertising using multi-armed bandit experiments. *Marketing Science*, 36(4):500–522, 2017.
- C. A. Scott. The effects of trial and incentives on repeat purchase behavior. *Journal of Marketing Research*, 13(3): 263–269, 1976.
- D. Simester, A. Timoshenko, and S. I. Zoumpoulis. Targeting prospective customers: Robustness of machine learning methods to typical data challenges. *Management Science*, 2019.
- D. Simester, A. Timoshenko, and S. I. Zoumpoulis. Efficiently evaluating targeting policies: Improving on champion vs. challenger experiments. *Management Science*, 66(8):3412–3424, 2020.
- T. Sunada. Customer learning and revenue-maximizing trial design. Technical report, Working Paper, 2018.
- A. Swaminathan and T. Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *International Conference on Machine Learning*, pages 814–823. PMLR, 2015.
- A. Swaminathan, A. Krishnamurthy, A. Agarwal, M. Dudik, J. Langford, D. Jose, and I. Zitouni. Off-policy evaluation for slate recommendation. In *Advances in Neural Information Processing Systems*, pages 3632–3642, 2017.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- T. J. VanderWeele. Surrogate measures and consistent surrogates. *Biometrics*, 69(3):561–565, 2013.
- S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 0(0):1–15, 2018. doi: 10.1080/01621459.2017.1319839.
- S. Wang and G. F. Özkan-Seely. Signaling product quality through a trial period. *Operations Research*, 66(2):301–312, 2018.
- J. Yang, D. Eckles, P. Dhillon, and S. Aral. Targeting for long-term outcomes. *arXiv preprint arXiv:2010.15835*, 2020.
- H. Yoganasimhan. Search personalization using machine learning. *Management Science*, 66(3):1045–1070, 2020.
- M. Zhu, Y. Yang, and C. K. Hsee. The mere urgency effect. *Journal of Consumer Research*, 45(3):673–690, 2018.

Appendices

A Appendix for Data Section

Trial Length	Variable	Mean	Standard Deviation	Min	25%	50%	75%	Max	Number of Observations
7 Days	Subscription	0.154	0.360	0.00	0.00	0.00	0.00	1.00	51,017
	Subscription length (all)	2.35	6.59	0.00	0.00	0.00	0.00	73	50,504
	Subscription length (subscribers)	16.19	8.70	0.0	10.00	18.00	23.00	73.00	7,322
	Revenue (all)	83	295	0	0	0	0	6,219	50,504
	Revenue (subscribers)	578	562	0	239	424	679	6,219	7,322
14 Days	Subscription	0.150	0.357	0.0	0.00	0.00	0.00	1.00	51,040
	Subscription length (all)	2.27	6.45	0.0	0.00	0.00	0.000	67	50,543
	Subscription length (subscribers)	16.09	8.48	0.0	10.00	18.00	22.00	67.00	7,138
	Revenue (all)	81	289	0	0	0	0	6,799	50,543
	Revenue (subscribers)	576	554	0	252	420	679	6,799	7,138
30 Days	Subscription	0.147	0.354	0.00	0.00	0.00	0.00	1.00	235,667
	Subscription length (all)	2.20	6.31	0.00	0.00	0.00	0.00	108	233,176
	Subscription length (subscribers)	15.96	8.36	0.00	10.00	17.00	22.00	108	32,073
	Revenue (all)	77	281	0	0	0	0	20,208	233,176
	Revenue (subscribers)	564	550	0	241	420	660	20,208	32,073

Table A1: Summary statistics of subscription, subscription length, and revenue outcomes in each trial length. All the revenue numbers are scaled by a constant factor to preserve the firm's anonymity.

Trial Length	Variable	Mean	Standard Deviation	Min	25%	50%	75%	Max	Number of Observations
7 Days	Total downloaded packages	1.0	0.0	1.0	1.0	1.0	1.0	4.0	51,017
	Indicator for software use	1.0	0.0	0.0	0.0	1.0	1.0	1.0	51,017
	Number of active days	2.0	2.0	0.0	1.0	1.0	2.0	7.0	45,810
	Ratio of active days	0.0	0.0	0.0	0.0	0.0	0.0	1.0	45,810
	Usage during trial	985	3,567	0.0	37	184	730	223,179	45,810
	Log usage during trial	5.0	3.0	0.0	4.0	5.0	7.0	12.0	45,810
	Dormancy length	1.0	0.0	0.0	0.0	1.0	1.0	1.0	45,810
14 Days	Total downloaded packages	1.0	0.0	1.0	1.0	1.0	1.0	4.0	51,040
	Indicator for software use	1.0	0.0	0.0	0.0	1.0	1.0	1.0	51,040
	Number of active days	2.0	3.0	0.0	1.0	1.0	3.0	14.0	45,902
	Ratio of active days	0.0	0.0	0.0	0.0	0.0	0.0	1.0	45,902
	Usage during trial	1,397	5,938	0.0	44	227	919	400,705	45,902
	Log usage during trial	5.0	3.0	0.0	4.0	5.0	7.0	13.0	45,902
	Dormancy length	1.0	0.0	0.0	0.0	1.0	1.0	1.0	45,902
30 Days	Total downloaded packages	1.0	0.0	1.0	1.0	1.0	1.0	4.0	235,667
	Indicator for software use	1.0	0.0	0.0	1.0	1.0	1.0	1.0	235,667
	Number of active days	3.0	4.0	0.0	1.0	2.0	4.0	30.0	211,802
	Ratio of active days	0.0	0.0	0.0	0.0	0.0	0.0	1.0	211,802
	Usage during trial	1,968	8,007	0.0	51	286	1,227	488,666	211,802
	Log usage during trial	5.0	3.0	0.0	4.0	6.0	7.0	13.0	211,802
	Dormancy length	1.0	0.0	0.0	0.0	1.0	1.0	1.0	211,802

Table A2: Summary statistics for usage features in each trial length.

Trial Length	Training Data	Test Data	Total
7 days	35,743	15,274	51,017
14 days	35,901	15,139	51,040
30 days	165,056	70,611	235,667
Total	236,700	101,024	337,724

Table A3: The number of observations for each trial length in each dataset.

B Appendix for ATE and Randomization Checks

B.1 ATE using Regressions

	Training Dataset		Test Dataset		All Dataset	
	No control	Control	No control	Control	No control	Control
7 Days	0.0064 (0.0021)	0.0064 (0.0018)	0.0082 (0.0032)	0.0069 (0.0027)	0.0069 (0.0017)	0.0065 (0.0015)
14 Days	0.0021 (0.0021)	0.0018 (0.0018)	0.0048 (0.0032)	0.0038 (0.0027)	0.0029 (0.0017)	0.0024 (0.0015)

Table A4: The coefficients of 7 days and 14 days in regression analysis. In the control case, we control for all the pre-treatment variables.

B.2 Randomization Checks

First, we show the distribution of treatment assignment across the five skill levels in Table A5. As we can see, the treatment assignment is pretty even; and none of the differences in assignment probabilities within each sub-category are significant. However, as discussed in the main text, this approach of testing for randomness in assignment within each sub-category is not considered the best practice for a variety of reasons; see Bruhn and McKenzie (2009) and Mutz et al. (2019).

Trial Length	Beginner	Experienced	Intermediate	Mixed	Unknown
7 days	0.6899	0.1295	0.0005	0.1077	0.0723
14 days	0.6887	0.1282	0.0006	0.1082	0.0743
30 days	0.6910	0.1269	0.0006	0.1075	0.0740

Table A5: The proportion of individuals with a given skill in different trial lengths.

Therefore, we present two more formal randomization checks. First, we regress the treatment variable (W_i) on the pre-treatment variables separately for both the 7 day and 14 day treatment using the 30 day treatment as the base in both regressions. We present these results in Table A6. The top panel shows the results of the regressions for the training data and the bottom panel shows the results for the test data. In all the four regressions, we see that the pre-treatment variables have no predictive power (all the R-squareds are zero) and F-statistics are not significant. This suggests that the pre-treatment variables are not systematically correlated with treatment assignment. Second, we regress the outcome variable (subscription decision) on

the treatment variable and all the pre-treatment variables and present the results in Table A7. The treatment effects are very similar to those in Table 5, which suggests that there are no issues with randomization.

Dataset	Comparison Treatment	R-Squared	Adj. R-squared	F-statistic	p-value of the F-statistic
Training	7 days	0.000377	0.000019	1.052991	0.356526
	14 days	0.000439	0.000081	1.225721	0.093690
Test	7 days	0.000865	0.000026	1.030504	0.406522
	14 days	0.000963	0.000125	1.149222	0.181734

Table A6: The results of regressing the treatment (W_i) on the pre-treatment variables. In all the regressions, the baseline is the 30 day treatment. The F-statistic tests the null hypothesis that coefficients are jointly statistically significantly different from zero.

Dataset	Trials Length	Coefficient	std err	t-stat	$P > t $	[0.025	0.975]
Training	14 days	0.0018	0.002	0.999	0.318	-0.002	0.005
	7 days	0.0064	0.002	3.581	0.000	0.003	0.010
Test	14 days	0.0038	0.003	1.385	0.166	-0.002	0.009
	7 days	0.0069	0.003	2.516	0.012	0.002	0.012

Table A7: Effect of trial length calculated by regressing the subscription decision on all the pre-treatment variables and trial length.

C Appendix for Mechanism

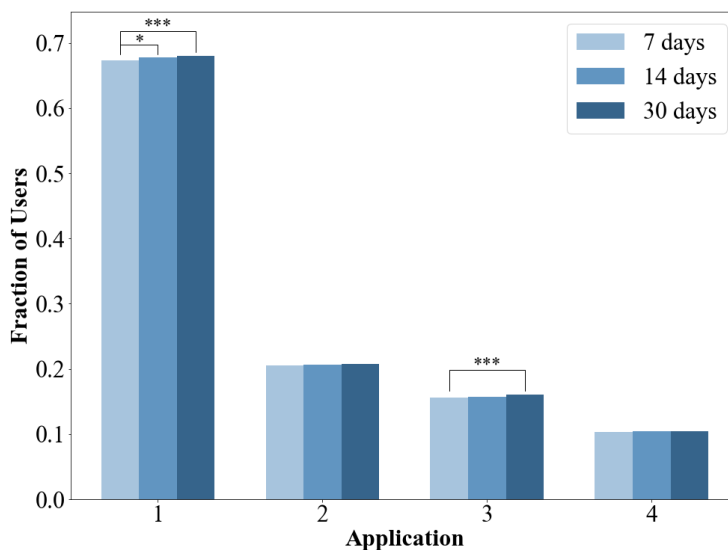


Figure A1: Fraction of users downloading products 1–4 under each trial length. Products 1 and 3 are complements.

D Lasso Implementation Details

Next, we briefly describe our implementation of the lasso model discussed in $\chi??$. First, note that we use the subscription indicator as outcome variable (Y_i). Next, in our setting, all the pre-treatment variables, X_i s, are categorical. So we transform each categorical variable into a set of dummy variables for each sub-category (also referred to as one-hot encoding in the machine learning literature). After this transformation, we have 82 dummies for the pre-treatment variables, three treatment dummies, and 246 first-order interaction variables. This gives us a total of 331 explanatory variables for the lasso model. Recall that there is one hyper-parameter that we need to tune in the lasso model, which is the L1 regularization parameter, λ . We use the standard cross-validation procedure implemented in the *glmnet* package in R. It searches over 98 different values of λ ranging from $1.8 \cdot 10^{-5}$ to $1.5 \cdot 10^{-1}$, and picks the one that gives us the best out-of-sample performance (based on five-fold cross-validation). In our case, it $\lambda = 3.1 \cdot 10^{-4}$. The final model achieves a MSE of 0.0933 in both the training and test data.

E Design of Alternative Personalized Policies

We discuss the design of a series of alternative personalized policies. First, in λ E.1, we present with the policies based on other outcome estimators. Next, in λ E.2, we discuss policies based on CATE estimators.

E.1 Policies based on Outcome Estimators

Recall that to design outcome estimators, we first need learn a model $f(x, w) = E[Y|X_i = x, W_i = w]$. Then, using estimate of the expected outcome, $\hat{f}(x = X_i, w)$, we obtain the optimal personalized policy (based on outcome estimator f) for observation i as:

$$\pi_f(X_i) = w, \quad \text{where } w = \arg \max_{w \in \mathcal{W}} \hat{f}(x = X_i, w) \quad (\text{A.1})$$

We now consider four commonly used outcome estimators (or models for f): (1) linear regression, (2) CART, (3) random forest, and (4) boosted regression trees. We focus on these because of a few key reasons. First, linear regression is the simplest and most commonly used method to model any outcome of interest. Second, lasso is worth exploring because it was designed to improve on the predictions from a linear regression by reducing the out-of-sample variance using variable selection. Third, CART is a semi-parametric way to model outcomes by partitioning the covariate space into areas with the highest within-group similarity in outcomes. Finally, both random forest and boosted regression trees are improvements of CART and have been shown to offer much higher predictive ability than simpler models such as CART, regression, or lasso. We discuss each of these briefly below.

E.1.1 Linear Regression

A linear regression model with first-order interaction effects can be used to predict an individual i 's observed outcome Y_i as a function of her pre-treatment variables, treatment variable, and the interaction of both, as follows:

$$Y_i = X_i\beta_1 + W_i\beta_2 + X_iW_i\beta_3 + \epsilon_i. \quad (\text{A.2})$$

β_1 is a vector that captures the effect of the pre-treatment variables (X_i) on the outcome. β_2 is a vector that captures the main effect of the different treatments on Y . Finally, the vector β_3 captures the interaction effect of treatment and pre-treatment variables. This interaction term is important because it helps us in personalizing the policy by capturing how the effectiveness of different treatments varies across individuals (as a function of their pre-treatment attributes).

However, in a high dimensional covariate space, linear regressions with first-order interaction effects will have a large number of explanatory variables. This usually leads to poor out-of-sample fit. That is, such regressions tend to have low bias, but high variance – they tend to overfit on the training sample but perform poorly in new samples, especially if the data generating process is noisy. Since our goal is to design optimal policies for counterfactual situations, the out-of-sample fit is the only metric of importance.¹ The lasso model

¹We can also add higher-order interaction effects into the regression model. However, we refrain from doing so because it will increase model complexity significantly and exacerbate this problem.

used in the main text addresses the above problem by learning a simpler model that uses fewer variables (Tibshirani, 1996; Friedman et al., 2010).

E.1.2 Tree-based Methods

Next, we discuss tree-based models, starting with CART. CART recursively partitions the covariate space into sub-regions, and the average of Y in a given region ($E(Y)$) is the predicted outcome for all observations in that region (Breiman et al., 1984). This type of partitioning can be represented by a tree structure, where each leaf of the tree represents an output region. A general CART model can be expressed as:

$$y = f(x, w) = \sum_{m=1}^M \rho_m I(x, w \in R_m), \quad (\text{A.3})$$

where x, w is the set of explanatory variables, R_m is the m^{th} region of the M regions used to partition the space, ρ_m is the predicted value of y in region m .

Trees are trained by specifying a cost function (e.g., mean squared error) that is minimized at each step of the tree-growing process using a greedy algorithm. To avoid over-fitting, usually a penalty term is added to the cost function, whose weight is a hyper-parameter learnt from the data using cross-validation.

While CART has some good properties², it often has poor predictive accuracy because it is discontinuous in nature and is sensitive to outliers. Therefore, even with regularization, CART models tend to over-fit on the training data and under-perform on out-of-sample data. We can address these problems using two different techniques – (1) Bagging and (2) Boosting. We discuss both of these techniques and the resulting estimators below.

In general, deep trees have high in-sample fit (low bias), but high out-of-sample variance because of over-fitting. However, we can improve the variance problem by bagging or averaging deep trees using bootstrapping. Ho (1995) formalized this idea and proposed random forests. Random forest usually consists of hundreds or thousands of trees, each of which is trained on a random sub-sample of columns and rows. Each tree is thus different from other trees and the average of these random trees is a better predictor than one single tree trained on the full data.

Another approach is to start with shallow trees. Shallow trees have poor in-sample fit (high bias), but low out-of-sample variance. Additively, adding a series of weak trees that minimize the residual error at each step by a process known as boosting can improve the bias problem, while retaining the low variance. This gives us boosted regression trees. Conceptually, boosting can be thought of as performing gradient descent in function space using shallow trees as the underlying weak learners (Breiman, 1998; Friedman, 2001).³ In this paper, we use XGBoost, a version of boosted trees proposed by Chen and Guestrin (2016) because it is

²CART can accept both continuous and discrete explanatory variables, is not sensitive to the scale of the variables, and allows any number of interactions between features (Murphy, 2012). It therefore has ability to capture rich nonlinear patterns in the data. Further, CART can do automatic variable selection, i.e., it uses only those variables that provide better accuracy in the prediction task for splitting.

³It should be noted that bagging is not a modeling technique; it is simply a variance reduction technique. Boosting, however, is a method to infer the underlying model $y = f(x; w)$. Thus, they are conceptually completely different.

superior to earlier implementations both in terms of accuracy and scalability.

Please see Appendix A.F for the set of hyper-parameters that need to be tuned in CART, random forest, and XGBoost.

E.2 Heterogeneous Treatment Effect Estimators

The second approach to designing a personalized policy is to first obtain consistent estimates of heterogeneous treatment effects for each pair of treatments for all users, and then use them to assign treatments. This method also follows a two-step procedure. In the first step, we can obtain consistent estimates of individual-level treatment effects, $\tau_{w_j, w_{j'}}(x)$, for each pair of treatments $w_j, w_{j'} \in \{0, \dots, W-1\}$. Under the standard assumptions of (1) unconfoundedness, (2) SUTVA, and (3) positivity, $\tau_{w_j, w_{j'}}(x)$, can be defined as (Rubin, 1974; Imbens and Rubin, 2015):

$$\tau_{w_j, w_{j'}}(x) = \mathbb{E} \left[Y(X_i, W_i = w_j) - Y(X_i, W_i = w_{j'}) \mid X_i = x \right]. \quad (\text{A.4})$$

So if we have W treatments, we have to build $\frac{W(W-1)}{2}$ pairwise models, where each model gives us an estimate of individual-level treatment effects for a given pair of treatments, $w_j, w_{j'}$.

In the second step, we use the estimated treatment effects to derive the optimal policy as⁴:

$$\pi(X_i) = w_j \quad \text{if and only if} \quad \forall j' \neq j \quad \tau_{w_j, w_{j'}}(x = X_i) \geq 0 \quad (\text{A.5})$$

Next, we discuss these methods are based on the potential outcomes framework and estimate the treatment effect for any pair of treatments $(w_j, w_{j'})$ at each point x as shown in Equation (A.4). However, this equation is not very useful in practical settings (where the covariate space is high-dimensional and data are finite) because we would not have sufficient observations at each X_i to estimate precise treatment effects. Therefore, the general idea behind modern heterogeneous treatment effects estimators is to pool observations that are close in the covariate space and estimate the conditional average treatment effect for sub-populations instead of estimating the treatment effect at each point in the covariate space. That is, we can modify Equation (A.4) as:

$$\tau_{w_j, w_{j'}}(x) = \frac{\sum_{X_i \in l(x), W_i = w_j} Y_i}{\sum \mathbb{1}[X_i \in l(x), W_i = w_j]} - \frac{\sum_{X_i \in l(x), W_i = w_{j'}} Y_i}{\sum \mathbb{1}[X_i \in l(x), W_i = w_{j'}]}, \quad (\text{A.6})$$

where $l(x)$ is the set of covariates that are fairly similar to x . Intuitively, for each point x , we use the observations in $l(x)$ to estimate treatment effects.

The main question in these methods then becomes how to find the optimal $l(x)$ around each x . On the one hand, if $l(x)$ is too small, then we will not have sufficient observations within $l(x)$, which would result in

⁴In practice, we can have situations where three or more pairs of the estimated treatments form a loop, i.e., $\hat{\tau}_{w_j, w_{j'}}(x = X_i) > \hat{\tau}_{w_{j'}, w_{j''}}(x = X_i) > \hat{\tau}_{w_{j''}, w_j}(x = X_i)$. This happens if the estimated treatment effects are noisy (usually due to lack of sufficient data). For these observations, we can simply assign the best average treatment (across all observations) as the policy-prescribed treatment. Further, there can be multiple optimal policies because two or more top treatments can have the same effect on the outcome. That is, for some combinations of X_i and $j' \neq j$, we can have: $\tau_{w_j, w_{j'}}(x = X_i) = 0$. However, notice that all the solutions give the same reward, i.e., the optimal reward is unique.

noisy estimates of treatment effects. On the other hand, if $l(x)$ is too large, then we will not capture all the heterogeneity in the treatment effects, which is essential for personalizing policy. Indeed, if $l(x)$ is the entire data, then we simply have one Average Treatment Effect for all users (which will give us one global policy). Thus, finding the optimal $l(x)$ involves the classic bias-variance trade-off.

Starting with Rzepakowski and Jaroszewicz (2012), a growing stream of literature has focused on developing data-driven approaches to finding the optimal $l(x)$ based on ideas from the machine learning literature. Among these methods, the recently developed causal tree and causal forest (or Generalized Random Forest) have been shown to have superior performance. So we focus on them.

E.2.1 Causal Tree

Causal tree builds on CART. The only difference is that instead of partitioning the space to maximize predictive ability, the objective here is to identify partitions with similar within-partition treatment effect. Athey and Imbens (2016) show that maximizing the variation in the estimated treatment effects (with a regularization term added to control complexity) achieves this objective. Their algorithm consists of two steps. In the first step, it recursively splits the covariate space into partitions. In the second step, it estimates the treatment effect within each partition ($l(x)$) using Equation (A.6). Intuitively, this algorithm pools observations with similar treatment effects into the same partition because splitting observations that have similar treatment effects does not increase the objective function (i.e., variation in the post-split treatment effect).

The main idea behind causal tree is that we can use recursive partitioning to estimate heterogeneous treatment effects if we can come up with an objective function that can identify partitions with similar within-partition treatment effect.⁵ Athey and Imbens (2016) show that maximizing the variation in the estimated treatment effects achieves this objective, which can be written as:

$$Var[\hat{\tau}(X)] = \frac{1}{N} \sum_{i=1}^N \hat{\tau}^2(X_i) - \left(\frac{1}{N} \sum_{i=1}^N \hat{\tau}(X_i) \right)^2 \quad (\text{A.7})$$

Since $\left(\frac{1}{N} \sum_{i=1}^N \hat{\tau}(X_i) \right)^2$ remains constant with respect to any possible next split, the objective function can be also described as maximizing the average of the square of estimated treatment effects. In practice, the algorithm chooses the split that maximizes $Var[\hat{\tau}(X)] - \zeta T$, where the second term is added for complexity control and is analogous to the regularization term in CART.

The algorithm consists of two steps. In the first step, it recursively splits the covariate space into partitions. In the second step, it estimates the treatment effect within each partition ($l(x)$) using Equation (A.6). Intuitively, this algorithm pools observations with similar treatment effects into the same partition because splitting observations that have similar treatment effects does not increase the objective function (i.e.,

⁵The first example of such a criterion was proposed by Hansotia and Rukstales (2002): for each potential split, the algorithm calculates the difference between the average outcome of the treatment and control in the right (Y_r) and left (Y_l) branches of the tree. Then, it selects the split with the highest value of $Y = j - Y_r - Y_l$. Other splitting criteria include maximizing the difference between the class distributions in the treatment and control groups (Rzepakowski and Jaroszewicz, 2012).

variation in the post-split treatment effect).

E.2.2 Generalized Random Forest

The causal tree algorithm nevertheless suffers from weaknesses that mirror those of CART. To resolve these issues, Wager and Athey (2018) proposes causal forest, which builds on random forests (Breiman, 2001). More broadly, Athey et al. (2019) show that the intuition from random forests can be used to flexibly estimate any heterogeneous quantity, $\theta(x)$, from the data, including heterogeneous treatment effects. They suggest a Generalized Random Forest (GRF) algorithm that learns problem-specific kernels for estimating any quantity of interest at a given point in the covariate space. For treatment estimation, the method proceeds in two steps. In the first step, it builds trees whose objective is to increase the variation in the estimated treatment effects. Each tree is built on a random sub-sample of the data and random sub-sample of covariates. In the second step, the algorithm uses the idea of a weighted kernel regression to calculate the treatment effect at each point x using weights from the first step.

The Generalized Random Forest (GRF) algorithm takes intuition from causal tree and combines it with ideas from predictive random forests to learn problem-specific kernels for estimating any quantity of interest at a given point in the covariate space. For the estimation of heterogeneous treatment effects, the method proceeds in two steps.

In the first step, it builds trees whose objective is to increase the variation in the estimated treatment effects. Each tree is built on a random sub-sample of the data and random sub-sample of covariates. At each step of the partitioning process (when building a tree), the algorithm first estimates the treatment effects in each parent leaf P by minimizing the R-learner objective function. This objective function is motivated by Robinson’s method (Robinson, 1988; Nie and Wager, 2017) and can be written as follows:

$$\hat{\tau}_P(\cdot) = \arg \min_{\tau} \left[\frac{1}{n_P} \sum_{i=1}^{n_P} \left((Y_i - \hat{m}^{(i)}(X_i)) - (W_i - \hat{\ell}^{(i)}(X_i)) \hat{\tau}(X_i) \right)^2 + n(\tau(\cdot)) \right], \quad (\text{A.8})$$

where n_P refers to the number of observations in the parent partition, $n(\tau(\cdot))$ is a regularizer that determines the complexity of the model used to estimate $\tau_P(\cdot)$ s, and $\hat{m}^{(i)}(X_i)$ and $\hat{\ell}^{(i)}(X_i)$ are out-of-bag estimates for the outcome and propensity score, respectively. While any method can be used for estimating $\hat{m}^{(i)}(X_i)$ and $\hat{\ell}^{(i)}(X_i)$, the GRF implementation uses random forests to estimate these values.

Then, it chooses the next split such that it maximizes the following objective function:

$$\frac{n_L}{n_P} \frac{n_R}{n_P} (\hat{\tau}_L - \hat{\tau}_R)^2, \quad (\text{A.9})$$

where n_L and n_R refer to the number of observations in the post-split left and right partitions, respectively. However, instead of calculating the exact values of $\hat{\tau}_L$ and $\hat{\tau}_R$ for each possible split (and for each possible parent), the causal forest algorithm uses a gradient-based approximation of $\hat{\tau}$ for each child node to improve compute speed; see Athey et al. (2019) for details. Thus, at the end of the first step, the method calculates weights, $\alpha_i(x)$, that denote the frequency with which the i -th training sample falls into the same leaf as x

in the first step. Formally, $\alpha_i(x)$ is given by $\alpha_i(x) = \frac{1}{B} \sum_{b=1}^B \alpha_{bi}(x)$, where $\alpha_{bi}(x) = \frac{\mathbf{1}(X_i \in l_b(x))}{|l_b(x)|}$, B is the total number of trees built in the first step, and $l_b(x)$ is the partition that x belongs to in the b -th tree.

In the second step, the algorithm uses the idea of a weighted kernel regression to calculate the treatment effect at each point x using weights $\alpha_i(x)$ as follows:

$$\hat{\tau}(x) = \frac{\sum_{i=1}^N \alpha_i(x) (Y_i - \hat{m}^{(i)}(X_i)) (W_i - \hat{\ell}^{(i)}(X_i))}{\sum_{i=1}^N \alpha_i(x) (W_i - \hat{\ell}^{(i)}(X_i))^2}, \quad (\text{A.10})$$

As with all supervised learning models, we need to do hyper-parameter optimization to prevent causal forest from overfitting. We refer readers to Appendix [XF](#) for details on this.⁶

⁶Athey and Imbens (2016) propose an additional approach to avoid over-fitting in causal tree and causal forest – honest splitting. The main idea behind honesty is to split the data into two parts and use one part for growing each tree (i.e., generating the partitions) and the other part for estimating the treatment effects given a partition. Since the two data-sets are independent of one another, there is a lower likelihood of over-fitting. However, honest splitting comes with its own costs – it reduces the amount of data we have for learning in both stages by half. In settings where the magnitude of the treatment effects is small (such as ours), honest splitting can adversely affect the performance of models. Indeed, we found that models based on honest splitting lead to worse policies compared to models without honest splitting in our setting. So we do not employ honest splitting in our analysis.

F Hyper-parameter Optimization for the Models Estimated

For each alternative model that we use, we describe the hyper-parameters associated with it and the optimal hyper-parameters that we derive after tuning. In all cases, we use five-fold cross-validation to optimize the hyper-parameters. We then train a model on the entire training data using the optimal hyper-parameters, and report the performance of this model on both the training and test data.

Linear regression does not use any hyper-parameters and hence does not require validation. In this case, we simply train the model on the full training data to infer the model parameters and report the model's performance on both the training and test data.

For the remaining tree-based models, we directly feed in the 82 dummy variables for the pre-treatment characteristics and the three treatment dummies. Specifically for CART, we use the package *rpart* in R, which implements a single tree proposed by (Breiman et al., 1984). We only need to pick the complexity parameter (ζ) using cross validation in this case. We search over 3886 different values for ζ ranging from $8.6 \cdot 10^{-11}$ to $1.7 \cdot 10^{-1}$, and derive the optimal complexity parameter as $5.4 \cdot 10^{-5}$.

For Random Forest, we use the package *sklearn* in Python. There are three hyper-parameters in this case – (1) n_{tree} , the number of trees over which we build our ensemble forest, (2) \max_f , the maximum number of features the algorithm try for any split (it can be either all the features or the squared root of the number of features), and (3) n_{min} , the minimum number of samples required to split an internal node.

The standard method for finding hyper-parameters is grid-search. However, grid-search is very costly and time-consuming when we have to tune many hyper-parameters. So we use the hyperopt package for tuning the hyper-parameters in this model. Hyperopt provides an automated and fast hyper-parameter optimization procedure that is less sensitive to researcher's choice of searched hyper-parameter values; see Bergstra et al. (2011, 2013) for details.

For each of these hyper-parameters, we now define the range over which we search as well as the optimal value of the hyper-parameter are shown below:

$$\begin{aligned} n_{tree} &\in [100, 1200] \text{ and } n_{tree} = 1000 \\ \max_f &\in \{n, \sqrt{n}\} \text{ and } \max_f = n \\ n_{min} &\in [10, 300] \text{ and } n_{min} = 70 \end{aligned}$$

XGBoost also has many hyper-parameters that need tuning. However, we found that our results is sensitive to only three of the parameters: α , η , and d_{max} . The first parameter, α , is an L1 regularization parameter, η is the shrinkage parameter or learning rate, d_{max} is maximum depth of trees. Again, we use the hyperopt package to search over a wide range of parameter values. The optimal values are shown below:

$$\begin{aligned} \alpha &\in \{0.1, 0.2, 0.5, 1, 2, 5, 10, 15, 20, 25\} \text{ and } \alpha = 20 \\ \eta &\in [0, 1] \text{ and } \eta = 0.59 \\ d_{max} &\in \{6, 8, 10, 12\} \text{ and } d_{max} = 12 \end{aligned}$$

Causal tree has two hyper-parameters that needs tuning – (1) the complexity parameter (ζ) and the

minimum number of treatment and control observations in each leaf (q). We use the cross-validation procedure in the "causalTree" package in R for tuning ζ . We manually tune q using grid-search over the range [100, 1000] in increments of 100. We search over all possible values of ζ for each q .

The optimal hyper-parameters for the three trees (one for each pair of treatments) are:

The tree for 7 and 14 days pair: $\zeta = 1.8e-05$ and $q = 100$.

The tree for 7 and 30 days pair: $\zeta = 8.0e-06$ and $q = 100$.

The tree for 14 and 30 days pair: $\zeta = 3.0e-06$ and $q = 900$.

Causal forest has five hyper-parameters that need to be tuned⁷: (i) $frac$, the fraction of data that is used for training each tree, (ii) $mtry$, the number of variables tried for each split, (iii) max_imb the maximum allowed imbalance of a split, (iv) imb_pen , a penalty term for imbalanced splits, (v) q , the minimum number of observations per condition (control, treatment) in each partition.⁸

We used the hyper-parameter optimization procedure available in the `grf` package for tuning these hyper-parameters. The optimal hyper-parameters for each model are shown below:

7-14 days pair: $frac = 0.5$, $max_imb = 0.11$, $imb_pen = 2.03$, $mtry = 13$, and $q = 1651$.

7-30 days pair: $frac = 0.5$, $max_imb = 0.13$, $imb_pen = 2.49$, $mtry = 1$, and $q = 4$.

14-30 days pair: $frac = 0.5$, $max_imb = 0.20$, $imb_pen = 5.11$, $mtry = 7$, and $q = 121$.

⁷For more information please visit <https://github.com/grf-labs/grf/blob/master/REFERENCE.md>

⁸The GRF algorithm estimates outcomes and propensities in the first step. However, in our setting, we do not need to estimate propensity scores since they are known and constant for each observation since treatment assignment is fully randomized. Our qualitative results remain the same if we instead estimate them from the data.

G Appendix for Empirical Results on Alternative Personalized Policies

Policy-prescribed Treatment	7	14	30	<i>reg</i>	<i>lasso</i>	<i>cart</i>	<i>r-forest</i>	<i>xgboost</i>	<i>c.tree</i>	<i>c-forest</i>
7 Days	1	0	0	0.521	0.689	1	0.444	0.697	1	0.911
14 Days	0	1	0	0.322	0.232	0	0.269	0.185	0	0.089
30 Days	0	0	1	0.157	0.079	0	0.287	0.118	0	0
Total	1	1	1	1	1	1	1	1	1	1

Table A8: Fraction of users assigned the 7-, 14-, and 30-day trials under counterfactual policies (in test data).

First, in Table A8, we show the distribution of the three treatments (7, 14, and 30 days) varies under the alternative policies.

Method	Mean Squared Error	
	Training Set	Test Set
Linear Regression	0.0932	0.0933
Lasso	0.0933	0.0933
CART	0.0916	0.0920
Random Forest	0.0904	0.0915
XGBoost	0.0905	0.0911

Table A9: Comparison of the predictive performance of the five outcome estimation methods. The MSE for any method are calculated as $\frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N}$, where \hat{y}_i is the prediction of y_i and N is the number of data-points in the data-set being considered.

Next, in Table A9, we present the MSE for the five outcome estimators on both training and test data. We find that linear regression and lasso are the worst in terms of model-fit, with XGBoost performing the best. This finding is consistent with earlier papers that have found XGBoost to be the best outcome prediction method in tasks involving prediction of human behavior (Rafieian and Yoganasimhan, 2020; Rafieian, 2019). However, the tree-based models – CART, Random Forest, and XGBoost – suffer from over-fitting in spite of hyper-parameter tuning. That is, their performance on the test data is worse than that on training data. For the heterogeneous treatment effects estimators, we cannot compare the estimates of treatment effects with any ground truth since we (as researchers/managers) do not know the true treatment effects.

Finally, in Figure A2, we show the CDFs of the estimated CATEs for the 7 vs. 14 day and 14 vs. 30 treatments for all the methods.

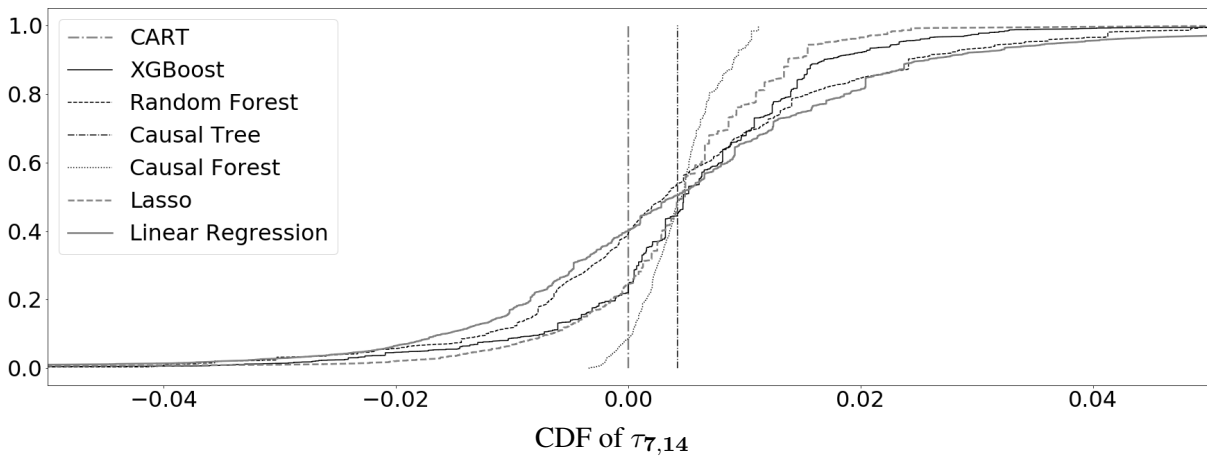
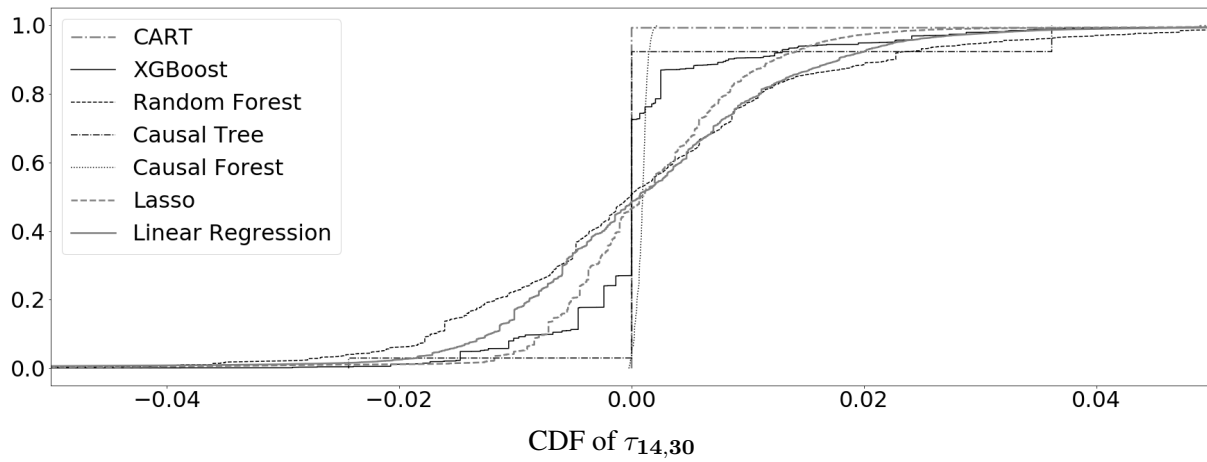


Figure A2: The CDF of estimated CATEs for 7 vs 14 days, and 14 vs 30 days of free trial from using different methods (for test data).

H Appendix for λ 6

This section is organized as follows. First, in λ H.1, we quantify the heterogeneity in the effect of trial length on post-treatment usage. Next, in λ H.2, we quantify the heterogeneity in the effect of usage on subscription. Then, in λ H.3, we show how these heterogeneous responses are consistent with the user’s pre-treatment demographic variables. Finally, in λ H.4, we provide additional evidence to rule out the demand cannibalization hypothesis.

H.1 Heterogeneity in Post-treatment Usage Across Segments

Outcome variable	Intercept	14-days optimal segment	30-days optimal segment	R^2	N
Total downloaded products	1.135 (0.001)	0.065 (0.002)	0.002 (0.003)	0.005	337,724
Indicator for software use	0.833 (0.001)	0.015 (0.002)	-0.03 (0.003)	0.001	303,514
Number of active days	2.75 (0.008)	1.116 (0.017)	0.445 (0.027)	0.014	303,514
Log usage during trial	4.992 (0.006)	0.43 (0.012)	-0.034 (0.019)	0.004	303,514
Dormancy length	17.465 (0.024)	-2.304 (0.049)	-0.953 (0.077)	0.007	303,514

Table A10: Regressions of different usage variables on the users’ optimal trial length. Each row denotes a separate regression, with the first column showing the outcome variable of that regression. Standard errors in parentheses.

Variable	Number of active days	Log usage during trial	Dormancy length
Intercept	1.676 (0.022)	4.732 (0.015)	4.732 (0.051)
14-days trial	0.519 (0.03)	0.148 (0.022)	5.343 (0.072)
30-days trial	1.427 (0.024)	0.341 (0.017)	17.09 (0.056)
14-days optimal segment	0.297 (0.044)	0.222 (0.031)	-0.422 (0.103)
14-days optimal segment 14-days trial	0.41 (0.062)	0.175 (0.044)	-0.791 (0.145)
14-days optimal segment 30-days trial	1.086 (0.048)	0.261 (0.034)	-2.519 (0.114)
30-days optimal segment	0.05 (0.068)	-0.158 (0.049)	-0.079 (0.162)
30-days optimal segment 14-days trial	0.176 (0.097)	0.118 (0.068)	-0.394 (0.228)
30-days optimal segment 30 days trial	0.523 (0.075)	0.15 (0.053)	-1.21 (0.178)
R -squared	0.044	0.008	0.346
Number of observations	303,514	303,514	303,514

Table A11: Regression of different usage features on the interaction of trial length, and optimal trial length. Standard errors in parentheses.

We start with Table A10, which shows how the three segments differ in their usage behavior based on five different regressions. Each regression in this table uses a specific usage feature as the outcome variable and the user’s segment (or optimal treatment) as the explanatory variable. The 7-day optimal segment is the baseline in all the regressions (and is denoted by the intercept).⁹ We find that the 14-day optimal users download and use the product more compared to the 7-day optimal users (positive coefficients in the regressions using log usage and number of active days) and they tend to remain active longer than both the 7- and 30-day optimal users (negative coefficient of dormancy length). In contrast, the 30-day optimal users do

⁹We do not control for the actual treatment assignment here since it is randomly assigned, and therefore orthogonal to a user’s optimal treatment.

not use the product significantly more than the 7-day optimal users, but they have somewhat shorter dormancy periods. Overall, this suggests that 7-day optimal users use the product the least and become inactive the soonest, while the 14-day optimal use it the most and are active for the longest.

Next, we examine how these three segments behave under different trial lengths. In Table A11, we present the results from three regressions, where the outcome variable in each regression is a usage feature (Number of active days, Log usage, or Dormancy length) and the explanatory variables are the user's optimal treatment, the actual treatment assigned to her, and the interaction effects. In all these regressions, the intercept refers to the baseline of 7-day optimal users when they are assigned the 7-day trial. We find that, compared to the 7-day optimal users, the 14-day optimal users use the product more and have shorter period of inactivity at the end, as their trial length increases. Next, consider the 30-day optimal users: these users use the product less when they get 7 days (compared to the 7-day optimal users). Further, when these users get 30 days, they use it more and their dormancy period is shorter. However, when these users are given 14 days, we don't see any significant improvement in their usage or reduction in their dormancy period. Thus, these users really need the longer 30-day trial to register higher usage and have shorter periods of dormancy.

H.2 Heterogeneity in the Effect of usage on Subscription

So far, we have shown how the three segments differ in their usage behavior as a function of trial length; i.e., we have focused on the left side of Figure 3. Now we examine the heterogeneity in the effect of usage on subscription across the three segments, i.e., we focus on the right hand side of Figure 3. Table A12 presents the results from regressing the user’s subscription outcome on her optimal treatment/segment and her usage features, and the interactions of these two sets of variables. There are two key findings here. First, For 7- and 30-day optimal users, the effect of usage on subscription is similar. However, notice that for the 14-day optimal segment, the effect of usage on subscription is much smaller (the interaction between 14-days optimal segment and log usage is negative). This implies that while longer trials have a positive impact on usage for 14-day optimal segment, more usage doesn’t lead to higher conversions in this segment. Second, for 7-days optimal users, dormancy length has a large negative effect on subscription. For the 14-days optimal users, the negative effect of dormancy length is still there, but is somewhat smaller. For the 30 day people, the negative effect is the smallest.

	coef	std err	z	<i>P</i> > z	[0.025	0.975]
Intercept	-2.3522	0.140	-16.838	0.000	-2.626	-2.078
14-days trial	0.0301	0.023	1.297	0.195	-0.015	0.076
30-days trial	0.1780	0.025	7.205	0.000	0.130	0.226
14-days optimal segment	0.2007	0.063	3.182	0.001	0.077	0.324
30-days optimal segment	0.0237	0.094	0.252	0.801	-0.161	0.208
Number of active days	0.0497	0.003	16.984	0.000	0.044	0.055
Number of active days 14-days optimal segment	-0.0126	0.004	-3.030	0.002	-0.021	-0.004
Number of active days 30-days optimal segment	0.0001	0.007	0.019	0.985	-0.013	0.013
Log usage during trial	0.0712	0.007	9.723	0.000	0.057	0.086
Log usage during trial 14-days optimal segment	-0.0303	0.012	-2.572	0.010	-0.053	-0.007
Log usage during trial 30-days optimal segment	0.0016	0.019	0.086	0.931	-0.035	0.038
Dormancy length	-0.0317	0.001	-27.513	0.000	-0.034	-0.029
Dormancy length 14-days optimal segment	0.0039	0.002	2.590	0.010	0.001	0.007
Dormancy length 30-days optimal segment	0.0102	0.002	4.344	0.000	0.006	0.015
Indicator for software use	-0.5910	0.047	-12.463	0.000	-0.684	-0.498
Indicator for software use 14-days optimal segment	0.2351	0.076	3.075	0.002	0.085	0.385
Indicator for software use 30-days optimal segment	0.0076	0.118	0.064	0.949	-0.224	0.240
Total downloaded products	0.6077	0.017	35.734	0.000	0.574	0.641
Total downloaded products 14-days optimal segment	-0.1124	0.027	-4.098	0.000	-0.166	-0.059
Total downloaded products 30-days optimal segment	-0.0640	0.045	-1.422	0.155	-0.152	0.024

Table A12: Regressing subscription on usage features interacted with optimal trial length. We also control for pre-treatment variables. The number of observations is 303,514, and the pseudo R-squared is 0.292.

H.3 Pre- and Post-Treatment Attributes of the Segments

We now show the distributions of the pre-treatment demographics and post-treatment subscription outcomes for the three segments in Tables A13 and A14.

Variable	Value	Population	Policy Assigned Treatment		
			30 Days	14 Days	7 Days
Country	<i>United States</i>	54.9%	55.9%	45.8%	57.9%
	<i>Germany</i>	8.9%	8.2%	16.0%	6.6%
	<i>Japan</i>	7.9%	9.7%	11.5%	6.4%
	<i>Other</i>	28.3%	26.2%	26.8%	29.1%
	<i>Total</i>	100%	100%	100%	100%
Operating System	<i>Windows 10</i>	29.0%	9.0%	11.9%	37.0%
	<i>Windows 7</i>	21.5%	46.8%	21.7%	18.5%
	<i>Windows 8.1</i>	14.1%	0.9%	22.6%	12.7%
	<i>El Capitan</i>	13.9%	30.8%	18.0%	10.6%
	<i>Yosemite</i>	13.4%	3.3%	22.0%	11.7%
	<i>Other</i>	8.2%	9.2%	4.0%	9.4%
	<i>Total</i>	100%	100%	100%	100%
Skill	<i>Beginner</i>	68.9%	33.8%	41.1%	82.3%
	<i>Experienced</i>	12.8%	41.6%	22.0%	6.4%
	<i>Mixed</i>	10.7%	2.0%	35.0%	3.6%
	<i>Unknown</i>	7.5%	21.8%	1.9%	7.7%
	<i>Intermediate</i>	0.1%	0.7%	0.0%	0.0%
	<i>Total</i>	100%	100%	100%	100%
Job	<i>Student</i>	28.1%	12.7%	13.1%	34.9%
	<i>Unknown</i>	22.0%	69.8%	16.2%	18.5%
	<i>Hobbyist</i>	20.0%	7.3%	20.0%	21.5%
	<i>Other</i>	29.8%	10.3%	50.7%	25.1%
	<i>Total</i>	100%	100%	100%	100%
Signup Channel	<i>Website</i>	81.6%	65.2%	76.1%	85.3%
	<i>App Manager</i>	8.2%	17.7%	5.3%	8.1%
	<i>Other</i>	10.2%	17.1%	18.7%	6.6%
	<i>Total</i>	100%	100%	100%	100%

Table A13: Distribution of users' pre-treatment attributes for the three segments: those assigned to the 7-day condition, those assigned to the 14-day condition, and those assigned to the 30-day condition. (For each categorical variable, we show the fractions only for the top few categories in the interest of space.)

Variable	Population	Policy Assigned Treatment		
		30 Days	14 Days	7 Days
Mean				
Subscription rate	14.8%	17.0%	29.1%	9.8%
Revenue	\$536	\$545	\$546	\$524
Retention (months)	16.1	17.3	16.0	16.1
Fraction				
Purchased Bundle				
<i>Bundle I</i>	55.9%	57.5%	55.4%	55.9%
<i>All inclusive</i>	21.5%	20.3%	21.6%	21.7%
<i>Single Product</i>	19.2%	16.7%	19.2%	19.8%
<i>Other</i>	3.4%	5.5%	3.8%	2.6%
<i>Total</i>	100%	100%	100%	100%
Subscription Type				
<i>Commercial</i>	79.1%	79.0%	81.0%	77.3%
<i>Education</i>	20.7%	20.6%	18.9%	22.7%
<i>Other</i>	0.1%	0.4%	0.1%	0.1%
<i>Total</i>	100%	100%	100%	100%

Table A14: Means and distributions of users' post-treatment attributes for the three segments: users assigned to the 7-day condition, users assigned to the 14-day condition, and users assigned to the 30-day condition.

H.4 Additional Evidence to Rule Out Demand Cannibalization

We now provide additional evidence to rule out the demand cannibalization hypothesis. In Figure A3, we show that beginners who use the product more when assigned to the 14- and 30-day condition are more likely to subscribe when they use the product more. This rules out the possibility that these users are less likely to subscribe when assigned to the 14 and 30-day condition because they have already obtained their use for the product. Similarly, Table A15 shows that beginners who use the product more are more likely to subscribe, even after controlling for all the other user-specific variables.

	coef	std err	z	$P > z $	[0.025	0.975]
Indicator for using the software	-0.5764	0.035	-16.432	0.000	-0.645	-0.508
Total downloaded packages	0.5714	0.012	46.194	0.000	0.547	0.596
Number of active days	0.0477	0.002	21.734	0.000	0.043	0.052
Log usage during trial	0.0728	0.005	13.688	0.000	0.062	0.083
Dormancy length	-0.0322	0.001	-33.883	0.000	-0.034	-0.030

Table A15: Regression of subscription on usage features and trial length, with all the pre-treatment variables included as controls (not shown in the table above) for beginner users.

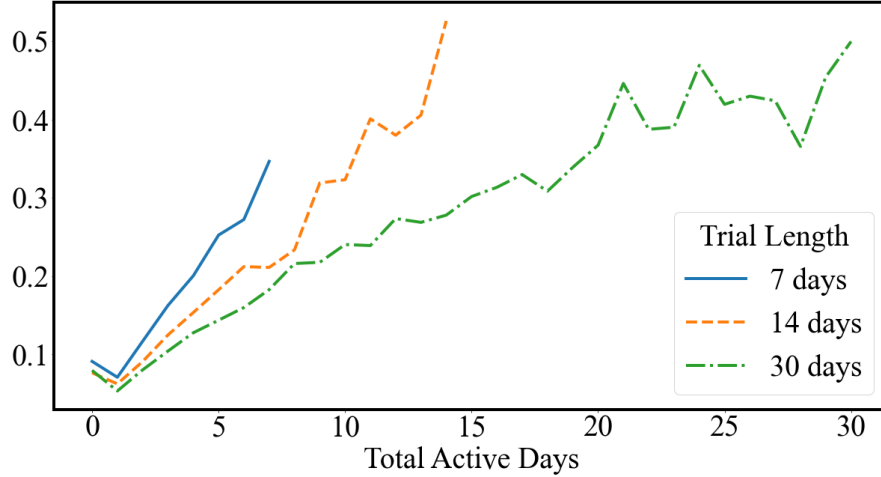


Figure A3: Subscription probability of *beginner* users with different total active days. The subscription probability increases as the total active days increases.

I Appendix for 7.2

Trial Length	Subscription Length						
	Mean	Std	Min	25%	50%	75%	Max
7 Days	15.13	9.31	0	8	16	22	73
14 Days	15.04	9.11	0	8	16	22	67
30 Days	14.81	9.05	0	8	16	22	108

Table A16: The summary statistics of the subscribed users' total months of subscription.

Trial Length	Subscribed Bundle					Subscription Type		
	1	All Inclusive	Single Product	4	5	Commercial	Education	Government
7 Days	55.24	22.02	19.44	1.68	1.62	78.57	21.30	0.13
14 Days	55.99	21.76	19.12	1.79	1.34	79.62	20.25	0.13
30 Days	55.98	21.64	18.95	1.82	1.62	79.02	20.85	0.13

Table A17: The fraction of each subscription bundle and type. We do not reveal the names of some of the bundles to preserve's the firm's anonymity.

J Appendix for 7.3

Personalized policy based on	7 Days	14 Days	30 Days	Total
Subscription	68.87	23.28	7.85	100
Subscription Length	84.17	15.25	0.58	100
Revenue	63.08	33.50	3.42	100

Table A18: The percentage of users allocated to each trial length in the three policies based on: (1) subscription, (2) subscription length, and (3) revenue.

Table A18 presents the proportion of users assigned to each trial length under the three different policies (optimized on the three different outcome variables). We see two interesting patterns here. First, when we personalize the policy to optimize subscription length, the policy has a tendency to assign shorter free trials more often. This is because users who get shorter trials are likely to subscribe sooner. The average number of days to subscription, from the start of the free-trial, is 121, 129, and 144 days for users who receive 7-, 14-, and 30- days trials, respectively. Further, we see that shorter trials lead to higher same-day subscriptions. 2.5% of users who received the 7-days trial and subscribed in the first day, whereas this number for the 14- and 30-days trials is 2% and 1.9%, respectively. Therefore, the subscription length is higher for users who received shorter trial lengths (see Table A16 in Appendix I). Hence, when a policy that optimizes subscription length will emphasize shorter trials. Next, we see that the policy designed to optimize revenues allocates a significantly larger proportion of users to the 14-days trial. This is because when we give 14 day trials, a slightly larger fraction of users subscribe to commercial licenses (Table A17 in Appendix I). Commercial licenses are significantly more expensive; so a revenue-optimizing policy tends to assign 14 days to a larger fraction of users in order increase the likelihood of commercial subscriptions. In sum, we see that the proportion of users assigned to different trial lengths can vary depending on the outcome being optimized.

References

- S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
- S. Athey, J. Tibshirani, S. Wager, et al. Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178, 2019.
- J. Bergstra, D. Yamins, and D. D. Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. *Journal of Machine Learning Research*, 2013.
- J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011.
- L. Breiman. Arcing Classifier. *The Annals of Statistics*, 26(3):801–849, 1998.
- L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen. *Classification and Regression Trees*. CRC press, 1984.
- M. Bruhn and D. McKenzie. In Pursuit of Balance: Randomization in Practice in Development Field Experiments. *American Economic Journal: Applied Economics*, 1(4):200–232, 2009.
- T. Chen and C. Guestrin. Xgboost: A Scalable Tree Boosting System. In *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794. ACM, 2016.
- J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.
- J. H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- B. Hansotia and B. Rukstales. Incremental value modeling. *Journal of Interactive Marketing*, 16(3):35–46, 2002.
- T. K. Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- K. Murphy. *Machine learning, a probabilistic perspective*, 2012.
- D. C. Mutz, R. Pemantle, and P. Pham. The perils of balance testing in experimental design: Messy analyses of clean data. *The American Statistician*, 73(1):32–42, 2019.
- X. Nie and S. Wager. Quasi-oracle estimation of heterogeneous treatment effects, 2017.
- O. Rafeian. Optimizing user engagement through adaptive ad sequencing. Technical report, Working paper, 2019.
- O. Rafeian and H. Yoganarasimhan. How does variety of previous ads influence consumer’s ad response? 2020.
- P. M. Robinson. Root-n-consistent semiparametric regression. *Econometrica: Journal of the Econometric Society*, pages 931–954, 1988.
- D. B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.
- P. Rzepakowski and S. Jaroszewicz. Decision trees for uplift modeling with single and multiple treatments. *Knowledge and Information Systems*, 32(2):303–327, 2012.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 0(0):1–15, 2018. doi: 10.1080/01621459.2017.1319839.