

Personalized Versioning: Product Strategies

Constructed from Experiments on Pandora *

Ali Goli

David H. Reiley

Hongkai Zhang

University of Washington

SiriusXM Pandora

SiriusXM Pandora

Abstract

The role of advertising as an “implicit price” has long been recognized by economists and marketers. However, the impact of personalizing implicit prices on firm profits and consumer welfare has not been studied. We first conduct a set of large-scale field experiments on Pandora by exogenously shifting the number of ads played per hour, i.e., the “ad load”, for over seven million users over a period of 18 months. We first show that while it takes a long time (more than a year) for the effect of ad load on consumption to stabilize, the treatment effect on subscriptions reaches steady-state much faster (less than six months). We then use a state-of-the-art machine learning model to examine the heterogeneous treatment effects of firm’s interventions on ad and subscription revenues. We next show that by reallocating ads across individuals, the firm can improve subscription profits by 7% without reducing total profits generated from advertising. To achieve the same subscription rate using a uniform ad-allocation policy, the firm would need to increase the number of ads served on the platform by more than 22%. Furthermore, the gains from personalization emerge quickly after implementation, as subscription behavior adapts to changing ad load relatively quickly. We also evaluate the welfare implications of personalizing implicit prices. Our results show that, on average, consumer welfare drops by 2% with the proposed personalization strategy, and the effect seems to be more pronounced for users that have a higher willingness to pay.

Keywords— Advertising, Personalization, Field Experiments, Heterogeneous Treatment Effects, Machine Learning

*Email addresses: agoli@uw.edu, dreiley@pandora.com, and hzhang@pandora.com. We thank JP Dubé, Pradeep Chintagunta, Günter Hitsch, Avi Goldfarb, Anit Rao, Hema Yoganarasimhan, Omid Rafieian, Simha Mummalaneni, and participants of the 2021 Quantitative Marketing and Economics and 2020 Marketing Science conference for their helpful comments and suggestions.

1 Introduction

The abundance of free online content creates a challenge for online-content providers to monetize their platforms. In the mid 1990s and early 2000s, to attract large audiences and generate advertising revenues, many firms offered online content for free (Edgecliffe-Johnson, 2009). As the industry matured, a number of content publishers experimented with subscription paywalls (Pérez-Peña and Arango, 2009). Although some firms, such as Netflix¹, have earned substantial profits through this strategy, the transition to a subscription-only model has not been especially easy for most firms. For instance, using data on a media publisher’s website visits, Chiou and Tucker (2013) show that instituting paywalls led to a 51% drop in online visits. The trade-off between viewership and subscription profits has led a number of firms, including Hulu, YouTube, Spotify, and Pandora, to adopt a hybrid approach, offering both an ad-supported free version and an ad-free subscription version. An interesting question is how those multiple versions should be designed and priced.

In the age of the Internet, digital products have become highly customizable, both in terms of content (e.g., Pandora’s personalized radio stations) and in terms of pricing. Although academic authors have discussed the returns to personalizing subscription prices (Shiller et al., 2013) or to engaging in fine-grained group-pricing strategies (Dubé and Misra, 2017; Smith, Seiler, and Aggarwal, 2022), these ideas are rarely implemented in industry (DellaVigna and Gentzkow, 2019; Hitsch, Hortacsu, and Lin, 2019; Bhatia, Moshary, and Tuchman, 2021). Firms such as Amazon and Staples have faced public backlash for experimenting with charging different prices to different customers (CNET, 2002; Valentino-DeVries, Singer-Vine, and Soltani, 2012).

Fear of customer backlash has led many firms to instead adopt “versioning” strategies (Shapiro, Carl, Varian, et al., 1998), where the seller offers each customer a menu of different product options, for example, ad-supported and paid subscriptions, and allows customers to self-select into choosing one of them. These versions can further be customized or made available based on what the publisher knows about its customers, such as when publishers release the hard copy when a new book is released and delay the release of paperback versions (Clerides, 2002). Versioning strategies have been used to discriminate in time or geographies by customizing the menu offered to the users, given the granular data that firms have on individuals this practice can also be personalized. The public appears to have a much more positive view of personalized product content or quality than of personalized prices

¹Netflix is also considering to offer an ad-supported version (Krouse and Vranica, 2022).

because many people feel that one group of consumers getting different explicit prices for the same product or service is unfair. Versioning and personalized pricing are two distinct strategies for the more general problem of price discrimination.

In 2015, the White House's Council of Economic Advisors report (CEA, 2015) noted that it is unclear which of these two strategies will become more prevalent in the era of Big Data:

“It is difficult to predict how big data will influence the prevalence of versioning. If it becomes easier to predict individual customers' willingness to pay and charge different prices for an identical product, versioning may be replaced by personalized pricing. On the other hand, versioning has the benefit of reducing concerns about inequity that arise with personalized pricing, and big data may facilitate versioning strategies based on “mass customization,” particularly for information goods that can be customized at relatively little incremental cost.”

Our goal is to demonstrate the role of advertising as an instrument for combining the two strategies into an idea that might be called “personalized versioning”: consumers choose between two versions of a product offering, one of which has quality personalization based on consumer characteristics. In this paper, we report the results and analysis of a field experiment conducted on Pandora during 2016-2017 that shifted audio ad load for over seven million users. During this period, Pandora offered two products: (i) the Pandora Plus subscription product, and (ii) the ad-supported product which was offered free of charge. Both of these products were *non-interactive radio* products, meaning the listener could not listen to an audio track on demand but could create stations based on a favorite artist or track, and personalize her stations by thumbing songs up and down. Both products used the same music catalog and user interface; the main difference between the two was that Pandora Plus was ad-free, whereas ad-supported Pandora listeners would encounter ads between tracks, when switching between stations, or when skipping tracks.² To manage the trade-off between ad and subscription revenues between these two versions of the product, Pandora has two levers: (a) setting subscription prices (explicit price) and (b) changing the number of ads served to ad-supported listeners (implicit price). Particularly, Pandora can change the number of ads served to each listener by personalizing the frequency of scheduled commercial interruptions (or “ad pods”) as well as the number of scheduled

²Pandora Plus offers some additional features including offline listening, the ability to replay songs, and higher-quality audio.

ads per pod.

Increasing the frequency or length of ad pods increases the opportunities to serve an ad, which we refer to as “ad capacity.” By contrast, the listener’s realized “ad load,” or the actual number of ads served per hour, also depends on the consumption level and advertisers’ demand for each listener. Our goal in this paper is to study the gains that arise from personalizing the ad allocation policy at the individual level. To achieve this goal, we need to understand how consumption and subscription decisions vary as a function of firm’s ad allocation strategy, which requires us to overcome the following challenges:

- (a) **Consumption endogeneity:** Even holding fixed the ad-serving strategy, the realized ad load experienced by different listeners correlates with their consumption behavior. For instance, the longer listeners stay on the platform, the harder it is for Pandora to fill their full ad capacity.³ In general, listeners’ product choice and consumption level decisions can both impact and be impacted by the realized ad load. These correlations create an endogeneity problem, which motivates our randomized experiment.
- (b) **Partial control over realized ad load:** The realized ad load not only depends on the firm’s ad-scheduling policy but also on listener behavior (discussed above) and on advertiser demand (some listeners are in higher demand than others). This partial control causes one listener to receive more ads than another even when both receive the same policy from the firm, and these heterogeneous differences need to be taken into account in the firm’s optimization problem.

We exploit a set of large-scale field experiments that exogenously shift the ad-pod frequency and length for over seven million Pandora listeners during 2016-2017. We then use a state-of-the-art machine learning model to learn the heterogeneous treatment effects of the firm’s interventions on the realized profits from ads and subscriptions. To achieve this we use a set of models to understand heterogeneity in realized ad load, consumers’ extensive-margin decisions (switching between outside option, plus, and ad-supported options), and consumers’ intensive-margin decisions (number of ad-supported hours consumed). To build and estimate these models, we combine insights from structural estimation with those from neural networks.⁴ We use split neural networks (Kim et al., 2017) to

³Given listener-level frequency caps standard in these ad campaigns, the longer listeners stay on the platform, the more difficult it is to fill their full ad capacity.

⁴For a few use cases of neural networks for estimating structural models, see Wei and Jiang (2020), and Kaji, Manresa, and Pouliot (2020).

impose exclusion restrictions that enable the model to better learn heterogeneous treatment effects. Our architecture is similar to Shalit, Johansson, and Sontag (2017), who use neural networks to predict individual-level outcomes across different counterfactuals. Shalit, Johansson, and Sontag (2017) and Farrell, Liang, and Misra (2021) show neural networks are effective in learning treatment heterogeneity and achieve comparable performance to direct methods for learning heterogeneous treatment effects such as causal forests and treatment-effect projection (Hitsch and Misra, 2018; Wager and Athey, 2018).

Subsequently, we solve the firm’s optimization problem and evaluate the impact of the prescribed policy using an inverse probability-weighted (IPW) estimator; see Horvitz and Thompson (1952) for its use and origins in statistics and Hitsch and Misra (2018), Simester, Timoshenko, and Zoumpoulis (2020), Rafeian and Yoganarasimhan (2021), and Yoganarasimhan, Barzegary, and Pani (2022) for a few recent examples of IPW estimators in marketing. Our results demonstrate that holding fixed the total number of ads served, the firm can improve subscription profits by 7% without any loss in total ad revenue. To achieve the same subscription rate with a uniform allocation strategy, the firm would have to serve 22% more ads, which would have a negative impact on hours listened. We then study the impact of the proposed policy on consumer welfare and show that, on average, consumer welfare declines by 2%. To the best of our knowledge, this study is the first to use a field experiment to evaluate the returns to personalizing product quality. Our results inform policymakers and firms regarding the implications and returns to personalizing price/quality of product offerings.

2 Literature Review and Contributions

In analyzing the potential that arises from “personalized versioning,” our findings contribute to four strands of academic research. First, we contribute to the literature that measures the returns to personalization. Researchers have studied personalization in a wide variety of contexts, examples include prices (Rossi, McCulloch, and Allenby, 1996; Shiller et al., 2013; Dubé and Misra, 2017), e-mail content (Ansari and Mela, 2003; Sahni, Wheeler, and Chintagunta, 2018), website design (Hauser et al., 2009), search rankings (Yoganarasimhan, 2020), promotions (Hitsch and Misra, 2018; Yoganarasimhan, Barzegary, and Pani, 2022), mobile advertising (Rafeian and Yoganarasimhan, 2021), and ad sequencing (Rafeian, 2019). While the research that involves personalizing advertisement measures returns to changing ad content or targeting ads to improve ad effectiveness or revenue from

ads, our work focuses on personalizing ad load as an instrument for improving subscription revenue.

Second, we add to the literature that models product quality as an endogenous decision. In a single-product setting, Spence (1975) shows that a monopolist may offer a higher or lower quality level than the social optimum. In a multi-product regime, Mussa and Rosen (1978) and Maskin and Riley (1984) demonstrate that to attract high-type customers, the monopolist has the incentive to degrade the quality of lower-end products, which creates a negative externality on customers with lower quality valuation. This finding relates to the “damaged goods” literature, where a firm has the incentive to “damage” a developed product to build a lower-quality version (Deneckere and Preston McAfee, 1996). One approach for implementing versioning is by bundling a good with a “bad” like waiting time, advertising, or search cost (Salop, 1977; Chiang and Spatt, 1982). McManus (2007), Clerides (2002), and Verboven (2002) provide evidence of versioning in specialty coffee, book publishing, and European auto industries, respectively. Crawford and Shum (2007) measure the extent of quality degradation in cable-television subscription bundles, and Crawford, Shcherbakov, and Shum (2015) study the welfare effects of endogenous quality choice. We show that personalizing product quality helps limit such distortions. In particular, we find that our proposed optimal personalization of ad load produces subscription benefits equivalent to a uniform increase in ad load (a degradation in quality) of about 22%. Previous research in marketing has shown that service-quality variation over time can improve profits by increasing customer retention, for some consumers, a phenomenon potentially explained by risk aversion in the consumer learning process (Sriram, Chintagunta, and Manchanda, 2015). Our findings show that another kind of variation in the quality of service – across consumers – can improve profits in a product line by inducing users to upgrade to higher-end products. The substitution between products offered in a product line, along with switching costs between products, presents yet another opportunity for firms to leverage changes in quality of service (ad load) as a screening mechanism.

Third, our results contribute to the literature that considers product-line strategy in offering free (ad-supported) and paid versions of information goods (Shapiro, Carl, Varian, et al., 1998). On the theoretical side, this literature extends the versioning framework in Mussa and Rosen (1978) for information goods that rely on both advertising and subscriptions as sources of revenue. Tåg (2009) shows that introducing an ad-free subscription decreases consumer welfare because the firm has the incentive to increase advertising in the ad-supported version to earn more profits from the paid

product. Researchers have studied the role of dynamics, consumer heterogeneity, competition, quality learning, and advertiser heterogeneity on the revenue model adopted by firms (Prasad, Mahajan, and Bronnenberg, 2003; Godes, Ofek, and Sarvary, 2009; Kumar and Sethi, 2009; Halbheer et al., 2014; Sato, 2019; Lin, 2020).

On the empirical side, Chiou and Tucker (2013) show introducing paywalls can dramatically reduce viewership. Lambrecht and Misra (2017) present evidence for counter-cyclical quality improvements to ESPN's free service. The authors argue consumers are heterogeneous in their valuation of the content, which may vary over time. This heterogeneity rationalizes a quality-discrimination mechanism along the time dimension. In this paper, we establish the trade-offs between ad and subscription revenue and then show that by personalizing the ad schedule (quality of service), the firm could improve subscription profits. Although the idea of using ads as a screening mechanism in freemium products is not new (Tåg, 2009; Sato, 2019), we are not aware of any paper that has empirically investigated the personalization of product quality, especially in the advertising context.

Another strand of empirical work in this area has established the negative impact of ads on media consumption (Becker and Murphy, 1993; Gentzkow, 2007; Goldstein et al., 2014). Wilbur (2008) uses observational data to document the negative impact of TV advertising on viewership. Huang, Reiley, and Riabov (2018) examine a previous field experiment in Pandora and document that it takes a long time to reach the steady state to measure the treatment effect of ad load on listenership. In this paper, however, our focus is to understand the trade-offs between ad and subscription revenues. We show that the treatment effect of ad load on subscriptions stabilizes much faster than the effect on listenership. This means that persistent short-run changes in ad load can induce users to substitute to the paid version. This phenomenon, coupled with switching costs between plus and ad-supported products, presents an opportunity for the firm to improve subscription profits through personalized temporary changes in implicit prices. Hence, we solve an optimization problem that enables the firm to reallocate the ad load across listeners and jointly maximize the profit from both subscriptions and ad revenue.

Finally, our findings are also relevant to the price-discrimination literature. Although the amount of advertising on ad-supported media is a quality of service measure, it can also be viewed as an *implicit price* that is charged in units of time rather than money. To the best of our knowledge, the returns to personalizing this implicit price and its welfare implications have not been studied before.

Theoretically, third-degree price discrimination could improve social welfare (Varian, 1985) or could even improve consumer surplus by expanding output (Cowan, 2012). Dubé and Misra (2017) examine the returns to an extreme form of third-degree price discrimination using a large-scale field experiment at Ziprecruiter. They show that while firm profits improve by about 10%, consumer surplus falls less than 1%. One of the main differences between our problem and a classical price-discrimination problem is the fact that the consumer has the option to pay both with time and money. Therefore, listeners are screened based on both their willingness to pay and their marginal value of time.⁵ This means that the correlation between willingness to pay in time and money units could influence the effectiveness of our personalization algorithm. For instance, income and the marginal value of time could be positively correlated (Aguiar, Hurst, and Karabarbounis, 2011; Aguiar, Hurst, and Karabarbounis, 2013). Furthermore, income is likely to be negatively correlated with price sensitivity. On the one hand, the algorithm has the incentive to move more ads toward higher-income individuals because they are less price sensitive and more likely to upgrade to the paid subscription. On the other hand, higher-income individuals may place a larger value on their time and are also more likely to churn when faced with more ads. Because of these trade-offs, the ad-allocation mechanism and its welfare implications are a priori ambiguous.

The rest of this paper is organized as follows. We first introduce a conceptual model to discuss the personalized versioning idea and illustrate trade-offs between ad and subscription revenues. We then discuss the field experiments conducted at Pandora Media and present reduced-form evidence to illustrate the impact of changing ad load on listeners' choices. Subsequently, we use a state-of-the-art machine learning model to learn the heterogeneity in response to changes in ad load among listeners. Which are then leveraged to reallocate ads to improve firm profits. Finally, we discuss the prescribed policy and its welfare implications.

3 Conceptual model

In this section, we present a toy model that illustrates how listeners choose between the outside options, consuming the ad-supported product, or using the paid subscription. We first consider the personalized versioning problem for a single listener and illustrate the trade-offs in optimizing ad load.

⁵The idea of using differences in valuation of time for optimizing menu offerings and its welfare implications has been discussed in Salop (1977), and Chiang and Spatt (1982). However, we are not aware of any empirical work that has investigated a personalized policy that leverages this heterogeneity.

Then we discuss the platform-level optimization problem where ad load for all listeners is optimized simultaneously.

3.1 Single-user optimization

Consider the following random utility model:⁶

$$u(z, p; \theta, \beta, \alpha) = \max \begin{cases} \epsilon_0 & \text{outside option,} \\ \theta - \alpha z + \epsilon_a & \text{ad-supported,} \\ \theta - \beta p + \epsilon_p & \text{paid subscription,} \end{cases}$$

where θ is the utility from consuming the product, z specifies the ad load⁷, and p is the subscription price. The ad-supported product is bundled with a “bad”, that is advertising, and users’ disutility per unit of advertising and payment are measured by α , and β , respectively. The parameters α and β reflect how time and money are valued by users, that is users with higher/lower values of α and β are more/less sensitive to ads and prices, respectively. Also, let γ and c be the revenue per ad and the marginal cost of offering the service, respectively. Since we are focusing on a single listener problem, we assume that the ad load z can change without impacting the price of impressions γ . Finally, ϵ_0 , ϵ_a , and ϵ_p are independent random variables that follow a type-I extreme value distribution.

Those who tend to have a higher willingness to pay in money terms are likely more sensitive when paying in time units. For instance, higher-income individuals tend to have lower price elasticity but higher marginal value of time (Aguiar, Hurst, and Karabarbounis, 2011; Aguiar, Hurst, and Karabarbounis, 2013). This confound can generate a negative correlation between price (β) and time sensitivity (α) in this setup. On the one hand, the monopolist has the incentive to serve fewer ads to more ad-sensitive users, that is larger α . On the other hand, the same users likely have a higher willingness to pay, that is smaller β , and are more likely to upgrade to the subscription service if they face higher prices in time terms.

Note that if the seller were to offer only the ad-supported version, customers with higher ad

⁶We refer the interested readers to the online appendix for a use case of versioning in screening heterogeneous users even in the absence of a random utility model.

⁷In this simplified model, we assume users can consume the service in exchange for listening to z ads. In our empirical exercise, we account for the fact that the intensive margin of consumption (number of hours) could vary across users and that possibility factors into the ad revenue. We also account for the fact that advertisers’ demand could vary across user segments.

sensitivity would receive fewer ads. However, the correlation structure between ad sensitivity (α) and price sensitivity (β) can lead to both higher or lower frequency of ads for more ad-sensitivity users. Let us consider the problem of personalizing the ad load z given a fixed price p for the subscription service. Let us assume the marginal cost of offering service is c , and revenue from serving each ad is γ . The problem the service provider faces is to

$$\underset{z}{\text{maximize}} \quad \underbrace{\frac{\overbrace{P_a(\theta, \alpha, \beta, z, p)}^{e^{\theta-\alpha z}}}{1 + e^{\theta-\alpha z} + e^{\theta-\beta p}}(\gamma z - c)}_{\text{expected profits from ads}} + \underbrace{\frac{\overbrace{P_s(\theta, \alpha, \beta, z, p)}^{e^{\theta-\beta p}}}{1 + e^{\theta-\alpha z} + e^{\theta-\beta p}}(p - c)}_{\text{expected profits from subscription}}. \quad (1)$$

The optimal ad load is determined by equating the marginal effect of increasing ad load on ads and subscriptions. As discussed above, on the one hand, a higher ad load leads to more profits from subscriptions and increases revenue conditional on being an ad-supported member; on the other hand, it lowers profits from the ad-supported users by increasing churn. Furthermore, the correlation structure between ad and price elasticity can lead to either higher or lower ad load for users with higher ad elasticity. To illustrate this trade-off, let us hold the price fixed and optimize the ad load for a set of given parameters $(\theta, \alpha, \beta, p, \gamma, c)$ while enforcing different correlation structures between α and β . Let us assume $\theta = 4$, $\beta = 2 - (0.1)\alpha$, $\gamma = 0.5$, $c = 1$, and $p = 5$, and let α vary between 1 and 2. The optimal ad load (z) as a function of ad sensitivity is strictly decreasing and is plotted in panel (a) of Figure 1. However, if the correlation structure between price and ad sensitivity is stronger, say, $\beta = 2 - (0.5)\alpha$, the optimal ad load could be a non-monotonic function of ad sensitivity as depicted in panel (b) of Figure 1. This example illustrates that in our multi-product setting, forming ex-ante predictions on which customer segments, for example, high/low income, bear the cost of personalizing the ad load is difficult.

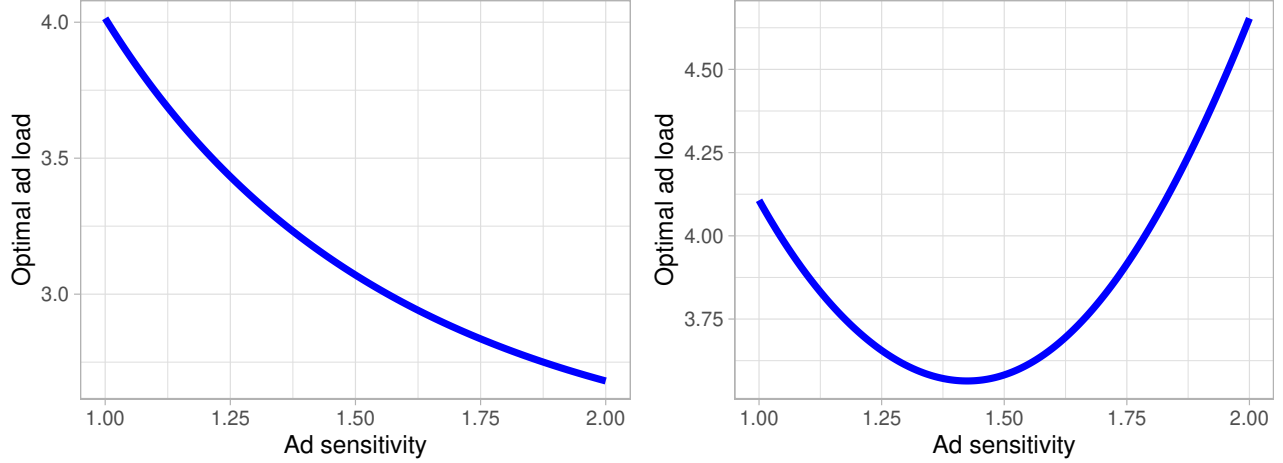


Figure 1: Optimal ad load as a function of ad sensitivity. Left panel: the price sensitivity $\beta = 2 - (0.1)\alpha$, the optimal ad load is a decreasing function of ad sensitivity because gains from subscriptions do not outweigh the losses. Right panel: the price sensitivity $\beta = 2 - (0.5)\alpha$, and α varies between 1 and 2. In this case, optimal ad load is a non-monotonic function of ad sensitivity, because users with higher ad sensitivity are less price sensitive, and uplift from subscriptions outweighs the losses.

3.2 Platform-level optimization

In section 3.1 we illustrated the trade-offs in optimizing the quality of service in a single-user setting. However, our objective is to implement personalized versioning for all listeners on the platform. In this scenario, the assumption that shifting z does not impact the price of impressions is not realistic. In other words, we need a supply-side model to measure how ad inventory size changes as a function of impressions' price. Our field experiments only shift ad load for a small portion of listeners without a significant impact on the overall ad inventory level. This means that we do not have exogenous variation in the price of impressions and cannot model the relationship between the total ad inventory size and the price of impressions. Therefore, we consider a more nuanced problem, that is can the firm improve profits by personalizing ad load while holding the total ad inventory fixed? In this scenario, the platform-level counterpart of problem (1) is:

$$\begin{aligned}
\underset{z_i}{\text{maximize}} \quad & \sum_i (\gamma z_i - c) \cdot P_a(\theta_i, \alpha_i, \beta_i, z_i, p) + (p - c) \cdot P_s(\theta_i, \alpha_i, \beta_i, z_i, p) \\
& \sum_i P_a(\theta_i, \alpha_i, \beta_i, z_i, p) \cdot z_i = \Gamma,
\end{aligned} \tag{2}$$

where i index users with heterogeneous preferences $(\theta_i, \alpha_i, \beta_i)$ with a utility function similar to (1). $P_s(\cdot)$ and $P_a(\cdot)$ are defined as in (1) and are the probability of being a subscriber, or an ad-supported listener, respectively. Γ is the total ad inventory size, and the rest of the parameters are defined similar to (1). The constraint on the total inventory size ties the ad load decisions across all users to each other. This constraint means that, unlike the single-user case, the marginal effects of changing ad load on profit from ads and subscriptions are not necessarily equal. In the subsequent sections, we discuss how to solve this problem while overcoming challenges such as the dependence of realized ad load on advertisers' preferences for different demographics, and the impact of intensive margin adjustments (hours consumed) on the number of ads shown to each user.

To summarize, our discussion in this section highlights that: (a) in our multi-product setting, a priori it is unclear which customer segments bear the cost of personalization (b) in absence of exogenous variation in the price of impressions, the proper counterfactual to consider is to personalize quality by re-allocating ads across listeners while holding the total ad inventory fixed, and (c) other nuances such as advertisers' preferences and intensive margin adjustments (hours consumed) play an important role in the platform-level optimization problem which we address in subsequent sections.

4 Field experiments at Pandora Media

Now that we have built a conceptual model to understand the trade-offs, we delve into the field experiments used in this study. At the time of this experiment, Pandora offered two tiers of products: (a) ad-supported and (b) plus. The ad-supported and plus versions are both “radio”⁸ products and have the same music catalog and user interface. Whereas the plus subscription is ad-free with a monthly subscription fee of \$4.99, ad-supported listeners are exposed to video/audio ads in exchange for using the service.

⁸The radio products offer quasi-audio-on-demand services as they personalize the radio stations to cater to listener preferences using the feedback (thumbs up/down, and skips) provided by the listeners. In the second quarter of 2017, Pandora started offering the premium service, which was an ad-free audio-on-demand product.

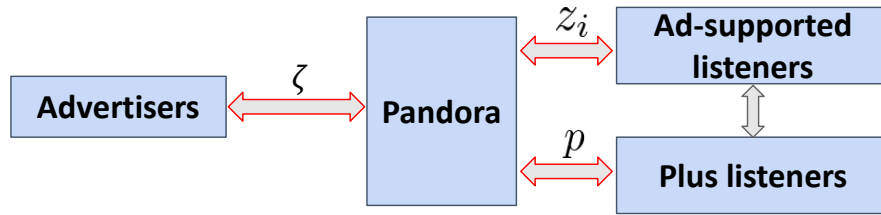


Figure 2: An illustration of variables that affect the revenues from ads and subscriptions. The overall ad inventory depends on the price of impressions (γ) and the ad-supported audience. The mix of ad-supported users is determined by the ad allocation policy that could be personalized (z_i). Finally, the price of the subscription plan (p), which is uniform across users can affect the substitution between the ad-supported and plus services.

Pandora has three main levers that impact listeners' consumption and its revenue: (a) price of impressions, (b) price of the subscription service, and (c) ad allocation strategy across users. Figure 2 illustrates these different levers. Our field experiments create exogenous variation in ad load by shifting (c) rather than affecting (a) or (b). This is facilitated by exogenously shifting the ad load for a small random subset of listeners, which means that the overall ad inventory or price of impressions is not impacted by the experiment. Unlike most digital ads that are sold via online auctions, audio ads are sold via forward contracts. Advertisers specify demographic targets, e.g., women 25-34 who live in New York. We do not observe the closing price of the individual advertising contracts and do not have exogenous variation in the price of impressions and the realized ad inventory. The subscription prices were also held fixed for all users in the course of the experiment.

The experiments shift the time spent consuming ads, by changing the number of audio ads played in each ad break (pod) and the frequency at which users become eligible to receive audio ads. Listeners become eligible for the first ad pod within the first five minutes of a listening session. The subsequent ad pods are delivered using a set of fixed timers, once an ad pod is delivered, the timer is reset. At the beginning of every track, the system checks to see if the user is eligible to receive an ad pod, that is if the timer is set. The length of an ad pod determines the number of ads (one, two, or three ads) that can be served in an ad break. The experiments start in July of 2016 and shift both *frequency* and *length* of ad pods using six experiment conditions and a control cell, which represents the default strategy employed by the firm.⁹ Figure 3 illustrates the ad delivery mechanism and the levers that

⁹Pandora conducted a previous ad-load experiment in 2014 and the findings from that experiment were reported in Huang, Reiley, and Riabov (2018). The objective of the previous field experiment was to measure the treatment effect of ad load on listenership. The goal of our study is to understand the treatment effect of ad load on the overall business (both subscriptions and ad revenue) and to investigate returns to personalization. To that end, we considered a larger range of ad load conditions in our experiment and we collected a diverse set of user features in the pre-treatment period which enable us to study the heterogeneous treatment effects of ad load on revenues from ads and subscriptions across users.

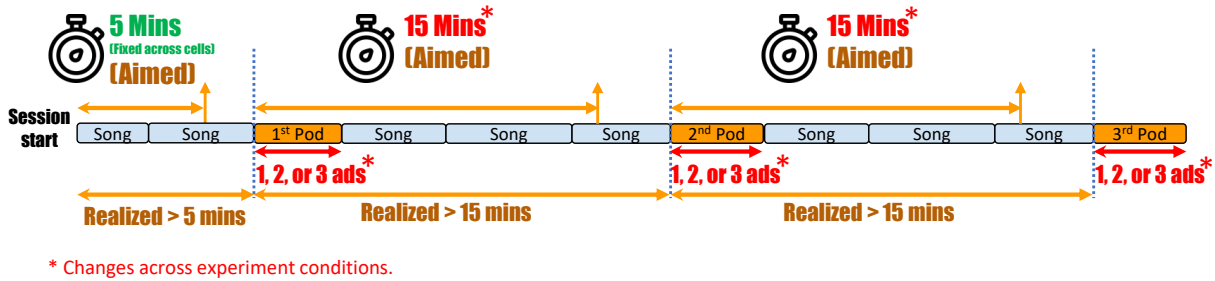


Figure 3: An illustration of the ad delivery mechanism. The listeners across all experiment conditions become eligible for the first ad pod after the first five minutes of each listening session. The experiment shifts the red components in the figure, which are the timers used after the first ad pod, and the number of ads served within each ad pod.

the experiment shifts. Random subsets of users were assigned to each treatment condition and the treatment persisted for about 18 months. The experiment cells are presented in Table 1 below:

Table 1: The experiment shifts the number of ads delivered in each interruption (pod length) and how often listeners become eligible for ad pods (pod frequency) in an hour. The size of each treatment cell is specified as the percentage of all listeners on the platform. Rows and columns correspond to pod frequency and size, respectively. In total, we have seven cells that include the six conditions shown below and the control condition which consists of 1% of the listeners.

		Audio ads per interruption		
		1	2	3
Audio ad interruptions per hour	3	1%		
	4		2%	0.5%
	5			0.5%
	6		0.5%	0.3%

The seven cells in our experiment consist of more than seven million listeners. From this point forward, we refer to each cell in our experiment as FxL, where F and L are *intended* pod frequency, and length, respectively. For instance, the 3x1 condition refers to a treatment where the timer in Figure 3 is set to 20 minutes. Therefore, the *intended* number of ad pods per hour in the 3x1 condition is three and each pod consists of one ad. The control condition is similar to the 4x2 condition, but the first ad pod within each listening session is constrained to have at most one ad. The control condition comprises 1% of the total listeners¹⁰ on Pandora. The treatment cells were determined based on the estimated effect sizes in the previous field experiments and in consultation with stakeholders. Note that as we move toward higher-capacity treatments such as the 6x3 condition the cell sizes shrink.

¹⁰ Assigning X% of total listeners to a certain condition means all existing listener ids in addition to those ids that may be created while the experiment is running will have an X% chance to be assigned to that condition and will be kept in that condition for the period of the study. We include all users in our analysis, rather than filtering them based on activity, as it yields the most conservative bound on returns to personalization and is a better representation of the effect of personalization on the entire platform. We have performed our analyses with the subset of users who had some activity in the 6 months period leading to the experiment or those who had logged in at least once within the first three months of the experiment and our results remain qualitatively similar.

Table 2: Comparing treatment and control groups across some of the pre-treatment features calculated during the first quarter of 2016. All features except for gender and zip code mean income are normalized such that the mean of control is equal to 100. Standard errors are reported in the parenthesis. The differences between treatment conditions and the control cell are statistically indistinguishable at $p < 0.05$ across all variables.

Variable	Experiment condition						Control
	3x1	4x2	6x2	4x3	5x3	6x3	
All hours	99.951 (0.351)	100.317 (0.249)	100.023 (0.497)	100.243 (0.499)	100.321 (0.499)	100.621 (0.644)	100 (0.351)
Ad supported hours	99.882 (0.385)	100.215 (0.273)	99.754 (0.542)	100.352 (0.549)	99.948 (0.542)	100.178 (0.704)	100 (0.385)
Thumbs	99.747 (0.584)	99.904 (0.413)	100.218 (0.833)	99.916 (0.826)	99.376 (0.821)	100.047 (1.062)	100 (0.586)
Thumbs up	99.888 (0.574)	100.112 (0.406)	100.317 (0.809)	100.007 (0.813)	99.426 (0.809)	99.959 (1.044)	100 (0.574)
Skipped tracks	99.863 (0.567)	100.104 (0.405)	100.341 (0.823)	100.247 (0.841)	99.547 (0.81)	100.883 (1.07)	100 (0.578)
Station changed	99.093 (1.04)	99.797 (0.74)	101.121 (1.596)	98.62 (1.387)	100.363 (1.545)	100.589 (1.838)	100 (1.053)
Age	99.97 (0.086)	99.964 (0.061)	99.97 (0.122)	99.967 (0.122)	100.05 (0.122)	100.014 (0.158)	100 (0.086)
Gender (Male = 1)	0.45 (0.001)	0.451 (0.001)	0.45 (0.002)	0.451 (0.002)	0.452 (0.002)	0.451 (0.002)	0.452 (0.001)
Zip code mean income	73,427.358 (66.191)	73,429.394 (46.87)	73,411.446 (93.814)	73,380.273 (93.76)	73,462.894 (93.454)	73,522.863 (120.934)	73,438.642 (66.062)

There are two reasons for this: (a) the treatment effects of high-capacity conditions are expected to be large and detectable even with small sample sizes, and (b) the cost of exposing listeners to high-capacity conditions such as 6x3 in terms of churn and listening hours is large and to mitigate this cost fewer users were assigned to high-capacity conditions.

We now illustrate that the randomization algorithm has achieved its goal and treatment assignment is not systematically correlated with any covariates of interest. We select the set of users across the six treatment conditions and control groups who were active in the first quarter of 2016 before the experiment and compare their age, gender, and some of the other key behavioral variables in the pre-treatment period in Table 2.¹¹ The differences between treatment cells and control across all variables are statistically indistinguishable from 0 at $P < 0.05$. Overall, this table shows the treatment and control groups are not systematically different, and confirms the treatment assignment has been random.

Note that the *ad capacity*, that is, the number of opportunities to show an ad per hour, in the FxL condition ends up being far less than F·L. For instance, a listener in the 6x3 condition ends up becoming eligible for an ad pod fewer than six times per hour, because the song endings do not

¹¹Due to our agreement with Pandora we cannot reveal the actual estimates for some of these features; therefore, we have normalized those features such that the sample average of the control group is equal to 100.

perfectly align with the timers. The experiments shift ad capacity by changing pod frequency and size, however, the realized number of ads shown to each listener (ad load) also depends on advertisers' demand. For instance, an increase in the ad capacity, that is the rate at which ads *can be shown* to a user, for a user that does not belong to an attractive demographic group for advertisers does not necessarily translate to an X% increase in the ad load, that is *the realized rate* of ads for each user.

Table 3 reports the average realized ad load, capacity, and fill rate during the first year of the experiment. Realized ad capacity is the number of opportunities that the ad delivery system determines a listener as being eligible to receive an ad, though not all these opportunities get filled as the system may not be able to fetch ads to serve users. The proportion of ad opportunities that were filled is referred to as the fill rate. As one would expect, the fill rate tends to fall as we move toward higher-capacity conditions; for example compare the 3x1 and 6x3 conditions.¹² Finally, note the realized ad load depends on both the realized ad capacity and the fill rate (advertisers' demand); therefore, an X% increase in realized ad capacity does not necessarily translate into an X% increase in realized ad load. This demonstrates the fact that firms need to account for the discrepancy between the *intended* and *realized* change in the implicit price, which leads to an additional layer of complexity relative to the traditional pricing problems. We refer the interested reader to online Appendix B for more discussion on the ad delivery mechanism and partial control over realized ad load/capacity.

Table 3: The realized ad load, capacity, and fill rate across experiment cells.

	Experiment condition						
	3x1	4x2	4x3	5x3	6x2	6x3	Control
Realized ad load	2.947 (0.006)	4.659 (0.007)	5.541 (0.008)	6.123 (0.011)	5.602 (0.008)	6.665 (0.023)	4.208 (0.008)
Realized ad capacity	3.512 (0.007)	6.326 (0.009)	8.289 (0.008)	9.35 (0.013)	7.789 (0.009)	10.347 (0.025)	5.56 (0.013)
Fill rate	0.853 (0)	0.738 (0)	0.676 (0.001)	0.665 (0.001)	0.723 (0.001)	0.657 (0.001)	0.763 (0)

¹²Note that users become eligible for the first ad pod within the first five minutes of a listening session, see Figure 3. Consequently, for users with a short listening session in the 3x1 condition, the realized ad capacity could be greater than 3.

5 Average treatment effects

In this section, we examine the impact of ad load on the different revenue sources and consumption patterns. To that end, we plot the change in the realized ad load, ad revenue, active users, subscription revenue, ad-supported hours, and all hours, that is, the sum of ad-supported and plus hours, across the highest and lowest ad-load arms relative to the control condition in Figure 4. The experimentation system ramps up/down ad load in a four-week period. As illustrated in the figures the ad load starts increasing in June 2016 and stabilizes by July 2016. The figures measure each outcome of interest as a percentage change relative to the control arm. For instance, in the 3x1 condition, all hours increase by about 2% relative to control by the end of the experiment. Higher ad load leads to higher revenues from ads and subscriptions, but it impacts both extensive and intensive margins of consumption which could affect profits.

To measure the average treatment effect of ad load on consumption and subscriptions, we normalize the outcomes to their average in control and scale them to measure differences relative to control:¹³

$$\tilde{\mathbf{Y}}_i = 100 \cdot \frac{\mathbf{Y}_i}{\frac{\sum_{j \in \mathcal{C}} \mathbf{Y}_j}{\mathcal{N}_{\mathcal{C}}}}, \quad (3)$$

where \mathcal{C} and $\mathcal{N}_{\mathcal{C}}$ are the set of users in the control condition, and the number of users in control, respectively. We then consider the following instrumental variable regression:

$$\tilde{\mathbf{Y}}_i = \alpha + \beta \cdot \mathcal{A}_i + \epsilon_i, \quad (4)$$

where i indexes listeners, and $\tilde{\mathbf{Y}}_i$ is a normalized outcome of interest as in (3), i.e., listening hours, activity dummy, and plus subscription dummy, for each user in a given week. \mathcal{A}_i is the average number of ads per hour delivered to listener i . We use the experiment condition dummies to instrument for \mathcal{A}_i and estimate the treatment effect of ad load on activity, listening hours, and plus subscriptions for all weeks in the 2016-2017 period.¹⁴ The estimates for β during each week along with the 95% confidence intervals are plotted in Figure 5. Note that the outcome $\tilde{\mathbf{Y}}_i$ is scaled by the average in

¹³Due to our agreement we cannot reveal the actual numbers for consumption or other metrics used here. Consequently, we normalize the outcomes relative to the average in control which means that we report the treatment effects as percentage changes relative to control rather than absolute differences, e.g., in dollars, counts, or hours.

¹⁴The estimates from reduced-form regressions, i.e., outcomes regressed directly on treatment dummies, are reported in online appendix C.

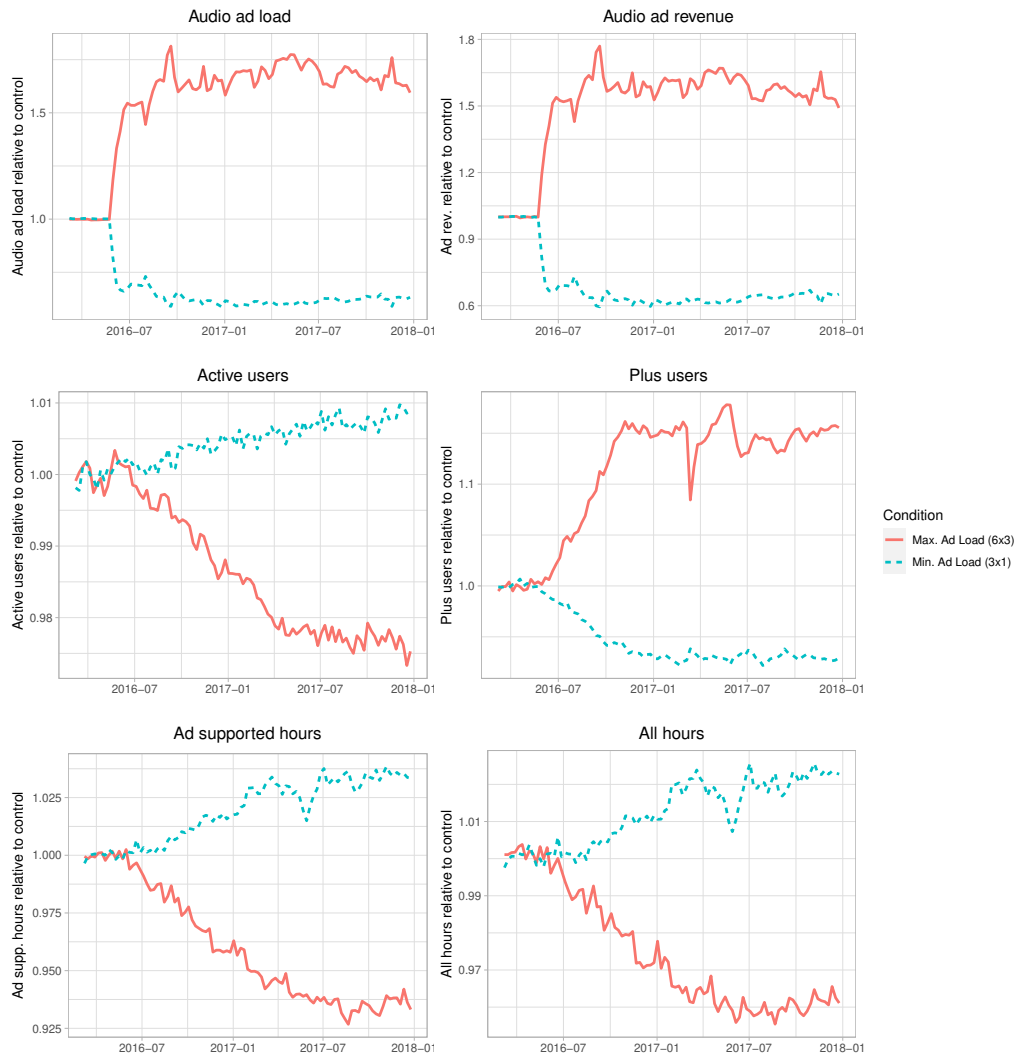


Figure 4: The effect of ad load on consumption, ad, and subscription revenues. The highest ad load condition (6x3) delivers about 50% more ads compared to the control. Even though the effect of ad load on listenership and activity grows over the course of the experiment it remains fairly small compared to the change in the level of ad load. Consequently, the increase in ad revenue remains fairly stable in the high ad load condition even though the listenership drops during the experiment. Plus subscription revenue grows by about 15% as users substitute to the ad-free version as the number of ads increases. Note that the impact of the treatment on ad and subscription revenues stabilizes fairly quickly post-treatment.

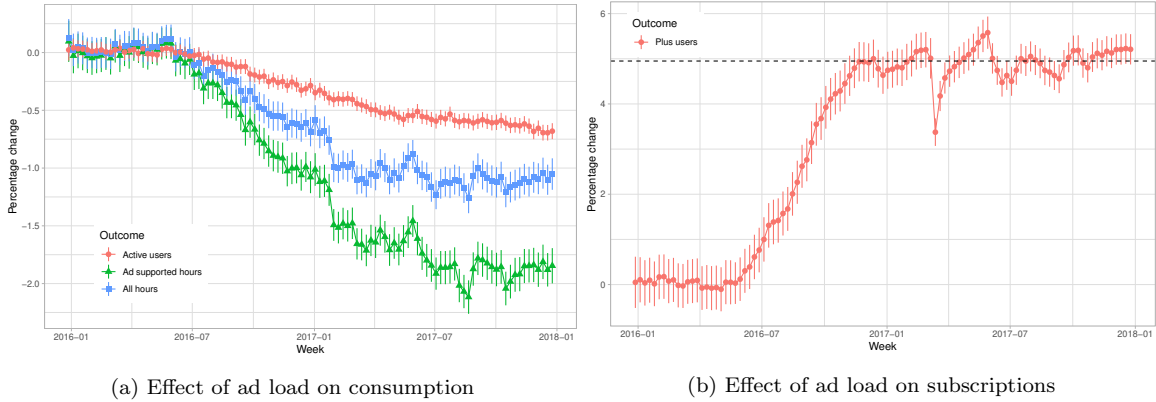


Figure 5: The instrumental variable regression estimates for the marginal effect of an additional ad per hour on consumption and plus subscriptions. While the treatment effect on consumption takes more than a year to stabilize, the effect of ad load on subscriptions stabilizes within six months of the experiment (December of 2016).

control, hence, β is interpreted as percentage changes in an outcome relative to control. As expected, we do not detect any treatment effect in the period before the experiments kicked off. Our results from Figures 4-5 establish two points:

- **Magnitude of the effect:** A one percent increase in ad load led to a 1.8% and 1% decrease in ad-supported and all hours after 18 months, respectively. This effect is about 0.7% for weekly active users. Finally, a one percent increase in ad load led to a five percent increase in plus subscribers over the same time period. These results show that the treatment effect of ad load on subscriptions is much larger than its effect on consumption and activity.
- **Time to reach steady state:** The treatment effect of ad load on listenership takes a long time to stabilize. For instance, the effect on weekly active users keeps growing even after a year into the experiment, while its impact on ad-supported hours and all hours takes more than a year to stabilize, see Figure 5. Nevertheless, the lift on ad revenues in the higher ad load condition (see Figure 4) remains stable in the course of the experiment. Lastly, the treatment effect of ad load on subscriptions stabilizes relatively quickly, i.e., by December of 2016 (within six months of the experiment).

Our findings so far show that the impact of ad load (quality of service) on plus subscription is large and materializes within six months of the experiments. Furthermore, the impact of ad load on ad revenue remains fairly stable throughout the experiment, since the impact on listenership in the same timeframe is significantly smaller than the increase in ad load.

6 Estimation models

Our goal is to study the impact of personalizing the quality of service on firm profits and consumer welfare. To optimize profits and evaluate welfare implications, we first need a demand model that reflects listeners' choice between being inactive (outside option), using the ad-supported version, or paying for the subscription service. Also, note the profit structure depends on listeners' subscription state. Although the profits generated by a paid subscriber are not a function of consumption intensity, the profits from an ad-supported user are a direct function of hours consumed and the *realized* ad load. In the following subsections, we first estimate a discrete-choice demand model where consumers choose between the outside option, ad-supported, and paid subscription. Then, conditional on being an ad-supported user, the number of hours consumed in a given period is estimated. Finally, as illustrated in section 4, the realized ad load depends on both the treatment condition and demand from advertisers. Hence, we construct a third model to account for the partial control problem discussed in section 4.¹⁵ Combining these models enables us to optimize profits and study the welfare implications of personalizing ad load.

6.1 Demand model

We estimate a nested logit model to reflect users' decision between two nests of options, namely, a degenerate nest that includes the outside option and another nest where the user decides between the ad-supported and paid versions. Consider the following discrete-choice demand model:

$$u(\boldsymbol{\tau}; \mathbf{x}) = \max \left\{ \begin{array}{ll} \epsilon_0 & \text{outside option,} \\ \underbrace{\theta(\mathbf{x}) + \sum_j \eta_j(\mathbf{x}) \mathbb{1}_{\{\tau=e_j\}} + \epsilon_a}_{v_1(\mathbf{x}, \boldsymbol{\tau})} & \text{ad-supported,} \\ \underbrace{\gamma(\mathbf{x}) + \epsilon_p}_{v_2(\mathbf{x})} & \text{paid version,} \end{array} \right. \quad (5)$$

where \mathbf{x} is a user-specific vector that includes exogenous and pre-experimental endogenous features at the listener level. The utility of consuming the outside option is normalized to zero, and the net utility of consuming the ad-supported and paid product in the control condition is captured by

¹⁵To allay concerns about selection affecting the performance in the conditional models, in the online appendix we estimate an alternative model that directly predicts the number of ads and show that our results remain the same.

$\theta(\mathbf{x}) + \sum_j \eta_j(\mathbf{x}) \mathbb{1}_{\{\boldsymbol{\tau} = \mathbf{e}_j\}}$, and $\gamma(\mathbf{x})$, respectively. This net utility consists of consumption utility and the disutility caused by ads and payment. One can impose further restrictions to disentangle these parts; however, our goal is to impose as few assumptions as possible.¹⁶ The experiment only affects the ad load across different conditions; therefore, the *treatment effect* only enters the utility of ad-supported product through $\eta_j(\mathbf{x})$. The treatment condition is represented by a binary vector $\boldsymbol{\tau}$, and \mathbf{e}_j is a unit vector whose j^{th} element is equal to 1. The treatment effect of condition j is denoted by $\eta_j(\mathbf{x})$, which measures the change in utility of consumption for each treatment arm relative to the control condition. Finally, $(\epsilon_0, \epsilon_a, \epsilon_p)$ follows a generalized extreme value (GEV) distribution that allows us to estimate a nested-logit model, where the probability of each option is as follows:

$$\mathbb{P}(Y = \text{outside, ad-supported, paid} | \boldsymbol{\tau}; \mathbf{x}) = \begin{cases} \frac{1}{1 + \exp(\Lambda \cdot IV)}, \\ \frac{\exp(\Lambda \cdot IV)}{1 + \exp(\Lambda \cdot IV)} \frac{\exp\left(\frac{v_1(\mathbf{x}, \boldsymbol{\tau})}{\Lambda}\right)}{\exp\left(\frac{v_1(\mathbf{x}, \boldsymbol{\tau})}{\Lambda}\right) + \exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)}, \\ \frac{\exp(\Lambda \cdot IV)}{1 + \exp(\Lambda \cdot IV)} \frac{\exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)}{\exp\left(\frac{v_1(\mathbf{x}, \boldsymbol{\tau})}{\Lambda}\right) + \exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)}, \end{cases}$$

where $v_1(\mathbf{x}, \boldsymbol{\tau})$, and $v_2(\mathbf{x})$ are defined as in equation (5). Also, $IV = \log\left(\exp\left(\frac{v_1(\mathbf{x}, \boldsymbol{\tau})}{\Lambda}\right) + \exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)\right)$ is the inclusive value of the nest, and $1 - \Lambda \in [0, 1]$ reflects the correlation structure inside the nest.

Functions $\gamma(\mathbf{x})$, $\theta(\mathbf{x})$, and $\eta_j(\mathbf{x})$ are parameterized as neural networks, which allows them to be represented as flexible functional forms of pre-treatment features. Note that a neural network with a terminal softmax activation layer is effectively a flexible logit, and by restricting the values fed to the terminal layer, we can create flexible structural models. Because we are estimating a nested logit model instead of a simple logit one, we need to change the terminal layer of the neural network to reflect the probability structure imposed by the nested logit with a tunable parameter Λ . We impose a few restrictions on the structure of the neural network. First, treatment dummies only enter the last layer of the neural network and are multiplied by coefficients $\eta_i(\mathbf{x})$ that are also parameterized as neural networks. This technique forces¹⁷ the neural network to use the information provided by

¹⁶Note that as one increases ad load listeners substitute to the outside option or the paid version. The variation in the ratio of substitution to the paid version relative to the outside option across different user segments helps us identify $\gamma(\mathbf{x})$. This means that even in the absence of price variation by imposing further structure on the utility function one can disentangle the utility from consumption and the disutility from payment. However, this is beyond the scope of this paper.

¹⁷Note the treatment effect tends to be very small, and if the treatment dummies are inserted in the input layer along with other features, they may get regularized out by the network. The fact that we use a number of shared layers to construct $\eta_j(\mathbf{x})$ and then spread into separate heads forces the model to use the information provided by treatment dummies and improves the statistical power of our algorithm; see Shalit, Johansson, and Sontag (2017) for more details.

treatment dummies and has been used in Shalit, Johansson, and Sontag (2017) and Farrell, Liang, and Misra (2021). Second, because features that may explain heterogeneity in treatment $\eta_i(\mathbf{x})$ may be very different from those that explain the cross-sectional heterogeneity ($\theta(\mathbf{x})$ and $\gamma(\mathbf{x})$), we use split neural networks (Kim et al., 2017) to separately fit values to each part. In particular, we let the network that learns $\eta_j(\mathbf{x})$, that is the “treatment effect,” be disjoint from the network that learns $\theta(\mathbf{x})$ and $\gamma(\mathbf{x})$, that is the utility from consuming the ad-supported product. Split neural networks have been mainly used for parallel computing and to boost the training process. In this application, however, the issue is that the heterogeneity in the treatment effect may be explained by different types of features that would explain the cross-sectional heterogeneity, and these exclusion restrictions help the network to learn treatment heterogeneity more efficiently. Note that even though we impose a split structure, the networks are trained jointly. A schematic view of the architecture is presented in Figure 6. Note the purple part that learns $\theta(\mathbf{x})$, and $\gamma(\mathbf{x})$ is separate from the green part, which is responsible for explaining the variation caused by the treatment ($\eta_i(\mathbf{x})$). Furthermore, change in the ad load affects only the utility of consuming the ad-supported version.

To train this model, we minimize a weighted negative log-likelihood function similar to Shalit, Johansson, and Sontag (2017), which jointly optimizes for the treatment effect $\tau(\cdot)$ and parameters that explain cross-sectional variation $\theta(\cdot)$ and $\gamma(\cdot)$:

$$\underset{\theta, \gamma, \eta, \Lambda}{\text{minimize}} \quad \frac{1}{N} \sum_i w_i \mathcal{L}(h(\theta(\mathbf{x}_i), \gamma(\mathbf{x}_i), \eta(\mathbf{x}_i)), y_i; \Lambda), \quad (6)$$

where $\mathcal{L}(\cdot, \cdot)$ is the negative log-likelihood for the nested-logit model, $\theta(\cdot)$, $\gamma(\cdot)$, and $\eta(\cdot)$ are functions parameterized by the neural networks. Finally, N is the total number of users, and w_i is an inverse propensity score for each treatment condition that is equal to $\frac{N}{\sum_{n=1}^N \mathbb{1}_{\{\tau_n = e_i\}}}$. Note that inverse probability weighting is often used to address selection issues; however, our goal is to balance different treatment conditions. For instance, if, by design, the size of one treatment cell is significantly larger than another cell, the optimization problem will have more incentive to fit the data to improve its prediction power for the larger cell. For instance, if a treatment group is twice as large as another treatment group, the same type of error is penalized twice for the larger treatment cell relative to the smaller one. However, our goal is to better learn the *differences* across treatment cells, and this weighting balances the prediction power across different counterfactual scenarios (Shalit, Johansson, and Sontag, 2017).

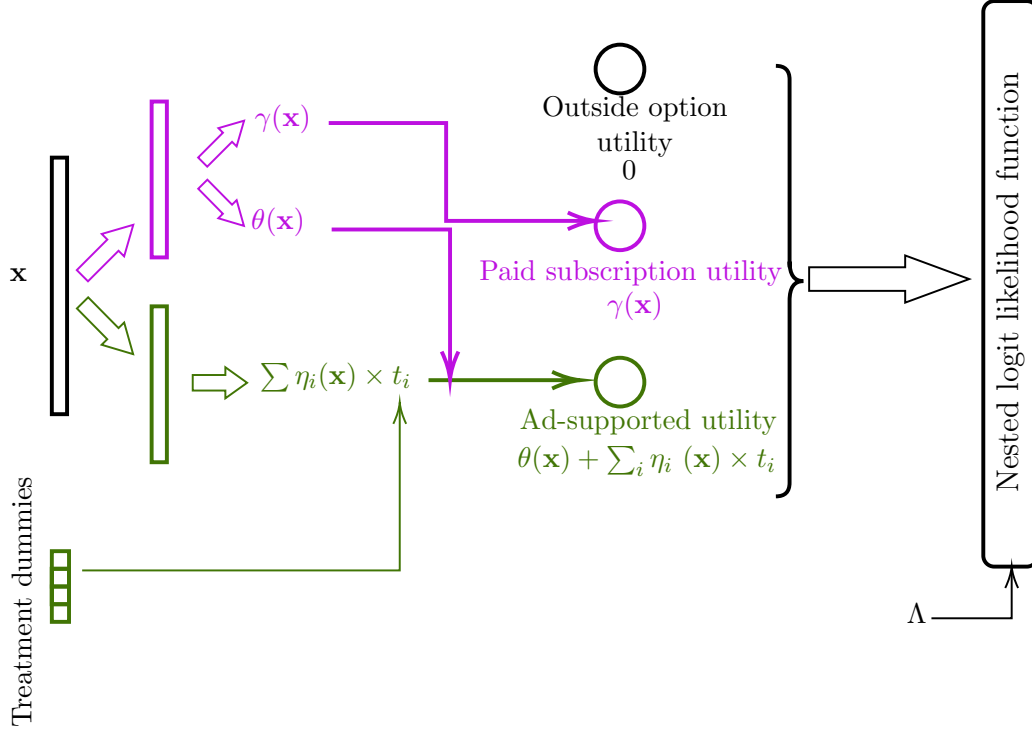


Figure 6: A schematic view of the split neural network architecture used for demand estimation. The purple and green parts of the network are separate. Also, note that treatment dummies enter right before the last layer. This restriction imposes structure on the network and forces it to learn the relationship between the output and ad load even though the amount of variation explained by the treatment dummies could be very small. Furthermore, the split structure of the network allows it to separately learn features that could explain the treatment effect $\eta_i(\mathbf{x})$ from other constructs that explain cross-sectional differences, that is $\theta(\mathbf{x})$ and $\gamma(\mathbf{x})$.

6.2 Intensive margin

A change in ad load not only affects the extensive margin of consumption and the choice across different products within the product line but also affects the intensive margin of consumption, that is hours spent listening to music conditional on being an ad-supported listener. Consider the following model:

$$\log(\mathbf{Y}) = \alpha(\mathbf{x}) + \sum_j \nu_j(\mathbf{x}) \mathbb{1}_{\{\tau=e_j\}} + \epsilon, \quad (7)$$

where \mathbf{x} is a set of exogenous and user-generated features collected during the pre-treatment period. \mathbf{Y} is the number of hours consumed conditional on being an active ad-supported user, $\alpha(\mathbf{x})$ represents the conditional expectation of $\log(\mathbf{Y})$ in the control condition, and the $\nu_i(\mathbf{x})$ capture the treatment effect of assigning a user with features \mathbf{x} to condition j on ad-supported hours relative to the control condition. Finally, ϵ is a random variable with a Normal distribution.

To estimate this model, we again resort to a split neural network model, where one split learns

$\alpha(\mathbf{x})$ and the other one fits $\nu_i(\mathbf{x})$. Note $\alpha(\mathbf{x})$ and $\nu_i(\mathbf{x})$ are estimated jointly; however, the networks that estimate them do not share weights and are allowed to have different parameters similar to the architecture presented in Figure 6. This model is learned by optimizing a weighted ℓ_2 loss counterpart of (6).

6.3 Partial control over realized treatment

As demonstrated in Table 3, the realized ad load is not necessarily equal to the intended ad load. In other words, the experiments only shift the ad capacity, which is the number of opportunities to show ads to each listener; however, the realized ad load depends on the advertisers' interest in different demographics. Therefore, the treatment depends on both the firm's actions and the advertisers' demand, and it needs to be accounted for in our optimization problem for reallocating ads. To that end, we estimate a model similar to that in section 6.2, with \mathbf{Y} being the realized ad load conditional on being an ad-supported listener. The rest of the parameters and the estimation procedure are similar to section 6.2.

7 Training and out-of-sample prediction power

In this section, we discuss how the models are trained and demonstrate that they have prediction power in the hold-out sample. We use the data from December 2016 for training our models. This choice is made for two reasons: (a) as demonstrated in Figures 4-5, the impact of the treatment on revenue sources (ads delivered and subscription) stabilizes by December of 2016, and (b) from a practical standpoint our goal is to select the earliest time where the effects are large enough to effectively sort users based on heterogeneous treatment effects of the intervention on ad and subscription revenues. In online appendix D, we show that the estimated heterogeneous treatment effects from December 2016 are persistent and have explanation power a year later (December of 2017). In other words, while the magnitude of the effects could fluctuate, the rank ordering of users in terms of lift in ads or subscription revenues remains consistent. The persistence in the rank ordering of heterogeneous treatment effects across users is important for practitioners because it means that a short-run six months experiment, which is much less costly to run than a long-run 18 months experiment, would suffice for learning reliable and persistent treatment effects to optimize revenues.

To train the models, we randomly divided the data set into halves. We train and validate the

models on one-half of the data. The other half (the hold-out sample) is used for demonstrating the prediction power of the model in this section and simulating counterfactuals in the next section, where we solve the firm’s optimization problem. In the remainder of this section, we first introduce some notation to describe the conditional average treatment effects (CATE) of ad load on the number of ads delivered and subscription status using the models introduced in Section 6. Next, we examine the heterogeneity in CATEs for ads and subscription revenues across treatment conditions. Finally, we show that the predicted CATEs do have explanation power in the hold-out sample.

Let \mathbf{x} be the set of pre-treatment outcomes and user features used for describing each individual.¹⁸ Also as before, let τ represent each of the seven experiment conditions. Note that the treatment assignment is persistent, and each listener can participate in only one of the treatment cells. The models introduced in section 6 enable us to predict user outcomes across “counterfactual” experiment conditions. We built three sets of models, and we use the following notation to refer to them:

- **Extensive margin:** Let $P_0(\mathbf{x}, \tau)$, $P_a(\mathbf{x}, \tau)$, and $P_s(\mathbf{x}, \tau)$ denote the conditional probability of choosing the outside option, ad-supported service, and the paid subscription as a function of pre-treatment user features \mathbf{x} and treatment condition τ . These conditional probabilities are the output from the estimated models in section 6.1.
- **Intensive margin:** Let $C(\mathbf{x}, \tau)$ be the expectation of the number of ad-supported hours consumed in a given period conditional on being an ad-supported user. $C(\mathbf{x}, \tau)$ would be the output of the model discussed in section 6.2.
- **Realized treatment:** Let $A(\mathbf{x}, \tau)$ be the conditional expectation of the realized ad load (number of ads per hour). $A(\mathbf{x}, \tau)$ is the output of the machine learning model discussed in section 6.3.

To demonstrate the effectiveness of our approach, we show that our models are able to detect heterogeneous treatment effects in terms of change in revenue from subscriptions and ads. First, we calculate the conditional average treatment effect on the lift in the subscription propensity for each condition relative to the control condition as:

$$\zeta_s(\mathbf{x}_i, e_j) = P_s(\mathbf{x}_i, e_j) - P_s(\mathbf{x}_i, e_0), \quad (8)$$

¹⁸We use over 120 different features including user consumption features, subscription status, session activity, and interaction with ads in the pre-treatment period. We also augment our data with census data to include socioeconomic and demographic information across different geographic areas.

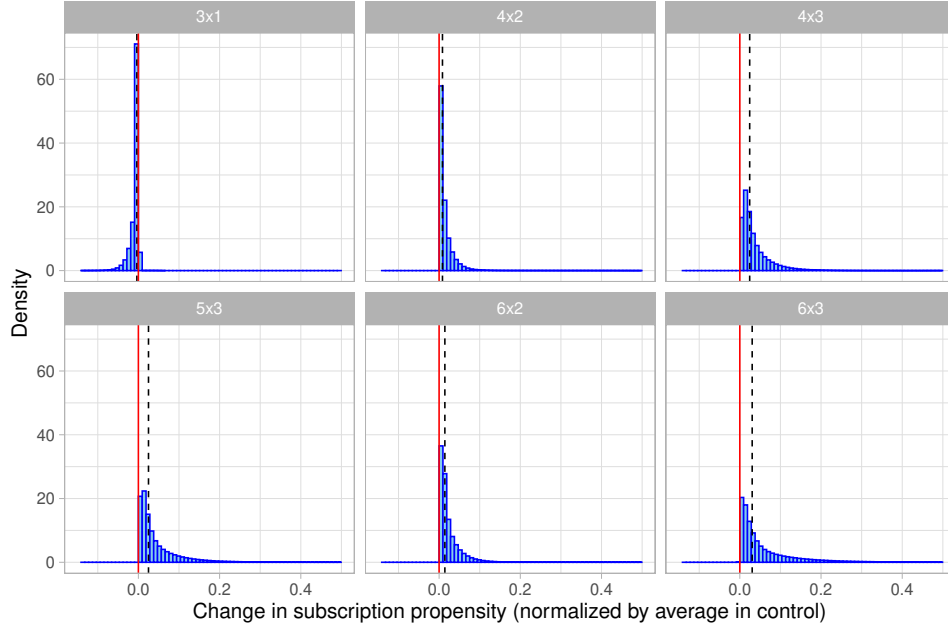
where $P_s(\mathbf{x}_i, e_j)$, and $P_s(\mathbf{x}_i, e_0)$ are the predicted probability of user i subscribing to the paid service when exposed to treatment condition j , and control, respectively. $\zeta_s(\mathbf{x}_i, e_j)$ is the expected lift in the propensity of subscription when a user with features \mathbf{x}_i is moved from control to condition j . Users' response to changes in ad load is likely heterogeneous and $\zeta_s(\mathbf{x}_i, e_j)$ captures the heterogeneous response in terms of subscribing to the paid service. For example, moving different types of users from say control to the 6x3 condition could have very different effects on the expected subscription revenues generated by them; Users with more disposable income or those who derive higher disutility from ads may be more likely to move to the subscription service. Figure 7a displays the histogram for predicted CATEs of treatment condition on subscription status relative to control, $\zeta_s(\mathbf{x}_i, e_j)$, for users in the hold-out sample.

Similarly, for examining the heterogeneous treatment effects on ad revenues, we calculate the conditional average treatment effect on the number of ads delivered for each condition relative to the control condition as:

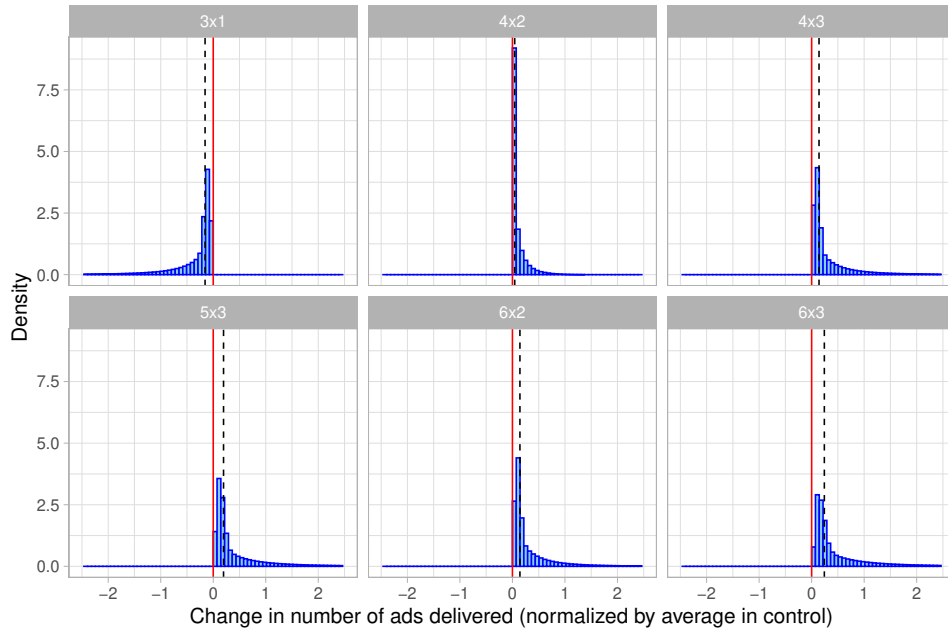
$$\zeta_a(\mathbf{x}_i, e_j) = P_a(\mathbf{x}_i, e_j) \cdot C(\mathbf{x}_i, e_j) \cdot A(\mathbf{x}_i, e_j) - P_a(\mathbf{x}_i, e_0) \cdot C(\mathbf{x}_i, e_0) \cdot A(\mathbf{x}_i, e_0), \quad (9)$$

where $P_a(\mathbf{x}_i, e_j)$, $C(\mathbf{x}_i, e_j)$, and $A(\mathbf{x}_i, e_j)$ are the predicted probability of user i being an active ad-supported user when exposed to treatment condition j , the number of hours consumed conditional on being an active ad-supported user, and the number of ads received per hour (ad load), respectively. $P_a(\mathbf{x}_i, e_0)$, $C(\mathbf{x}_i, e_0)$, and $A(\mathbf{x}_i, e_0)$ are defined similarly for the control condition, and e_0 is a vector of all zeros. $\zeta_a(\mathbf{x}_i, e_j)$ is the expected lift in the number of ads shown to a user with features \mathbf{x}_j when moved from control to condition j . Note that moving different types of users from say control to the 6x3 condition could have very different effects on the expected ad revenue generated from them. In particular, the lift in ad revenues depends on: (i) how a listener might react in terms of switching out of the ad-supported service, $P_a(\mathbf{x}_i, e_j)$, (ii) the hours spent listening to music conditional on being an ad-supported user, $C(\mathbf{x}_i, e_j)$, and (iii) listener's attractiveness for advertisers, $A(\mathbf{x}_i, e_j)$.¹⁹ Figure 7b presents the histogram of predicated CATEs on number of ads delivered, $\zeta_a(\mathbf{x}_i, e_j)$, in the hold-out sample.

¹⁹The conditional models $C(\mathbf{x}_i, e_j)$ and $A(\mathbf{x}_i, e_j)$ are trained on the selected sub-sample of ad-supported users who had non-zero consumption. While we evaluate the models on the hold-out sample and use inverse propensity weighted estimates which are model-free, one might worry that this selection may degrade the performance of the prediction model. In online appendix E, we show that using a model that directly learns the number of ads delivered provides similar performance.



(a) Subscriptions



(b) Ads

Figure 7: The histogram of predicted change in subscription propensity and number of ads delivered for each treatment arm relative to control in the hold-out sample. The median treatment effect across each condition is represented by the dashed line, and the solid red line represents zero. Due to our data protection agreement with Pandora, we normalize the lifts by dividing them by the sample average in the control group.

We now demonstrate that the conditional average treatment effects on subscription propensities and the number of ads delivered, $\zeta_s(\mathbf{x}_i, e_j)$ and $\zeta_a(\mathbf{x}_i, e_j)$, do indeed have prediction power in the hold-out sample. We break users in the hold-out sample into five quintiles based on the predicted

treatment effects of the 6x3 relative to the control condition for both subscriptions and ads²⁰. Then, we use the hold-out sample to estimate the *realized* lift in the subscription propensity and number of ads delivered in the 6x3 condition relative to control across these quintiles. The results are presented in Figure 8 and show the models are indeed able to capture the heterogeneity in the lift in subscription and ad revenues across users. Note the models used for making these predictions, that is sorting users into five groups, were not trained on the hold-out sample, and this prediction power on the hold-out sample demonstrates that the model was able to learn meaningful patterns that generalize beyond the training set.

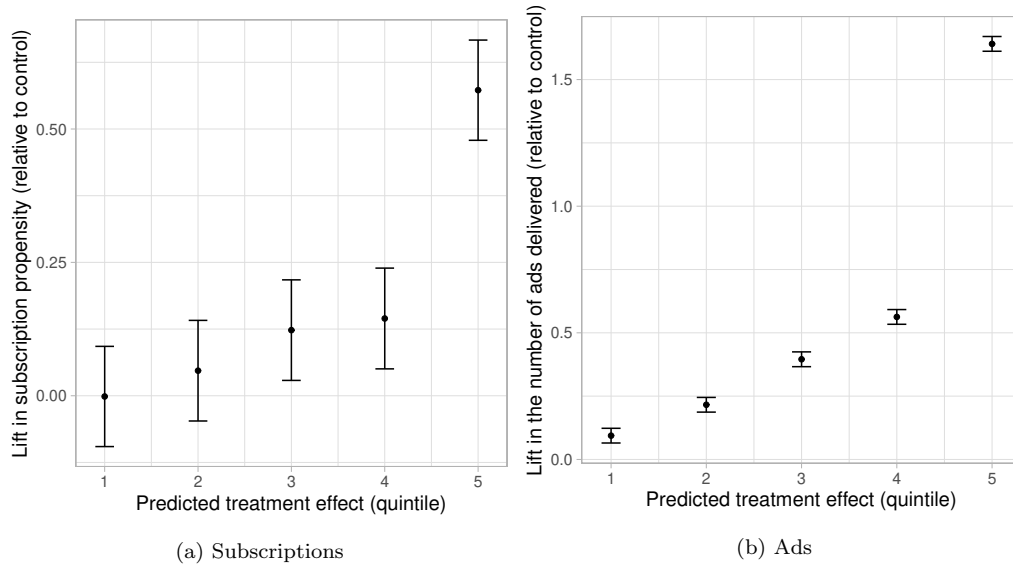


Figure 8: Realized lift in subscription propensity and ads delivered in the 6x3 condition relative to control as a function of predicted treatment effect quintile in the hold-out sample.

Our results in Figures 8 illustrate two interesting points:

- **Lift in subscription propensity:** The lift in subscription propensity for the bottom quintile in Figure 8a is centered around zero, which means the model has identified customer segments who are very unlikely to subscribe in response to an increase in ad load, i.e., when moved from control to the highest ad load condition (6x3).
- **Inelastic demand for the ad-supported service:** The lift in the number of ads delivered in Figure 8b is always positive and statistically different from zero including for those users who are in the bottom quintile. This reflects that demand for the ad-supported service is fairly inelastic, since the decline in consumed hours or active users, see Figure 4, due to higher ad load does not

²⁰Recall that the 6x3 condition is the most extreme treatment condition in our experiment with the largest treatment effect.

lead to lower ad revenues even for the bottom quintile in Figure 8b. In online appendix D, we show that these patterns persist even in December 2017, 18 months after the experiment started.

These patterns imply that moving all users from control to the highest ad load condition (6x3) would improve the expected revenues from both subscriptions and ads. However, this change would mean that the firm should also sell more ads on the advertiser side to increase the ad inventory size. This decision would affect the price of ad impressions, see the levers in Figure 2, and also impacts ad revenues. Unfortunately, we do not observe the closing price of the individual advertising contracts nor do we have exogenous variation in the price of impressions and the realized ad inventory. Note that our ad load experiments do not shift the overall ad inventory size as they only affect a small portion of the overall user base. Consequently, we embark on answering a more nuanced problem, which is whether holding the ad inventory fixed can the firm improve profits by reallocating ads across individuals.

8 Optimizing profits

To optimize profits, one needs to understand the heterogeneity both in terms of change in subscription rates and the number of ads delivered. To construct estimates for the number of ads delivered and subscription rates, we use the sets of models presented in section 6, namely, a discrete-choice model that reflects the user's decision between the outside option, ad-supported consumption, and the subscription service, a model that determines the number of hours consumed conditional on using the ad-supported service, and a model that estimates the realized ad load across different treatment conditions. Throughout this section, we rely on the notation developed in section 7.

Recall from Figure 2 that the profits depend on three sets of decisions: (a) allocation of ads across listeners, (b) price of the paid subscription, and (c) price of ad impressions. The exogenous variation in our experiment comes from (a). The price of subscription service (b) was held fixed throughout the experiment. The firm can change the price of ad impressions (c) and that would affect its ad inventory size. As discussed above, even with a fixed ad inventory and without exogenous variation in (b) and (c), the firm faces an implicit pricing problem that involves allocating ads across individuals while satisfying constraints imposed by advertisers. We show that the firm can increase subscription profits by reallocating ads across individuals while satisfying the constraints imposed by advertisers. Given the ad inventory (Γ), the firm's optimization problem translates to:

$$\begin{aligned} \underset{\boldsymbol{\tau}}{\text{maximize}} \quad & \sum_i m_s P_s(\mathbf{x}_i, \boldsymbol{\tau}_i) \\ & \sum_i P_a(\mathbf{x}_i, \boldsymbol{\tau}_i) C(\mathbf{x}_i, \boldsymbol{\tau}_i) A(\mathbf{x}_i, \boldsymbol{\tau}_i) = \Gamma, \end{aligned} \quad (10)$$

where i indexes users and m_s is the margin from subscriptions. The objective function is the expected profits from subscriptions across all users. Γ is the total number of ads available in the inventory, and the constraint ensures that the ads in the inventory are served. Note that the constraint accounts for the fact that moving different types of users (\mathbf{x}_i) across treatment conditions ($\boldsymbol{\tau}_i$) has heterogeneous impact on ads served both through heterogeneous changes in activity/consumption, captured by $P_a(\mathbf{x}_i, \boldsymbol{\tau}_i)$ and $C(\mathbf{x}_i, \boldsymbol{\tau}_i)$, and also heterogeneity in the realized ad load due to differences in advertisers' demand for different types, which is captured by $A(\mathbf{x}_i, \boldsymbol{\tau}_i)$. The optimization problem (10) assigns each user to one of the seven cells to maximize profits while serving Γ ads. In other words, $\boldsymbol{\tau}_i$ is a binary vector that takes one of the seven values $\{e_0, e_1, \dots, e_6\}$, where e_0 is the vector of all zeros, and e_i is a unit vector whose i^{th} element is non-zero. This discrete optimization problem aims at assigning each user to one of the seven conditions, that is 7^N different combinations in total.

The functions above are constructed using the estimates of models developed in section 6. For instance, note $P_s(\mathbf{x}, \boldsymbol{\tau}) = \frac{\exp(\Lambda.IV)}{1+\exp(\Lambda.IV)} \frac{\exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)}{\exp\left(\frac{v_1(\mathbf{x}, \boldsymbol{\tau})}{\Lambda}\right) + \exp\left(\frac{v_2(\mathbf{x})}{\Lambda}\right)}$, where $v_1(\mathbf{x}, \boldsymbol{\tau})$, $v_2(\mathbf{x})$, and Λ are estimated parameters from the model described in section 6.1. The rest of the functions are also outputs from estimated models presented in section 6. Note (10) is a discrete non-convex optimization problem with a non-convex constraint, and even finding a local optimum is NP-Hard for general continuous non-convex problems in the worst case (Murty and Kabadi, 1985). In its current form, the problem is intractable. To approach it, we use the Lagrangian relaxation of (10), which yields

$$\underset{\boldsymbol{\tau}}{\text{maximize}} \quad \sum_i m_s P_s(\mathbf{x}_i, \boldsymbol{\tau}_i) + \lambda(\Gamma) \left(\sum_i P_a(\mathbf{x}_i, \boldsymbol{\tau}_i) C(\mathbf{x}_i, \boldsymbol{\tau}_i) A(\mathbf{x}_i, \boldsymbol{\tau}_i) \right), \quad (11)$$

where $\lambda(\Gamma)$ is the marginal impact of ads on subscriptions. Different values of $\lambda(\Gamma)$ lead to different ad inventory sizes. Note that given λ , we can swap the maximization with the summation in 11 and

the problem can now be decoupled across users:

$$\sum_i \max_{\boldsymbol{\tau}_i} \underbrace{m_s \cdot P_s(\mathbf{x}_i, \boldsymbol{\tau}_i) + \lambda \cdot (P_a(\mathbf{x}_i, \boldsymbol{\tau}_i)C(\mathbf{x}_i, \boldsymbol{\tau}_i)A(\mathbf{x}_i, \boldsymbol{\tau}_i))}_{f(\mathbf{x}_i, \lambda, \boldsymbol{\tau}_i)}, \quad (12)$$

which is simply equivalent to evaluating $f(\mathbf{x}_i, \lambda, \boldsymbol{\tau}_i)$ for each user i across different treatment conditions and choosing the maximum.

Given λ , the complexity of problem (12) is $\mathcal{O}(N)$ compared to $\mathcal{O}(7^N)$ for (10). Note the shadow price $\lambda(\Gamma)$ corresponding to each value of Γ is not a priori known and one needs to resolve problem (10) for different shadow-price values λ to find the corresponding ad inventory level Γ . However, this task can be done easily with binary search techniques. Given a fixed λ , an ad load assignment policy $\mathcal{P}_\lambda(\mathbf{x}_i)$ is defined as:

$$\mathcal{P}_\lambda(\mathbf{x}_i) = \operatorname{argmax}_{\boldsymbol{\tau}_i} f(\mathbf{x}_i, \lambda, \boldsymbol{\tau}_i). \quad (13)$$

To evaluate the performance of the personalization policy $\mathcal{P}_\lambda(\mathbf{x}_i)$, we use inverse propensity weighted estimates similar to Hitsch and Misra (2018), Simester, Timoshenko, and Zoumpoulis (2020), and Yoganarasimhan, Barzegary, and Pani (2022). This means that our counterfactual estimates of the personalized policy's performance are based on realized outcomes and randomization in the hold-out sample and are model-free. Thus, while our models may be imperfect, our estimates of lift in profits, ad load, subscriptions, or other performance metrics related to the policy are measured using inverse propensity weighted estimates which are not extrapolations of a structural model. We construct an inverse probability weighted estimator for different parameters interest y as follows:

$$\hat{\Pi}_y(\mathcal{P}_\lambda) = \frac{1}{N} \sum_{i=1}^N \sum_{j=0}^6 w_j \mathbb{1}_{\{\boldsymbol{\tau}_i=e_j\}} \mathbb{1}_{\{\mathcal{P}_\lambda(\mathbf{x}_i)=e_j\}} y_i, \quad (14)$$

where $w_j = \frac{N}{\sum_{i=1}^N \mathbb{1}_{\{\boldsymbol{\tau}_i=e_j\}}}$ is the inverse propensity for each treatment condition and is a fixed number for each of the seven treatment conditions because we have a randomized control trial. The randomized treatment assignment for user i is denoted by $\boldsymbol{\tau}_i$, and $\mathcal{P}_\lambda(\cdot)$ is a policy function that assigns each user i to one of the seven treatment conditions. Our goal is to evaluate the performance of a policy $\mathcal{P}_\lambda(\cdot)$ devised by solving (13). The product of the two indicator functions in (14) filters out observations for which the assignment rule \mathcal{P}_λ and the randomized treatment $\boldsymbol{\tau}_i$ coincide. Finally, y_i is an outcome of

interest for user i , for example, subscription status or the number of ads received at a given point in time. Equation (14) then provides a consistent estimator of expectation of y_i under a given assignment rule \mathcal{P}_λ .

To illustrate gains to personalizing ad load, we vary λ and solve (13) for users in the hold-out sample to get a policy \mathcal{P}_λ . Then, we evaluate the average number of ads realized under the assignment rule \mathcal{P}_λ and the expected profits from the subscription service during December of 2016 under this policy using (14). Essentially, for each λ , we get a point on the 2D plane whose x coordinate is the average number of ads realized, and its y coordinate is the expected profits from subscriptions across individuals. We vary the shadow price λ to generate points across the Pareto frontier. The results are plotted in Figure 9. Each purple dot in Figure 9 reflects the performance of the policy \mathcal{P}_λ for different values of λ , which translates to different levels of ads served on the platform. Each pink dot represents the performance of a personalized ad-allocation strategy. Our agreement with Pandora prevents us from sharing the actual numbers in dollar terms, and the performance has been reported relative to the control condition. First, note that the uniform ad allocation policies, e.g., 3x1, control, 4x2, and so forth, are co-locating around the straight dashed line. The personalized counterpart of the control condition, that is the pink dot inside the star with the same x coordinate as the control, leads to 7 percent more subscribers than the control condition. Furthermore, the distance between the straight line, which corresponds to the performance of a uniform policy, and the star is about 22%. Therefore, if the firm's goal was to achieve the same number of subscribers using a uniform ad-load strategy, it would have to increase its ad load by more than 20%.

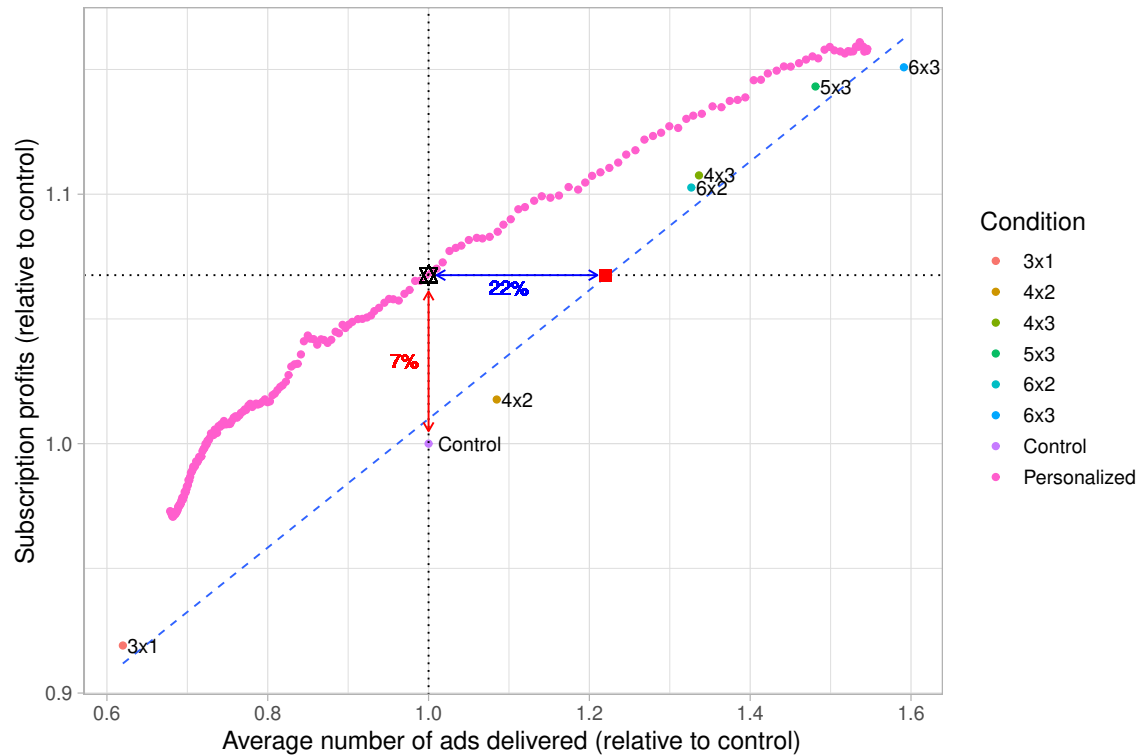


Figure 9: Change in subscription profits as a function of the number of ads served. The uniform allocation policies co-locate around the straight dashed line. Each pink dot represents the performance of a personalized assignment rule. Holding fixed the number of ads served, the personalized assignment strategy dominates uniform ad-load strategies that the firm experimented with. The star is the personalized counterpart of the control condition and leads to 7 percent more subscribers than the control condition.

So far, we have examined the gains from personalization after a six-month period, that is during December of 2016. To find the minimum time required for such gains to materialize, we compare the control condition with its personalized counterpart throughout time. Note (14) provides a consistent estimator of any outcome of interest in any given week. Figure 10 compares the 6x3 and 3x1 conditions, and the personalized counterpart of the control condition relative to the control; that is each outcome of interest is measured relative to the control. The results demonstrate that the control and its counterpart lead to similar realized ad load. However, the counterpart increases the Plus subscription rate by 7%. Whereas the impact on subscriptions manifests within three to four months after personalizing the ad load, the effect on all and ad-supported hours in the same time period seems to be negligible. This finding indicates a dynamic optimization model may be beneficial here. However, due to the persistent nature of treatment in our experiments, we cannot evaluate the benefits of dynamic implicit pricing using the current randomized control trial.

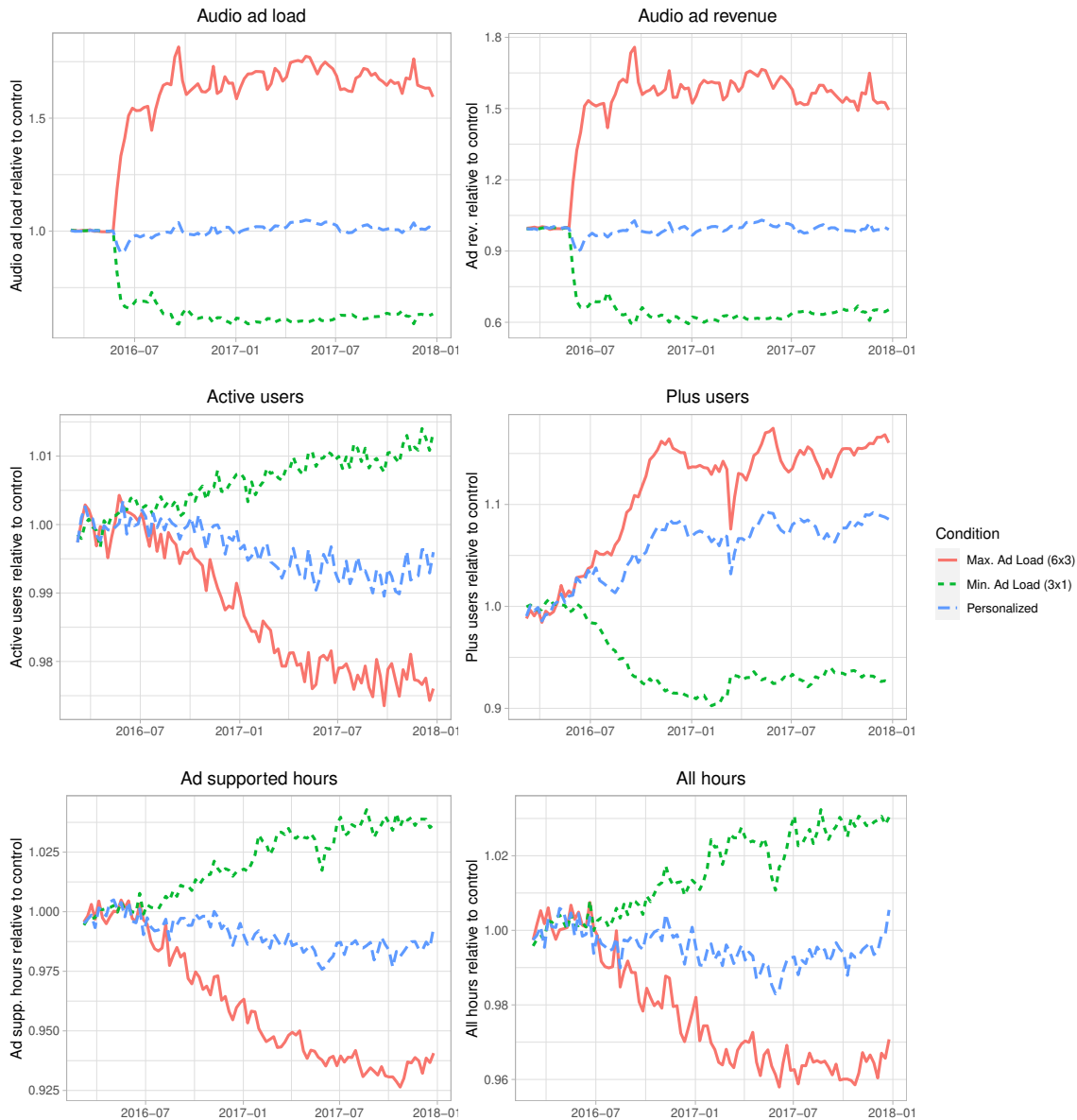


Figure 10: The effect of personalization throughout time. The top-left panel shows the personalized counterpart of the control condition is delivering approximately the same amount of ads as control, whereas the 6x3 condition is delivering 60% more ads relative to control. The top-right figure shows the personalized counterpart increases the number of subscribers by 7%, and this gain is expected to materialize within three months of implementing this policy. Note the impact on ad-supported hours or all hours within three months of implementation seems to be negligible.

We now illustrate the underlying mechanism that enables the algorithm to improve the firm's profits. Pandora offers two types of products: the subscription service (high-tier) and the ad-supported (low-tier) product. When the menu of products cannot be personalized, the problem is similar to the one discussed in Mussa and Rosen (1978) and Deneckere and Preston McAfee (1996). In the absence of personalization, the seller has the incentive to lower the quality of the low-tier (ad-supported) product

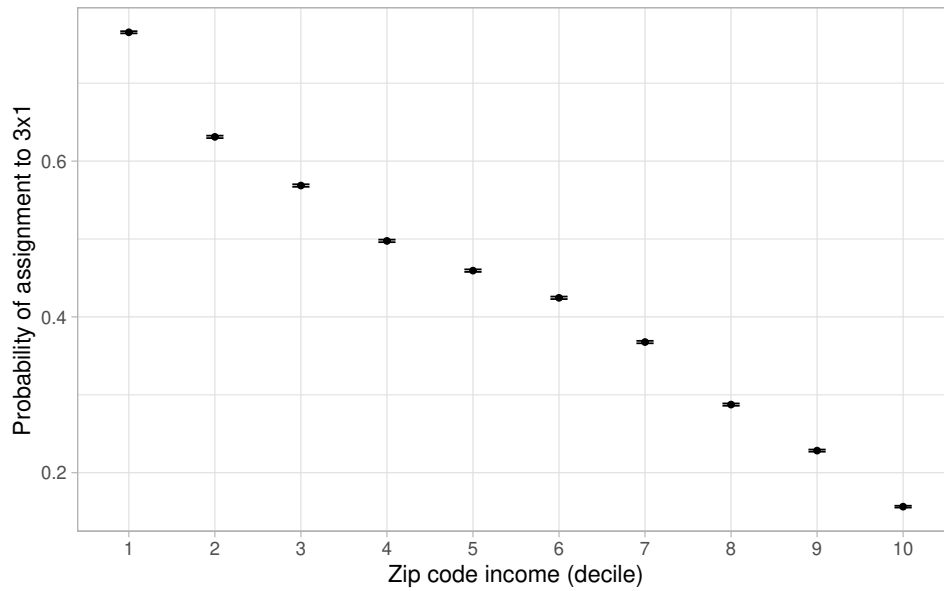
for everybody to make adopting the high-tier product worthwhile for those who have a higher willingness to pay. However, personalizing ad load limits the distortion to high willingness to pay customers, and the “implicit price” for other segments falls (quality improves) at the expense of this segment. In the next section, we investigate how our policy allocates ads across listeners of different age and income groups. We also investigate the welfare effects of the personalized policy.

9 Distribution of welfare

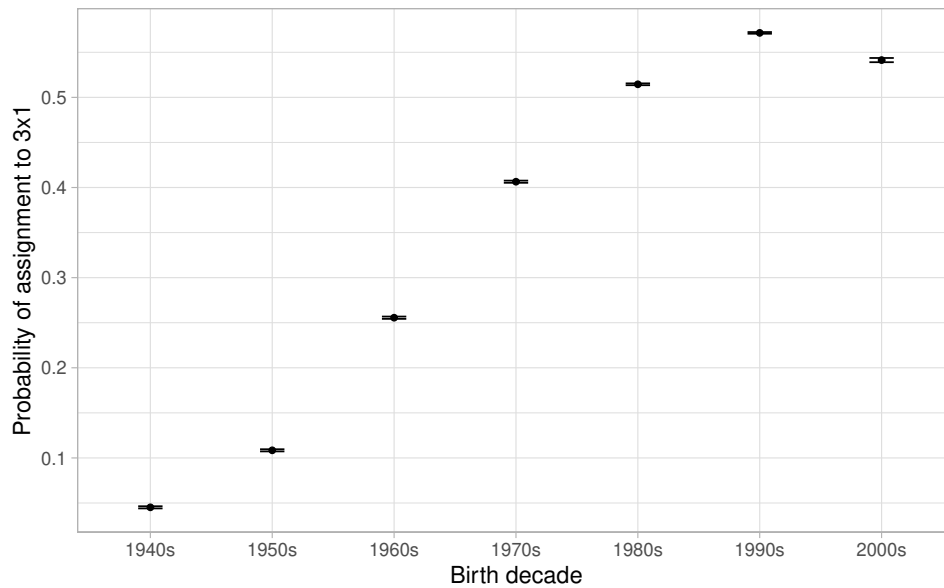
As we mentioned before, the welfare implications of personalizing ad load are a priori ambiguous due to possible correlation between willingness to pay in time and money units. For instance, whether higher-income users would be assigned to higher or lower-ad-load conditions would be unclear, whereas, in a single-product pricing problem, one would expect low-income users to face lower prices because they are more price sensitive. To examine the proposed policy, we compare the allocation of ads in the control condition with its personalized counterpart, that is the personalized assignment rule that serves the same number of ads as the control condition (marked by the star in Figure 9). Note the policy assigns each user to one of the seven conditions, and the only experiment arm that has a lower ad load than control is the 3x1 condition. The propensity of being assigned to the lower-ad-load condition, that is higher quality of service for the ad-supported product, tends to be monotonically increasing as a function of zip code income and decreasing as a function of age; see Figure 11. Both younger individuals and those residing in lower-income zip codes tend to be more price sensitive, and if the algorithm were to charge prices in dollar amounts, it would likely charge them a lower price. In our case, the company is following a uniform pricing scheme for the paid service; however, our personalization algorithm induces wealthier individuals to upgrade to the paid service by personalizing the ad load of the ad-supported product and provides a better quality of service to other demographic groups.

We now illustrate the impact of our ad-allocation algorithm on consumer welfare by comparing the control condition with its personalized counterpart, that is the personalized algorithm that delivers the same overall number of ads. The overall utility for an individual with features \mathbf{x} who is assigned to treatment condition τ is equal to:

$$U(\mathbf{x}, \tau) = \log \left(1 + \left[\exp\left(\frac{v_1}{\lambda}\right) + \exp\left(\frac{v_2}{\lambda}\right) \right]^\lambda \right), \quad (15)$$



(a)



(b)

Figure 11: Probability of assignment to lower ad load than control as a function of price sensitivity. Note the total number of ads for this assignment was set to be equal to the control condition, and those assigned to the 3x1 condition are effectively receiving a higher quality of service on the ad-supported product. (a) Users from lower-income zip codes tend to be more likely to receive an ad-load “discount.” (b) Older users tend to be less likely to receive a discount. The algorithm seems to be adjusting the quality of service for users who have a higher willingness to pay to make converting incentive compatible for them.

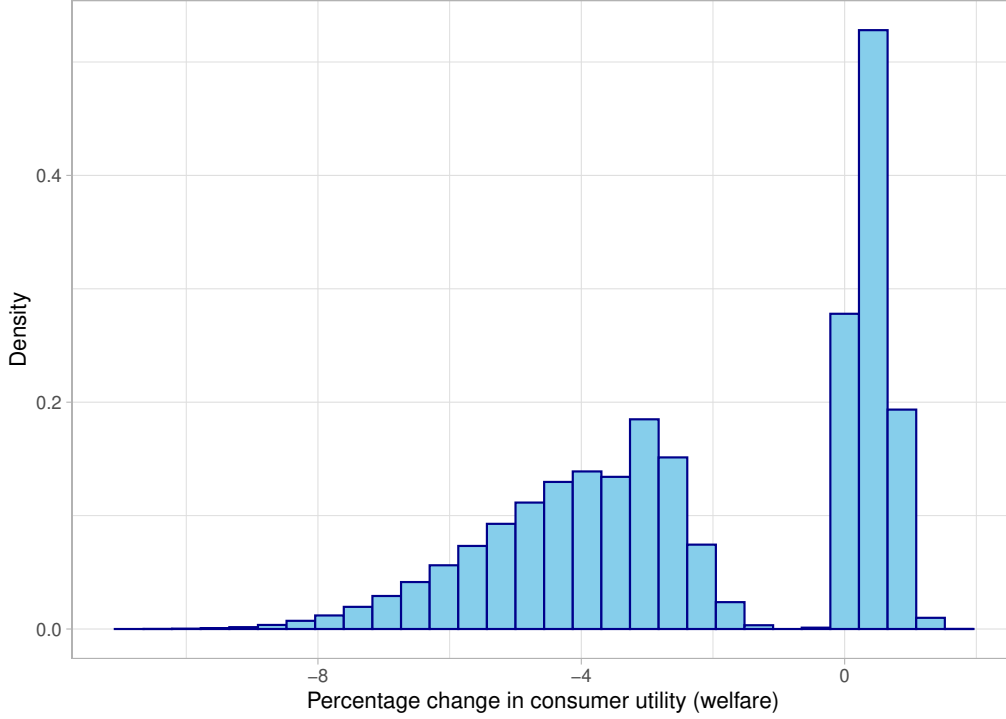


Figure 12: The impact of ad-load personalization on consumer welfare. The figure compares the percentage change in consumer utility across the control condition and its personalized counterpart. On average, personalizing ad load lowers consumer utility by 2%.

where v_1 and v_2 are the utilities associated with the ad-supported, and paid-subscription products, respectively. Recall that v_1 and v_2 are functions of \mathbf{x} and $\boldsymbol{\tau}$ and were defined in (5). To study the impact of our personalization model on consumer surplus, we can compare the percentage change in the utility of users in the control condition relative to its personalized counterpart. In particular, we examine the following construct:

$$\Delta U_{\mathcal{P}} = 100 \times \frac{U(\mathbf{x}, \mathcal{P}(\mathbf{x})) - U(\mathbf{x}, \text{control})}{U(\mathbf{x}, \text{control})},$$

where $\mathcal{P}(\mathbf{x})$ denotes the policy that is the personalized counterpart of the control condition. The distribution of $\Delta U_{\mathcal{P}}$ is plotted in Figure 12. On average, consumer utility drops by -2% and utility improves for 41.2% of users.

We now illustrate the impact of this policy on users from different age and income groups in Figure 13. The results demonstrate the loss in consumer utility is more pronounced for older users and those from higher-income zip codes. This observation is consistent with our prior findings in Figure 11 that showed younger users and those from lower-income zip codes tend to be more likely to

be assigned to lower-ad-load conditions.

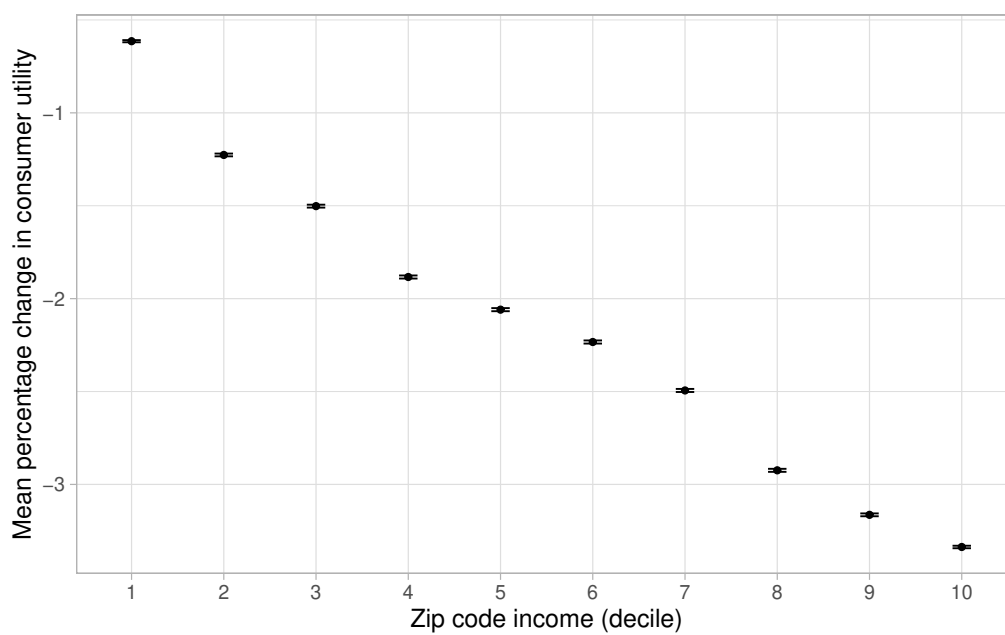
10 Discussion and conclusions

The advent of big data and large-scale data-processing technologies has allowed firms to optimize services, ads, and prices at the individual level. Although a large body of literature has focused on the implications of personalized pricing, the impact of personalizing the product itself is overlooked. The public perceives price discrimination to be unfair and companies have largely avoided such practices fearing a consumer backlash. Given these limitations, whether big data is going to be employed for personalizing pricing or versioning instead is still unclear. To the best of our knowledge, this study is the first empirical paper to investigate returns to personalized product versioning.

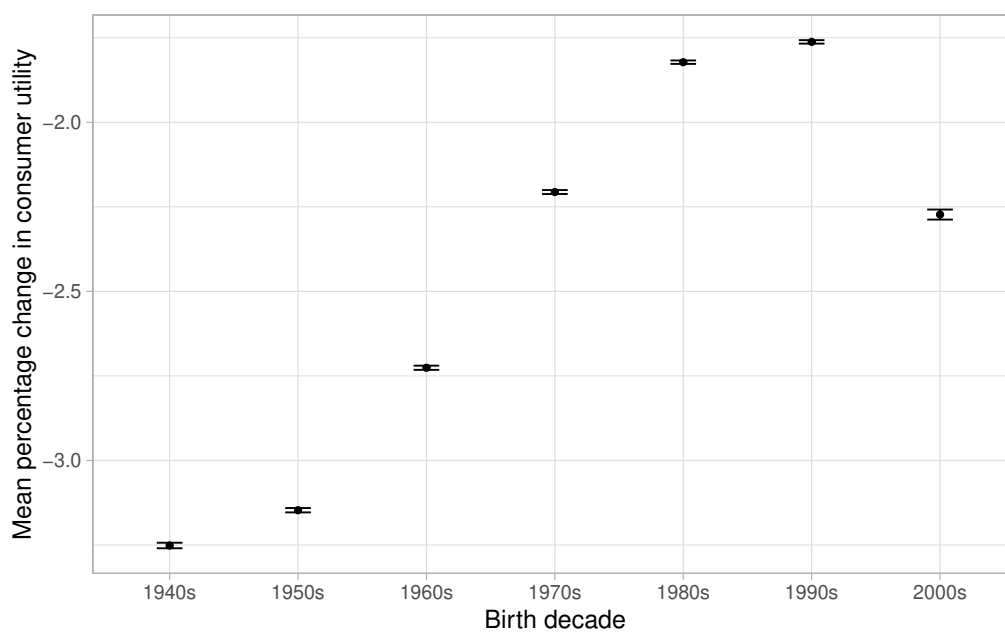
Although advertising can be used as an instrument to implement versioning for many content providers, including YouTube, Spotify, or Pandora, the idea of personalized versioning applies more broadly to other freemium business models. We provide two such examples here. First, in the online newspaper industry, the number of free pages or the amount of free content available to users can be used for versioning. Second, among cloud storage services, consider Dropbox's free plan, which offers 2GBs of free storage to users. However, users may be eligible for a wide variety of free storage promotions, including student discounts or offers available to users who purchase HP or Samsung devices (Martinez, 2014). We are not sure how targeted these strategies are, but the amount of free space offered to users on Dropbox is surely not uniform. We are not aware if companies have used these features to experiment with targeted versioning strategies, but as big data and experiments gain popularity, personalized versioning strategies may become an alternative to personalized pricing.

Our field experiments at Pandora present a unique opportunity to study the impact of personalized versioning in a product line. The availability of large-scale pre-treatment features allows us to segment users and prescribe personalized ad schedules. Our study highlights the importance of conducting field experiments and data-collection efforts for designing reliable prescriptive strategies. We also highlight challenges that take place in causal inference in two-sided platforms including partial control over realized outcomes or treatment exposure. The fact that Pandora has allowed us to share the details of their experiments and analyze the data to evaluate counter-factual strategies²¹ is, unfortunately, an

²¹The personalized versioning algorithm developed in this paper is not adopted by Pandora and we used inverse probability weighting to evaluate its performance using the randomization in the data.



(a)



(b)

Figure 13: The average percentage change in consumer utility. The impact of the policy across (a) different income levels and (b) users of different age.

exception in the industry, not a norm. We hope that efforts by firms such as Pandora, Yahoo, eBay, and Ziprecruiter (Lewis and Reiley, 2014; Blake, Nosko, and Tadelis, 2015; Dubé and Misra, 2017) promote transparency of firm-sponsored research.

Our results show that to achieve the same level of subscribers in the absence of a personalized ad-scheduling strategy, the firm needs to increase its ad load by more than 20%. This finding shows personalization can both improve firm profits and the average quality of service. We also find that gains from ad-load personalization materialize quickly. In particular, within three months of implementing the personalized counterpart of the control condition, the profits from subscriptions increase by 7%. Interestingly, the short-term impact of this strategy on the overall consumption of the ad-supported service is negligible. This finding, combined with switching costs between products, presents an opportunity for firms to investigate returns to dynamic optimization of implicit prices. Although some evidence shows firms change their quality of service in time due to demand seasonality (Lambrecht and Misra, 2017), we believe studying the trade-offs between personalized and time-varying quality of service adjustments is a fruitful area for future research. Finally, changing ad load could affect the click-through rate of ads or, in general, their effectiveness. This effect adds an additional layer of complexity for platforms that are compensated based on conversions or click-through rates. Although studying how ad effectiveness changes as a function of the number of ads in online platforms is beyond the scope of this paper, we acknowledge it could play an important role in the firm's decision to adopt the personalization algorithm discussed here.

References

- Aguiar, M., Hurst, E., and Karabarbounis, L. (2013). Time use during the great recession. *American Economic Review* 103 (5): 1664–96.
- Aguiar, M. A., Hurst, E., and Karabarbounis, L. (2011). Time use during recessions. Tech. rep. National Bureau of Economic Research.
- Ansari, A. and Mela, C. F. (2003). E-customization. *Journal of marketing research* 40 (2): 131–145.
- Becker, G. S. and Murphy, K. M. (1993). A simple theory of advertising as a good or bad. *The Quarterly Journal of Economics* 108 (4): 941–964.
- Bhatia, N., Moshary, S., and Tuchman, A. (2021). Investigating the Pink Tax: Evidence Against a Systematic Price Premium for Women in CPG. *Available at SSRN 3882214*.
- Big data and differential pricing / Executive Office of the President of the United States, Council of Economic Advisors. (2015).
- Blake, T., Nosko, C., and Tadelis, S. (2015). Consumer heterogeneity and paid search effectiveness: A large-scale field experiment. *Econometrica* 83 (1): 155–174.
- Chiang, R. and Spatt, C. S. (1982). Imperfect price discrimination and welfare. *The Review of Economic Studies* 49 (2): 155–181.
- Chiou, L. and Tucker, C. (2013). Paywalls and the demand for news. *Information Economics and Policy* 25 (2): 61–69.
- Clerides, S. K. (2002). Book value: intertemporal pricing and quality discrimination in the US market for books. *International Journal of Industrial Organization* 20 (10): 1385–1408.
- CNET (Jan. 2002). Now showing: random DVD prices on Amazon. URL: <https://www.cnet.com/news/now-showing-random-dvd-prices-on-amazon/>.
- Cowan, S. (2012). Third-Degree Price Discrimination and Consumer Surplus. *The Journal of Industrial Economics* 60 (2): 333–345.
- Crawford, G. S., Shcherbakov, O., and Shum, M. (2015). The welfare effects of endogenous quality choice in cable television markets.
- Crawford, G. S. and Shum, M. (2007). Monopoly quality degradation and regulation in cable television. *The Journal of Law and Economics* 50 (1): 181–219.
- DellaVigna, S. and Gentzkow, M. (2019). Uniform pricing in us retail chains. *The Quarterly Journal of Economics* 134 (4): 2011–2084.

- Deneckere, R. J. and Preston McAfee, R. (1996). Damaged goods. *Journal of Economics & Management Strategy* 5 (2): 149–174.
- Dubé, J.-P. and Misra, S. (2017). Scalable price targeting. Tech. rep. National Bureau of Economic Research.
- Edgecliffe-Johnson, A. (May 2009). Media wants to break free. URL: <https://www.ft.com/content/d0960f18-4303-11de-b793-00144feabdc0?fbclid=IwAR06SW3pqxG14RVQVKqjg02Nap-xF9gNlpVcycTuGRstKV>
- Farrell, M. H., Liang, T., and Misra, S. (2021). Deep neural networks for estimation and inference. *Econometrica* 89 (1): 181–213.
- Gentzkow, M. (2007). Valuing new goods in a model with complementarity: Online newspapers. *American Economic Review* 97 (3): 713–744.
- Godes, D., Ofek, E., and Sarvary, M. (2009). Content vs. advertising: The impact of competition on media firm strategy. *Marketing Science* 28 (1): 20–35.
- Goldstein, D. G., Suri, S., McAfee, R. P., Ekstrand-Abueg, M., and Diaz, F. (2014). The economic and cognitive costs of annoying display advertisements. *Journal of Marketing Research* 51 (6): 742–752.
- Halbheer, D., Stahl, F., Koenigsberg, O., and Lehmann, D. R. (2014). Choosing a digital content strategy: How much should be free? *International Journal of Research in Marketing* 31 (2): 192–206.
- Hauser, J. R., Urban, G. L., Liberali, G., and Braun, M. (2009). Website morphing. *Marketing Science* 28 (2): 202–223.
- Hitsch, G. J., Hortacsu, A., and Lin, X. (2019). Prices and promotions in us retail markets: Evidence from big data. Tech. rep. National Bureau of Economic Research.
- Hitsch, G. J. and Misra, S. (2018). Heterogeneous treatment effects and optimal targeting policy evaluation. *Available at SSRN 3111957*.
- Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association* 47 (260): 663–685.
- Huang, J., Reiley, D., and Riabov, N. (2018). Measuring Consumer Sensitivity to Audio Advertising: A Field Experiment on Pandora Internet Radio. *Available at SSRN 3166676*.
- Kaji, T., Manresa, E., and Pouliot, G. (2020). An adversarial approach to structural estimation. *arXiv preprint arXiv:2007.06169*.

- Kim, J., Park, Y., Kim, G., and Hwang, S. J. (2017). SplitNet: Learning to semantically split deep networks for parameter reduction and model parallelization. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org: 1866–1874.
- Krouse, S. and Vranica, S. (July 2022). Netflix partners with Microsoft for New Advertising-backed option. URL: <https://www.wsj.com/articles/netflix-partners-with-microsoft-to-launch-advertising-supported-plan-11657738975>.
- Kumar, S. and Sethi, S. P. (2009). Dynamic pricing and advertising for web content providers. *European Journal of Operational Research* 197 (3): 924–944.
- Lambrecht, A. and Misra, K. (2017). Fee or free: When should firms charge for online content? *Management Science* 63 (4): 1150–1165.
- Lewis, R. A. and Reiley, D. H. (2014). Online ads and offline sales: measuring the effect of retail advertising via a controlled experiment on Yahoo! *Quantitative Marketing and Economics* 12 (3): 235–266.
- Lin, S. (2020). Two-sided price discrimination by media platforms. *Marketing Science* 39 (2): 317–338.
- Martinez, J. (Oct. 2014). Acer and HP to preload Dropbox on all laptops, PCs. URL: <https://www.techradar.com/news/internet/cloud-services/acer-and-hp-to-preload-dropbox-on-all-laptops-pcs-1270862>.
- Maskin, E. and Riley, J. (1984). Monopoly with incomplete information. *The RAND Journal of Economics* 15 (2): 171–196.
- McManus, B. (2007). Nonlinear pricing in an oligopoly market: The case of specialty coffee. *The RAND Journal of Economics* 38 (2): 512–532.
- Murty, K. G. and Kabadi, S. N. (1985). Some NP-complete problems in quadratic and nonlinear programming. Tech. rep.
- Mussa, M. and Rosen, S. (1978). Monopoly and product quality. *Journal of Economic theory* 18 (2): 301–317.
- Pérez-Peña, R. and Arango, T. (Apr. 2009). They Pay for Cable, Music and Extra Bags. How About News? URL: <https://www.nytimes.com/2009/04/08/business/media/08pay.html>.
- Prasad, A., Mahajan, V., and Bronnenberg, B. (2003). Advertising versus pay-per-view in electronic media. *International Journal of Research in Marketing* 20 (1): 13–30.

- Rafieian, O. (2019). Optimizing user engagement through adaptive ad sequencing. Tech. rep. Technical report, Working paper.
- Rafieian, O. and Yoganarasimhan, H. (2021). Targeting and privacy in mobile advertising. *Marketing Science* 40 (2): 193–218.
- Rossi, P. E., McCulloch, R. E., and Allenby, G. M. (1996). The value of purchase history data in target marketing. *Marketing Science* 15 (4): 321–340.
- Sahni, N. S., Wheeler, S. C., and Chintagunta, P. (2018). Personalization in email marketing: The role of noninformative advertising content. *Marketing Science* 37 (2): 236–258.
- Salop, S. (1977). The noisy monopolist: Imperfect information, price dispersion and price discrimination. *The Review of Economic Studies* 44 (3): 393–406.
- Sato, S. (2019). Freemium as optimal menu pricing. *International Journal of Industrial Organization* 63: 480–510.
- Shalit, U., Johansson, F. D., and Sontag, D. (2017). Estimating individual treatment effect: generalization bounds and algorithms. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org: 3076–3085.
- Shapiro, C., Carl, S., Varian, H. R., et al. (1998). Information rules: a strategic guide to the network economy. Harvard Business Press.
- Shiller, B. R. et al. (2013). First degree price discrimination using big data. Brandeis Univ., Department of Economics.
- Simester, D., Timoshenko, A., and Zoumpoulis, S. I. (2020). Efficiently evaluating targeting policies: Improving on champion vs. challenger experiments. *Management Science* 66 (8): 3412–3424.
- Smith, A. N., Seiler, S., and Aggarwal, I. (2022). Optimal Price Targeting. *Marketing Science*.
- Spence, A. M. (1975). Monopoly, quality, and regulation. *The Bell Journal of Economics*: 417–429.
- Sriram, S., Chintagunta, P. K., and Manchanda, P. (2015). Service quality variability and termination behavior. *Management Science* 61 (11): 2739–2759.
- Tåg, J. (2009). Paying to remove advertisements. *Information Economics and Policy* 21 (4): 245–252.
- Valentino-DeVries, J., Singer-Vine, J., and Soltani, A. (Dec. 2012). Websites Vary Prices, Deals Based on Users’ Information. URL: <https://www.wsj.com/articles/SB100014241278873237772045781893918138815>
- Varian, H. R. (1985). Price discrimination and social welfare. *The American Economic Review* 75 (4): 870–875.

- Verboven, F. (2002). Quality-based price discrimination and tax incidence: evidence from gasoline and diesel cars. *RAND Journal of Economics*: 275–297.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association* 113 (523): 1228–1242.
- Wei, Y. and Jiang, Z. (2020). Estimating Parameters of Structural Models Using Neural Networks. *USC Marshall School of Business Research Paper*.
- Wilbur, K. C. (2008). A two-sided, empirical model of television advertising and viewing markets. *Marketing science* 27 (3): 356–378.
- Yoganarasimhan, H. (2020). Search personalization using machine learning. *Management Science* 66 (3): 1045–1070.
- Yoganarasimhan, H., Barzegary, E., and Pani, A. (2022). Design and evaluation of optimal free trials. *Management Science*.

Online appendix

A. conceptual model for user heterogeneity and versioning

In this section, we illustrate the benefits of personalized versioning with a set of examples. First, we consider a simple utility model that reflects the discrete choice between the outside options, consuming the ad-supported product, or using the paid subscription. Consider the following model:

$$u(z, p \mid \theta, \beta, \alpha) = \max \begin{cases} 0 & \text{outside option,} \\ \theta - \alpha z & \text{ad supported,} \\ \theta - \beta p & \text{paid version,} \end{cases} \quad (16)$$

where θ is the utility from consuming the product, z specifies the ad load²², and p is the subscription price. In the ad-supported condition, users are effectively paying with their time by listening to ads. Therefore, parameters α and β reflect how time and money are valued by users, that is users with higher/lower values of α and β are more/less sensitive to ads and prices, respectively. Also, let γ and c be the revenue per ad and the marginal cost of offering the service, respectively. Finally, let $\theta_\alpha = \frac{\theta}{\alpha}$ and $\theta_\beta = \frac{\theta}{\beta}$. If $\theta_\beta > z$ and $\frac{\theta_\beta}{\theta_\alpha} > \frac{z}{p}$, the user picks the ad-supported version. And if $\theta_\alpha > p$ and $\frac{\theta_\beta}{\theta_\alpha} < \frac{z}{p}$, the paid version is purchased; otherwise, the outside option is preferred. Figure 14 illustrates the decision regions for different types when price is set to p and ad load is equal to z .

²²In this simplified model, we assume users can consume the service in exchange for listening to z ads. In our case study, we account for the fact that the intensive margin of consumption (number of hours) could vary across users, and that possibility factors into the ad revenue.

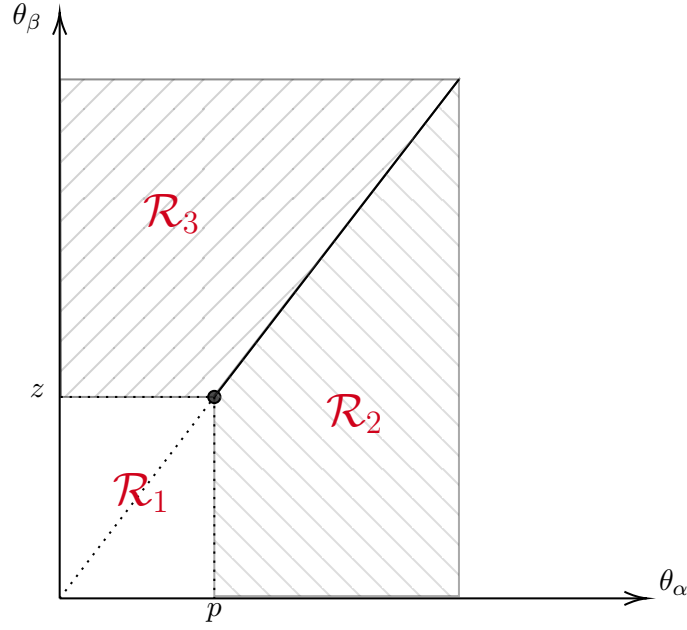


Figure 14: Decision regions for different types $(\theta_\alpha, \theta_\beta)$ for price vector (p, z) , where p and γz are assumed to be larger than c . Types that lie in \mathcal{R}_1 pick the outside option, types in \mathcal{R}_2 subscribe for the paid service, and those in \mathcal{R}_3 choose the ad-supported version.

A monopolist that can perfectly discriminate along both ad load and price dimensions will maximize its profits for each type $(\theta_\alpha, \theta_\beta)$. Note that a listener of type $(\theta_\alpha, \theta_\beta)$ has a willingness to consume at most $z = \theta_\alpha$ ads, and pay price $p = \theta_\beta$. For a listener with $\gamma\theta_\alpha > \max(\theta_\beta, c)$, the monopolist will only offer the ad-supported service with $(p, z) = (\infty, \theta_\alpha)$, and for those with $\theta_\beta > \max(\gamma\theta_\alpha, c)$, only the subscription service is offered with $(p, z) = (\theta_\beta, \infty)$, and when none of these conditions are satisfied serving the customer would not be worthwhile and $(p, z) = (\infty, \infty)$. Figure 15 demonstrates the decision regions for a monopolist based on the type of products sold, and the profits over regions \mathcal{R}_2 and \mathcal{R}_3 are equal to:

$$\Pi^* = \int_{\theta_\beta=c}^{\infty} \int_{\theta_\alpha=0}^{\frac{1}{\gamma}\theta_\beta} (\theta_\beta - c) f(\theta_\alpha, \theta_\beta) d\theta_\alpha d\theta_\beta + \int_{\theta_\alpha=\frac{c}{\gamma}}^{\infty} \int_{\theta_\beta=0}^{\gamma\theta_\alpha} (\gamma\theta_\alpha - c) f(\theta_\alpha, \theta_\beta) d\theta_\beta d\theta_\alpha, \quad (17)$$

where $f(\theta_\alpha, \theta_\beta)$ is the joint density of $(\theta_\alpha, \theta_\beta)$.

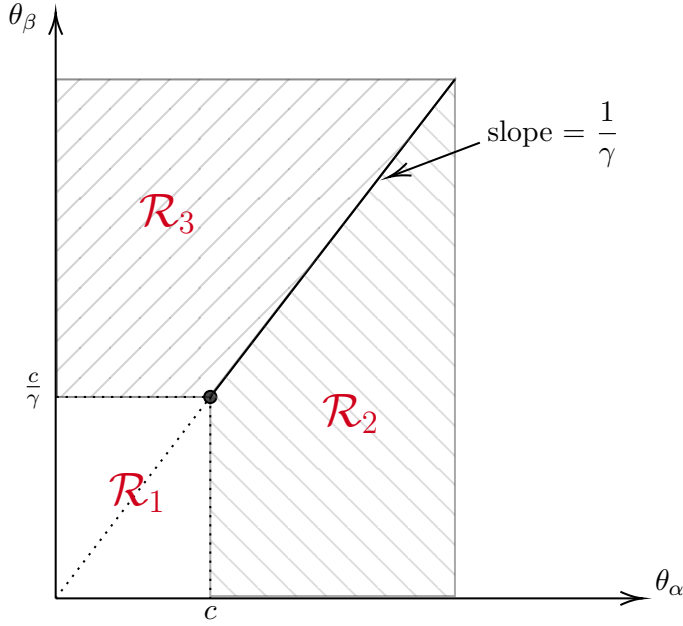


Figure 15: Decision regions for different types $(\theta_\alpha, \theta_\beta)$ for a price-ad load-discriminating monopolist. Serving customers in region \mathcal{R}_1 is not worthwhile; those in \mathcal{R}_2 will purchase the paid subscription and the rest will use the ad-supported service.

The results above demonstrate that when the monopolist has full information, he will only make one of the products available to each customer. However, these results rely on the crucial assumption that the monopolist can accurately observe the type of each listener $(\theta_\alpha, \theta_\beta)$. Our goal is to illustrate the benefits of offering a personalized menu of products when the seller has partial information about the type of consumers or randomness exists in choices made by the customers. We call this practice “personalized versioning,” which is akin to combining second- and third-degree price discrimination. The seller uses its information about the type of each customer to offer a personalized menu rather than only one product or price.

Let us consider a few simple examples to illustrate this idea:

- **Separable types:** Consider a monopolist (he) who needs to provide service to a customer (she). The monopolist knows that with probability ρ_1 , the customer is of type $(\theta_\alpha^{(1)}, \theta_\beta^{(1)})$, and with probability $\rho_2 = 1 - \rho_1$, she is of type $(\theta_\alpha^{(2)}, \theta_\beta^{(2)})$. If $\theta_\alpha^{(1)} > \theta_\alpha^{(2)}$ and $\theta_\beta^{(2)} > \theta_\beta^{(1)}$ (see Figure 16 for visual illustration), it is optimal for the monopolist to offer the paid version when the realized type is 1 and offer the ad-supported product when the realized type is 2. In this case, by offering a menu $(p, z) = (\theta_\alpha^{(1)}, \theta_\beta^{(2)})$, the firm can extract monopoly profits regardless of the type of customer. In particular, if the customer is of type 1, she purchases the paid version, whereas

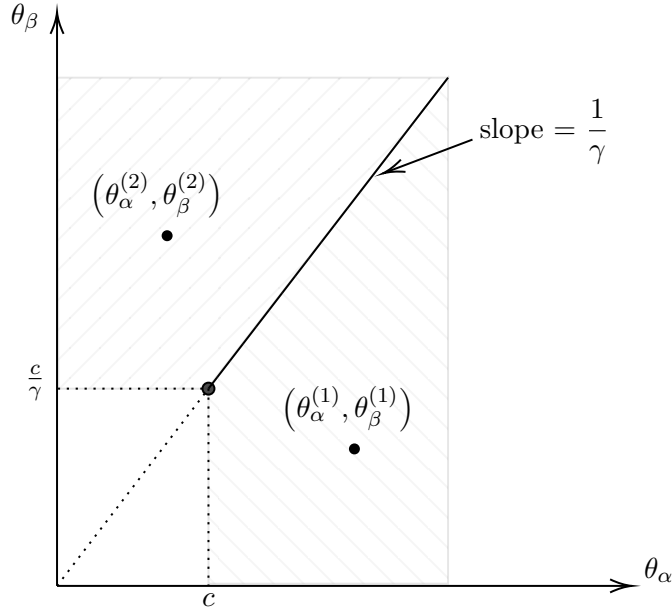


Figure 16: Perfectly separable condition. The seller is uncertain if the consumer is of type 1 or 2, but offering a menu with $(p, z) = (\theta_\alpha^{(1)}, \theta_\beta^{(2)})$ yields profits that are equal to the case where the seller has full information.

if she is of type 2, she uses the ad-supported product. In other words, in this case, using a menu can fully separate types from each other and resolves the uncertainty.

- **Inseparable types:** Let us now consider a more nuanced case. The customer can be of one of two types with probabilities ρ_1 and $\rho_2 = 1 - \rho_1$. This time $(\theta_\alpha^{(2)}, \theta_\beta^{(2)}) > (\theta_\alpha^{(1)}, \theta_\beta^{(1)})$; see Figure 17 for visual illustration. In this case, the optimal menu corresponds to one of the five red dots in Figure 17. Depending on the values of γ , c , ρ_1, ρ_2 , $(\theta_\alpha^{(1)}, \theta_\beta^{(1)})$, and $(\theta_\alpha^{(2)}, \theta_\beta^{(2)})$, either of these can be the optimal menu to offer. The only case where both products are offered is when $(p, z) = \left(\theta_\alpha^{(1)}, \theta_\alpha^{(1)} \frac{\theta_\beta^{(2)}}{\theta_\alpha^{(2)}} - \epsilon \right)$. Note that in this case, the existence of type 1 imposes a positive externality on the ad load, that is the quality of service when type 2 is realized. In other words, if the seller were certain the customer is of type 2, he would have the incentive to increase the ad load. However, in this case, because the customer is served under both conditions the ad load cannot be increased to more than $\theta_\alpha^{(1)} \frac{\theta_\beta^{(2)}}{\theta_\alpha^{(2)}}$ to make it incentive compatible for the type 2 customer to use the ad-supported service²³.

The examples above illustrate that the combination of personalization and versioning can make use of both the partial information that the firm may have and the customer's private information. We refer to this approach as personalized versioning. Note the application of personalized versioning

²³Recall from Figure 15 that if $\frac{\theta_\beta}{\theta_\alpha} > \frac{1}{\gamma}$, the seller is better off providing the ad-supported service.

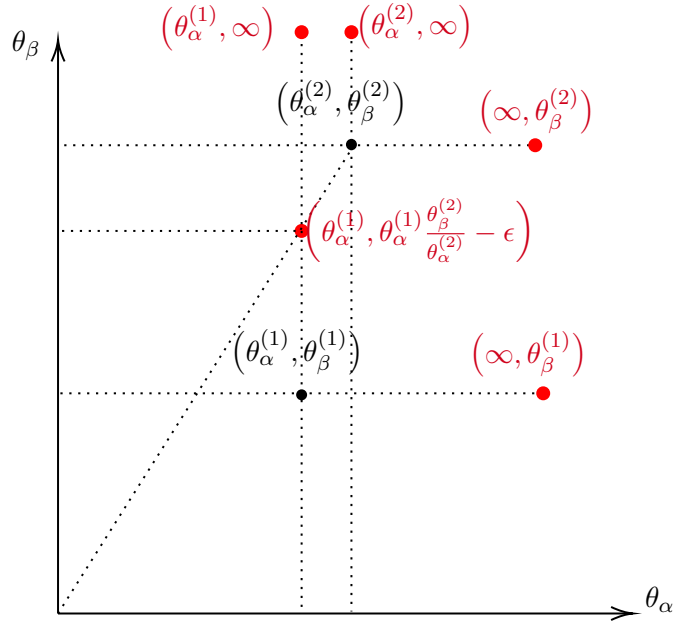


Figure 17: Inseparable condition. The uncertainty in the types cannot be perfectly dealt with by using a personalized menu. The optimal menu (p, z) is one of the five red dots, and depending upon the types and realization probabilities either one can be optimal.

is not limited to the case where uncertainty is present in parameter estimates, but also applies in random utility models even when parameter uncertainty is neglected. In section 3 we used a random utility model that corresponds to (16) to illustrate the trade-offs involved in personalizing ad load.

B. Ad delivery mechanism and partial control over realized ad load/capacity

In section 4, we discussed the field experiment and presented results on randomization checks, and realized changes in ad load across experiment cells. In this section, we discuss the realized effect on ad load and capacity in more detail.

Our experiment consists of seven cells: 3x1, 4x2, 6x2, 4x3, 5x3, 6x3, and control. The experiment shifts the timers in the ad delivery system. For instance, in the 3x1 condition, the timer is set to 20 minutes, and this number is equal to 15 minutes for the control, 4x2, and 4x3 conditions. The only difference between the 4x2 and the control condition is that the first ad pod²⁴ within each listening session in the control condition is constrained to be of length one. Note that a listener in the 6x3 condition ends up becoming eligible for an ad pod fewer than six times per hour (every 15 minutes), because the song endings do not perfectly align with the timers, see Figure 3. For example, the ad capacity, that is, the number of opportunities to show an ad per hour, in the 6x3 condition ends up being far less than $6 \times 3 = 18$. Figure 18 presents the realized ad capacity across different experiment conditions. Note that the ad capacity for the FxL condition is often much smaller than $F \cdot L$, since the song endings do not align with the timers. There is only one exception and that is the 3x1 condition. As highlighted in Figure 3, the listeners across all conditions become eligible for the first ad pod within the first five minutes of each listening session, therefore, the realized ad capacity may end up being larger than three in the 3x1 condition, especially for listeners that have short listening sessions.

When a listener (she) becomes eligible to receive an ad, the ad delivery system makes a request to fetch an ad for her. If she belongs to a demographic group that is attractive to advertisers it is easier to fetch ads for her. However, if the system runs out of ads to show to her those opportunities (ad capacity) are not filled. For instance, a listener may be eligible to receive eight ads in a given hour, but the system may end up fetching only five ads in the ad inventory that could be played for her. In our data, we record both the number of opportunities (ad capacity) and the realized number of ads (ad load) delivered in an hour for every user.

To summarize, pod length and frequency determine *ad capacity* rather than ad load. Figure 3 and Figure 18 illustrate that the realized ad capacity for the FxL condition ends up being far less than $F \cdot L$ as the song endings do not perfectly align with the timers. Furthermore, the realized number of ads shown to each listener (ad load) also depends on advertisers' demand. As we move to more extreme

²⁴The listeners across all experiment conditions become eligible for the first ad pod after the first five minutes of each listening session.

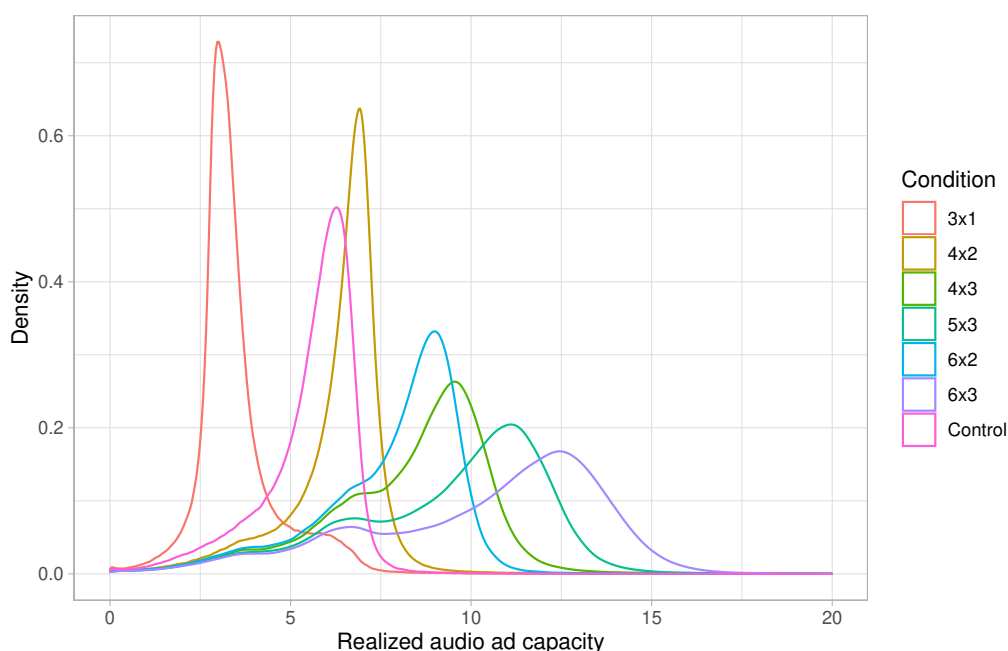


Figure 18: Realized audio ad capacity across different treatment arms

conditions such as the 6x3 condition it becomes more difficult to shift both ad capacity and ad load, see Table 3. Figure 19 depicts the density of realized ad load for users in different treatment cells. Although higher-ad-capacity conditions have a higher realized ad load, the distribution becomes more dispersed as the capacity increases, compare Figures 18 and 19. This finding is indicative of the fact that filling higher capacities for users tends to be more difficult because running out of ads to serve in the higher-capacity conditions is more probable. Table 3 in the body illustrates these findings by reporting the average realized ad load, capacity, and fill rate during the first year of the experiment.

To further demonstrate the partial control problem, we plot the lift in ad load between the control and 6x3 condition across different consumer groups based on their pre-treatment ad load in Figure 20. The figure shows the increase in ad load in the 6x3 condition relative to the control condition is not uniform across different listener groups. The lift in ad load is more pronounced for consumers who received more ads in the pre-treatment period. This heterogeneity in the lift reflects the role of advertisers' demand in the realized ad load and shows that the additional capacity is more likely to be filled for those consumers who are more attractive to advertisers. This demonstrates the fact that firms need to account for the discrepancy between the *intended* and *realized* change in the implicit price, which leads to an additional layer of complexity relative to the traditional pricing problems.

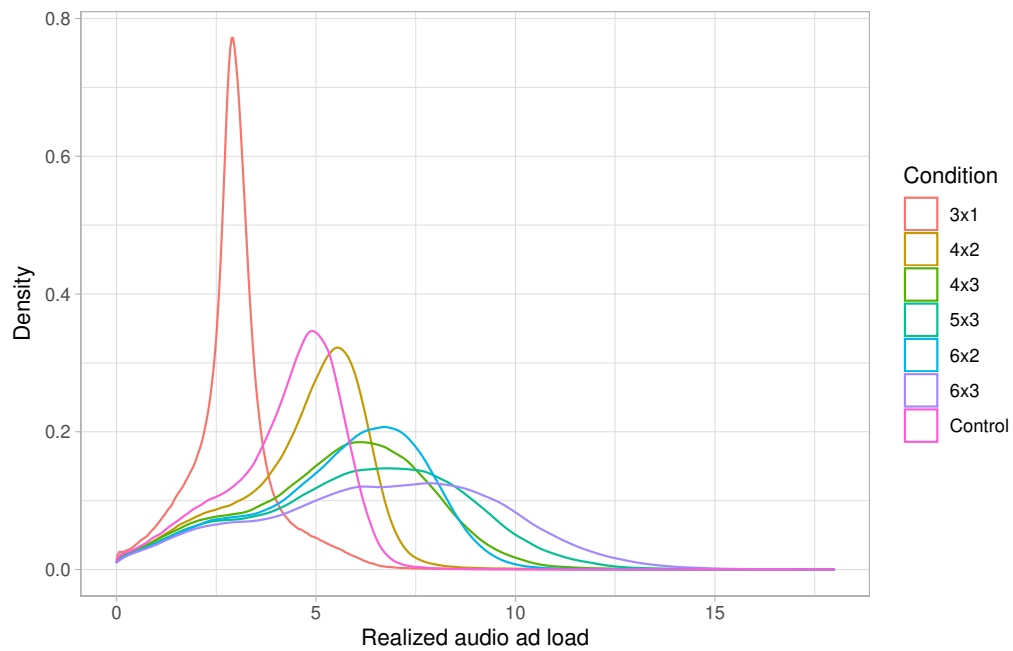


Figure 19: Realized ad load across different treatment arms

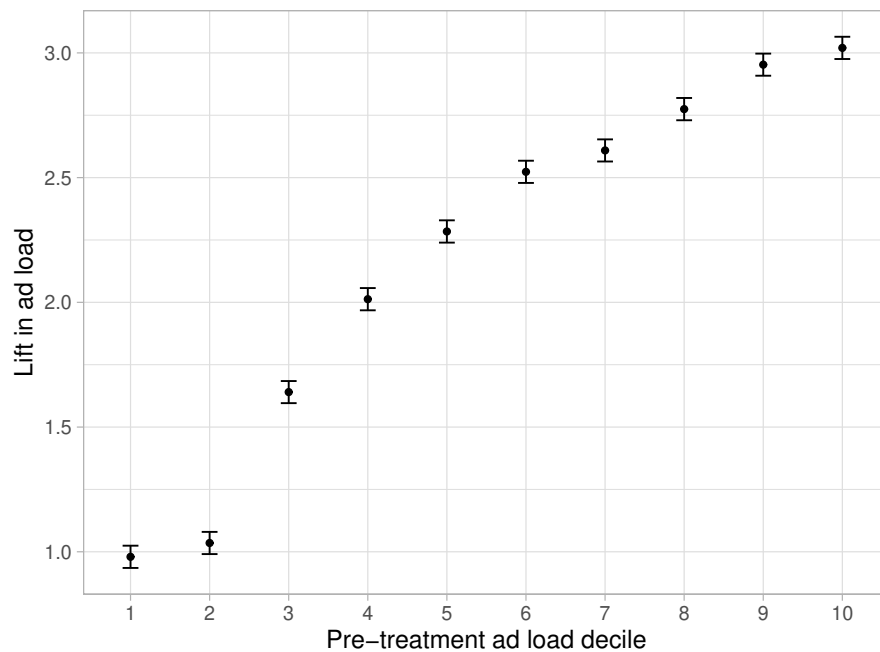


Figure 20: Realized lift in ad load (ads/hours) in the 6x3 condition relative to control as a function of pre-treatment ad load. Note the lift in ad load could vary drastically across different groups, due to differences in the attractiveness of different segments for advertisers.

C. Reduced-form regressions

In Figure 5 of section 5, we illustrated the average treatment effect of changing ad load on consumption and subscription status using a series of instrumental variable regressions. In this section, we present the reduced-form regressions that correspond to those IV regressions. In particular, we regress the normalized outcomes (3) directly on the treatment dummies using the following specification:

$$\tilde{\mathbf{Y}}_i = \alpha + \sum_{j=1}^6 \beta_j \cdot \mathbb{1}_{\{\tau=e_j\}} + \epsilon_i, \quad (18)$$

where i indexes listeners, and $\tilde{\mathbf{Y}}_i$ is a normalized outcome of interest, i.e., listening hours, activity dummy, and plus subscription dummy, for each user in a given week. We present the results of these analyses for the last week of each quarter after the experiment kicked off in Tables 4-9. Similar to our pooled IV results in Figure 5, our results here show that the treatment effect of ad load changes remains stable after Q4 2016, however, the impact on consumption takes longer than a year to stabilize.

Table 4: Activity and subscription status in the last week of Q3-2016 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.218)	100.000*** (0.234)	100.000*** (0.110)	100.000*** (0.809)
3x1	0.288 (0.308)	0.814** (0.330)	0.364** (0.156)	-4.188*** (1.143)
4x2	-0.375 (0.267)	-0.522* (0.286)	-0.106 (0.135)	3.913*** (0.990)
6x2	-1.144*** (0.377)	-1.546*** (0.405)	-0.413** (0.191)	6.959*** (1.400)
4x3	-0.878** (0.377)	-1.185*** (0.405)	-0.430** (0.191)	8.956*** (1.401)
5x3	-0.862** (0.377)	-1.794*** (0.404)	-0.471** (0.190)	12.886*** (1.399)
6x3	-1.415*** (0.453)	-2.169*** (0.486)	-0.536** (0.229)	14.129*** (1.682)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

Table 5: Activity and subscription status in the last week of Q4-2016 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.219)	100.000*** (0.238)	100.000*** (0.110)	100.000*** (0.706)
3x1	1.003*** (0.309)	1.544*** (0.336)	0.462*** (0.155)	-5.418*** (0.998)
4x2	-0.615** (0.268)	-0.785*** (0.291)	-0.210 (0.134)	2.293*** (0.865)
6x2	-1.499*** (0.379)	-2.196*** (0.412)	-0.447** (0.190)	8.126*** (1.223)
4x3	-1.410*** (0.379)	-2.163*** (0.412)	-0.564*** (0.190)	11.858*** (1.223)
5x3	-1.834*** (0.379)	-3.156*** (0.412)	-0.843*** (0.190)	14.945*** (1.222)
6x3	-2.478*** (0.455)	-3.888*** (0.495)	-1.051*** (0.228)	17.941*** (1.469)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

Table 6: Activity and subscription status in the last week of Q1-2017 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.215)	100.000*** (0.233)	100.000*** (0.109)	100.000*** (0.591)
3x1	1.998*** (0.303)	2.946*** (0.329)	0.644*** (0.154)	-5.581*** (0.835)
4x2	-0.468* (0.263)	-0.691** (0.285)	-0.460*** (0.134)	2.738*** (0.723)
6x2	-1.794*** (0.372)	-2.512*** (0.403)	-0.993*** (0.189)	8.267*** (1.023)
4x3	-1.861*** (0.372)	-2.955*** (0.403)	-1.233*** (0.189)	10.620*** (1.023)
5x3	-2.617*** (0.372)	-4.516*** (0.403)	-1.419*** (0.189)	14.838*** (1.022)
6x3	-3.156*** (0.447)	-5.027*** (0.484)	-1.867*** (0.227)	17.215*** (1.229)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

Table 7: Activity and subscription status in the last week of Q2-2017 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.219)	100.000*** (0.243)	100.000*** (0.112)	100.000*** (0.548)
3x1	2.134*** (0.310)	3.401*** (0.344)	0.678*** (0.158)	-5.615*** (0.774)
4x2	-0.215 (0.268)	-0.334 (0.298)	-0.345** (0.137)	2.795*** (0.671)
6x2	-1.928*** (0.379)	-3.040*** (0.421)	-1.077*** (0.194)	9.398*** (0.949)
4x3	-2.031*** (0.380)	-3.096*** (0.421)	-1.475*** (0.194)	10.937*** (0.949)
5x3	-2.939*** (0.379)	-4.673*** (0.421)	-1.764*** (0.194)	15.182*** (0.948)
6x3	-3.759*** (0.456)	-6.054*** (0.506)	-2.267*** (0.233)	16.315*** (1.140)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

Table 8: Activity and subscription status in the last week of Q3-2017 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.225)	100.000*** (0.249)	100.000*** (0.114)	100.000*** (0.566)
3x1	2.179*** (0.318)	3.392*** (0.352)	0.554*** (0.161)	-5.457*** (0.800)
4x2	-0.532* (0.276)	-0.835*** (0.305)	-0.759*** (0.140)	3.955*** (0.693)
6x2	-1.510*** (0.390)	-2.590*** (0.431)	-1.298*** (0.197)	10.684*** (0.981)
4x3	-1.758*** (0.390)	-3.082*** (0.431)	-1.528*** (0.198)	12.532*** (0.981)
5x3	-2.814*** (0.390)	-5.043*** (0.431)	-2.185*** (0.197)	16.985*** (0.980)
6x3	-3.527*** (0.469)	-6.174*** (0.518)	-2.349*** (0.237)	18.201*** (1.178)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

Table 9: Activity and subscription status in the last week of Q4-2017 across different treatment arms relative to control.

	<i>Dependent variable:</i>			
	All hours	Ad-supported hours	Active user	Subscription rate
	(1)	(2)	(3)	(4)
Control	100.000*** (0.234)	100.000*** (0.262)	100.000*** (0.119)	100.000*** (0.575)
3x1	2.114*** (0.331)	3.173*** (0.370)	0.897*** (0.168)	-5.802*** (0.813)
4x2	-0.569** (0.287)	-1.238*** (0.321)	-0.635*** (0.145)	4.100*** (0.705)
6x2	-1.411*** (0.405)	-2.491*** (0.454)	-1.072*** (0.206)	10.695*** (0.997)
4x3	-1.761*** (0.406)	-3.447*** (0.454)	-1.738*** (0.206)	13.386*** (0.997)
5x3	-2.330*** (0.405)	-5.082*** (0.453)	-2.166*** (0.205)	17.077*** (0.996)
6x3	-3.604*** (0.487)	-6.424*** (0.545)	-2.369*** (0.247)	18.932*** (1.197)
Observations	7,350,278	7,350,278	7,350,278	7,350,278
<i>Note:</i>			*p<0.1; **p<0.05; ***p<0.01	

D. Persistence in heterogeneous treatment effects

In section 7 we showed that our models are able to sort users based on the magnitude of the treatment effect on ad and subscription revenues using data from December 2016. In this section, we illustrate that the heterogeneous treatment effects detected in this time period are persistent. In particular, we show that if one sorts users in the hold-out sample based on the predicted treatment effect on subscription and ad revenues, the ordering does not only explain the lift in December 2016, see Figure 8, but it also has explanation power in December 2017.

In section 7 we divided the listeners in the hold-out sample into five quintiles based on the predicted treatment effects of the 6x3 relative to the control condition for both subscriptions and ads. We then demonstrated that the *realized lift* in subscriptions and ads were indeed correlated with predicted treatment effects in Figure 21. We replicated this analysis but instead of using December 2016 data for calculating we used data from a year later, that is during December 2017. We present these results in Figure 21. Note that the magnitude of the effects could be slightly different from those in Figure 8, but the ordering remains consistent. This means that the ordering of users based on treatment effects predictions using 2016 data is still valid a year later in 2017. This is very interesting and shows that it is possible to use the short-run 6 months data to detect heterogeneous treatment effects that are persistent.

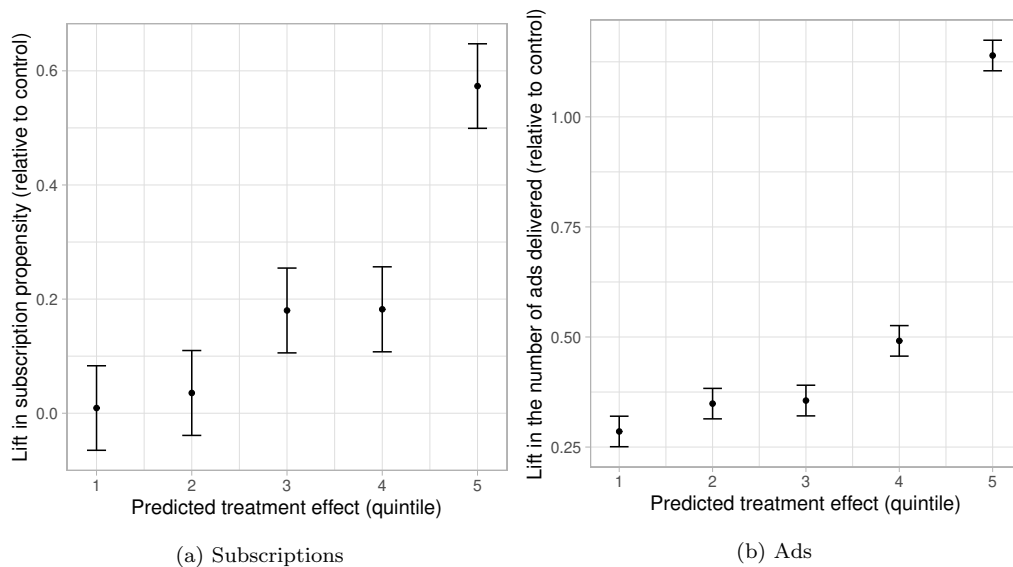


Figure 21: Realized lift in subscription propensity and ads delivered in the 6x3 condition relative to control in December of 2017 on the hold-out sample as a function of predicted treatment effect quintile using December 2016 as training.

E. A direct estimation model for ads delivered

In section 7, we combined three models to predict the number of ads played to each listener under different treatment conditions (counterfactual states). A model that predicts the probability of being an active ad-supported listener $P_a(\mathbf{x}, \boldsymbol{\tau})$, a model that predicts the number of listening hours conditional on being active and ad-supported $C(\mathbf{x}, \boldsymbol{\tau})$, and finally a model that predicts the ad load conditional on being active and ad-supported $A(\mathbf{x}, \boldsymbol{\tau})$ across different conditions. The predicted number of ads for an individual with pre-treatment features \mathbf{x} and in treatment condition, $\boldsymbol{\tau}$ was calculated as:

$$\text{predicted \# of ads delivered} = P_a(\mathbf{x}, \boldsymbol{\tau}) \cdot C(\mathbf{x}, \boldsymbol{\tau}) \cdot A(\mathbf{x}, \boldsymbol{\tau}). \quad (19)$$

Since models $C(\mathbf{x}, \boldsymbol{\tau})$ and $A(\mathbf{x}, \boldsymbol{\tau})$ are trained on listeners conditional on being ad-supported and active one might be concerned about selection. However, note that we only use these models to solve for the allocation policy (12), and we *do not* use the predictions of the models to evaluate the realized number of ads or subscriptions in the counterfactuals. In particular, the performance of the model in terms of lift in subscriptions and delivered ads are calculated using inverse propensity weighted estimates in the hold-out sample, which is similar to the approach used by Hitsch and Misra (2018), Rafieian and Yoganarasimhan (2021), Simester, Timoshenko, and Zoumpoulis (2020), and Yoganarasimhan, Barzegary, and Pani (2022). This means that even if the errors are correlated across the first model that learns the probability of being ad-supported and the models used for predicting listening hours and ad load, this correlation does not affect the counterfactual estimates of model performance, e.g., Figures 9-10, because those results are based on inverse propensity weighted estimates and do not depend on the predicted values themselves. Nevertheless, here we show that we obtain similar performance if we use a model that directly predicts the number of ads delivered instead of the approach used in (19).

Let \mathbf{Y}_i be the number of ads played to listener i during the training period (December of 2016). Note $\mathbf{Y}_i = 0$ for paid listeners or inactive users in that time period. We model \mathbf{Y}_i as follows:

$$\mathbf{Y}_i = \mathbb{1}_{f(\mathbf{x}_i, \boldsymbol{\tau}_i) > 0} \cdot f(\mathbf{x}_i, \boldsymbol{\tau}_i), \quad (20)$$

where $f(\mathbf{x}, \boldsymbol{\tau})$ is parameterized as a neural network. A schematic view of the neural network's structure is presented in Figure 22. We use a weighted objective similar to (6) with mean-squared error as the

loss to train the model. Let $\mathcal{A}(\mathbf{x}, \boldsymbol{\tau})$ be the model that predicts the number of ads served to each individual. Using this model instead to predict the number of ads served to each individual transforms the optimization problem (13) into:

$$\sum_i \underset{\boldsymbol{\tau}_i}{\text{maximize}} \underbrace{m_s \cdot P_s(\mathbf{x}_i, \boldsymbol{\tau}_i) + \lambda \cdot \mathcal{A}(\mathbf{x}_i, \boldsymbol{\tau}_i)}_{f(\mathbf{x}_i, \lambda, \boldsymbol{\tau}_i)}. \quad (21)$$

Similar to our approach in section 8, we shift λ in equation (13) to obtain different personalized policies albeit using the joint model $\mathcal{A}(\mathbf{x}_i, \boldsymbol{\tau}_i)$ for predicting the number of ads served. The counterpart of the Pareto frontier in Figure 9 for the new set of personalized policies using the joint model is presented in Figure 23. The performance remains similar to our approach in section 8 and the personalized counterpart of control achieves about a 7% gain in subscription profits relative to control.

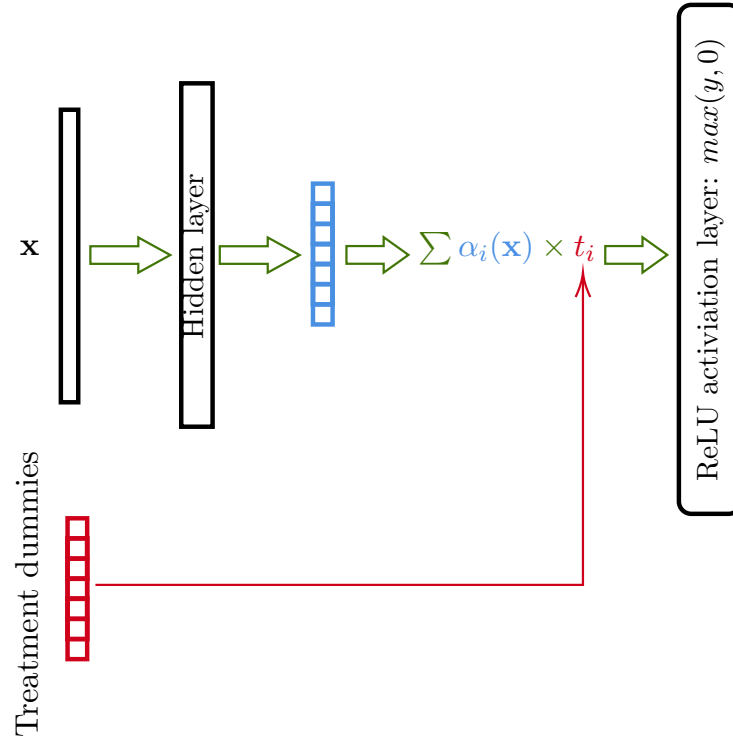


Figure 22: A schematic view of the neural network architecture used for estimating the realized number of ads delivered across experimental conditions.

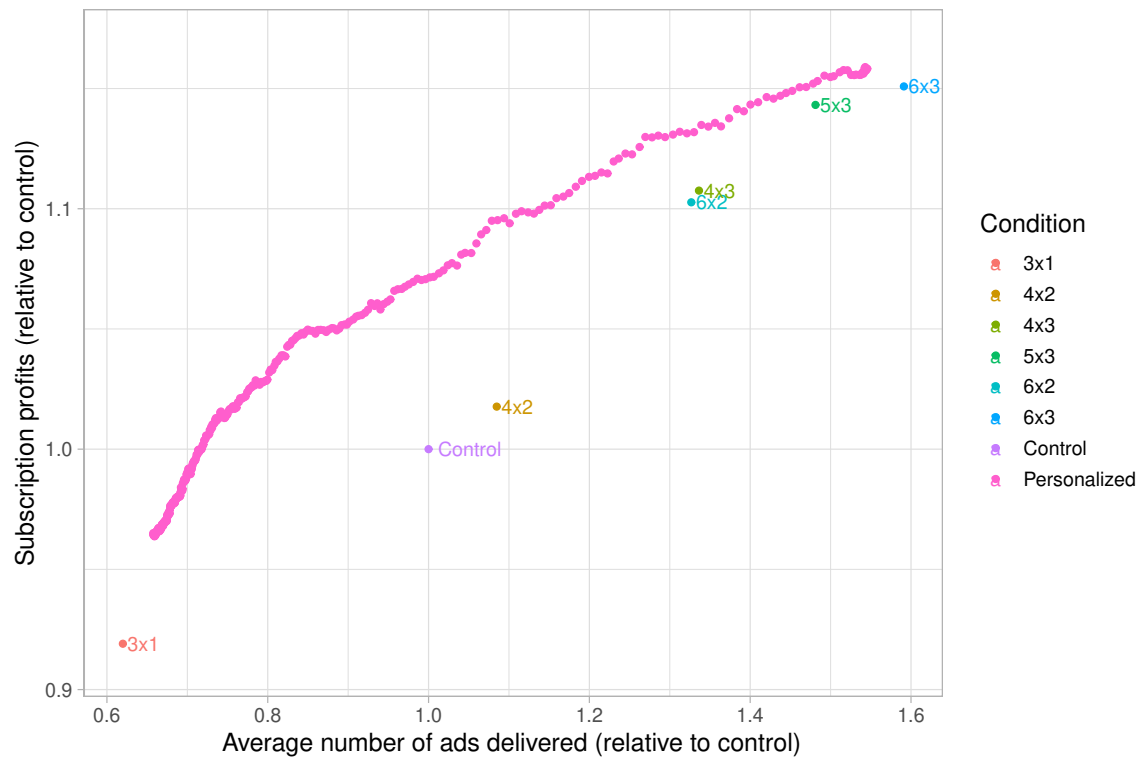


Figure 23: Change in subscription profits as a function of the number of ads served using a joint prediction model for the number of ads served. The performance of the personalized policies is similar to Figure 9.