



A Large-Scale Constrained Joint Modeling Approach for Predicting User Activity, Engagement, and Churn With Application to Freemium Mobile Games

Trambak Banerjee^a, Gourab Mukherjee^a, Shantanu Dutta^b, and Pulak Ghosh^c

^aDepartment of Data Sciences and Operations, University of Southern California, Los Angeles, CA; ^bDepartment of Marketing, University of Southern California, Los Angeles, CA; ^cDepartment of Decision Sciences and Information Systems, Indian Institute of Management Bangalore, Bangalore, India

ABSTRACT

We develop a *constrained extremely zero inflated joint* (CEZIJ) modeling framework for simultaneously analyzing player activity, engagement, and dropouts (churns) in app-based mobile freemium games. Our proposed framework addresses the complex interdependencies between a player's decision to use a freemium product, the extent of her direct and indirect engagement with the product and her decision to permanently drop its usage. CEZIJ extends the existing class of joint models for longitudinal and survival data in several ways. It not only accommodates extremely zero-inflated responses in a joint model setting but also incorporates domain-specific, convex structural constraints on the model parameters. Longitudinal data from app-based mobile games usually exhibit a large set of potential predictors and choosing the relevant set of predictors is highly desirable for various purposes including improved predictability. To achieve this goal, CEZIJ conducts simultaneous, coordinated selection of fixed and random effects in high-dimensional penalized generalized linear mixed models. For analyzing such large-scale datasets, variable selection and estimation are conducted via a distributed computing based split-and-conquer approach that massively increases scalability and provides better predictive performance over competing predictive methods. Our results reveal codependencies between varied player characteristics that promote player activity and engagement. Furthermore, the predicted churn probabilities exhibit idiosyncratic clusters of player profiles over time based on which marketers and game managers can segment the playing population for improved monetization of app-based freemium games. Supplementary materials for this article, including a standardized description of the materials available for reproducing the work, are available as an online supplement.

ARTICLE HISTORY

Received May 2018
Revised March 2019

KEYWORDS

CEZIJ; Constrained joint modeling; Freemium behavior; Large-scale longitudinal data analysis; Mobile apps; Zero-inflation.



1. Introduction

Mobile games have become an integral part of modern life (Koetsier 2015). While their almost ubiquitous presence is increasingly reshaping the recreational, socialization, educational and learning media (Statista 2018; Hwong 2016, chaps. 1 and 3; Garg and Telang 2013), the monetization policies associated with these new mobile apps is rapidly revolutionizing the digital marketing and advertisement space in information systems (Appel et al. 2019; Liu, Au, and Choi 2014). As such mobile games (as per industry standards formally defined as any app-based game played on an Internet enabled mobile device such as tablets, phones, etc.) currently comprise 42% of the market share of global gaming products (McDonald 2017) and more than 800,000 mobile games were available for download in the iOS App Store alone, with approximately 400 new submissions arriving each day (Pocket Gamer 2018). To understand how quickly the gaming market is growing, a new industry study from Spil Games (Diele 2013) reports that 1.2 billion people are now playing games worldwide, with 700 million of those online. The unprecedented growth and popularity of mobile games has resulted in a market with some very unique consumer characteristics (Boudreau, Jeppesen, and


Miric 2019). It is an extremely crowded market with significant proportion of revenue accumulated through advertisement based on free products (Appel et al. 2019). Specifically, app retention rates are much lower than the observed retention rates in classical products and services, with reports suggesting that more than 80% of all app users churn (dropout) within the first quarter (Perro 2016; MarketingCharts 2017). The freemium business model (Niculescu and Wu 2011), which offers a certain level of service without cost and sells premium add-on components to generate revenue, is a popular strategy for monetization of these mobile games. As such, industry reports indicate that more than 90% of the mobile games start as free, and more than 90% of the profits currently come from products that began as free (AppBrain 2017; Taube 2013). User characteristics in freemium models differ in fundamental aspects from traditional marketing models. This necessitates development of new analytical methods for modeling freemium behavior.


1.1. Freemium Model: Player Activity and Engagement

In the freemium market, firms initially attract customers with free usage of their products, with the expectation that free

CONTACT Gourab Mukherjee  gmukherj@marshall.usc.edu  Department of Data Sciences and Operations, University of Southern California, 3670 Trousdale Pkwy, Los Angeles, CA 90089-0809.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/r/JASA.

 These materials were reviewed for reproducibility.

 Supplementary materials for this article are available online. Please go to www.tandfonline.com/r/JASA.

© 2019 American Statistical Association

usage will lead customers to engage in future purchase of premium components. However, customers can always remain free users and never need to enjoy the premium components of the product. This is an important distinction with non-freemium business models, where customers *must* purchase to use the product. While the free to use part of freemium products helps to attract the consumer base quickly (Kumar 2014), managers are uncertain on whether and how freemium can generate profits (Needleman and Loten 2012) as majority of the consumers do not use the premium part of freemium. As such, unless a game is very popular, in-app purchases (IAP) contribute an insignificant proportion of its revenue. Mobile marketing automation firm Swrve (Swrve 2016) found that over 48% of all in-game revenue are derived from 0.19% of all players, which is a tiny segment. While in-game (direct) revenue is important, there are several indirect ways of monetizing the free users by involving them to *engage* with the game via social media (through facebook or twitter likes and posts of game achievements, inviting social media friends to join game, watching, liking or posting youtube videos related to the game) or the app center. To measure the *daily engagement* of a player, we judiciously combine her IAP (direct source of revenue) with her varied involvements with the game in media (indirect source of monetization), under the notion that purchase is the highest form of engagement. We define a player's *daily activity* as the time she spends playing the game in the day. Positive daily activity does not always lead to positive engagement. It is commonly believed that as a game grows with increasing and prolonged player *activities*, it will have more positive as well as higher engagement values.

For game managers it is extremely important to accurately measure player *activity*, *engagement*, and their codependencies. Also, varied retention strategies are often used to curb high churn rates and their effects need to be properly analyzed. Here, we develop a *constrained extremely zero inflated joint* (CEZIJ) modeling framework that provides a disciplined statistical program for jointly modeling player activity, engagement, and churn in online gaming platforms. Our proposed framework captures the codependencies between usage (activity), direct and indirect revenue (engagement), and dropouts (which is a time-to-event) and provides a systematic understanding of how the dependent variables influence each other and are influenced by the covariates. Furthermore, the CEZIJ framework can be used to predict the activity, engagement, and attrition of new players. The ability to forecast behavior of new players is critical for managers, as this enables them to better predict the effectiveness of their gaming platform in engaging customers and thus attract future advertisers to their platform.

1.2. Joint Modeling of Player Characteristics

Our joint modeling framework uses generalized linear mixed effect models (GLMM) and relies on a joint system of equations that model the relationships between activity, engagement, and churn. In the activity equations, we separately assess whether consumers are active (i.e., play the game) and the extent of their activity through the amount of time they spend playing

the game. Engagement is modeled by the probability of having positive engagement and by a conditional model on the positive engagement values. In the churn equations, we account for permanent churn identified as those players who are not active for more than 30 consecutive days. Our modeling systems addresses the complex interdependencies between (1) the decision to use the free product, (2) how much time will be spent using the free product, (3) the decision to engage, (4) the extent of engagement, and (5) the decision to churn. That is, the joint equation system comprehensively uncovers positive, negative, or zero codependencies among activity, engagement, and churn in freemium markets. In recent times, joint modeling of multiple outcomes have received considerable attention (Rizopoulos 2012). Many applications consider the modeling of single or multiple longitudinal outcomes and a time-to-event outcome (e.g., Jiang 2007; McCulloch 2008; Rizopoulos, Verbeke, and Lesaffre 2009; Rizopoulos, Verbeke, and Molenberghs 2010; Banerjee, Carlin, and Gelfand 2014; Rizopoulos and Lesaffre 2014). Our motivation for jointly modeling the drivers of player gaming traits and dropout arises from the fact that there is heterogeneity across player's outcomes and one must combine these effects by correlating the multiple responses. Since these responses are measured on a variety of different scales (viz. time spend in hours, revenue in dollars), a flexible solution is to model the association between different responses by correlating the random heterogeneous effects from each of the responses. Such a model not only gives us a covariance structure to assess the strength of association between the responses, but also offers useful insights to managers, since despite huge popularity of mobile games among users, managers are not certain whether freemium is profitable. Furthermore, it is important for managers to understand how activity and engagement are related to player churn. While customers who frequently use the free product could be more satisfied, thus reducing their probability of churn (Gustafsson, Johnson, and Roos 2005), free usage could be related to a greater probability of churn as there is little switching cost for customers due to their lower perceived value (Yang and Peterson 2004). Earlier studies used simpler models for churn that are independent of the purchase rate (Jerath, Fader, and Hardie 2011). Here we model churn allowing for possible codependencies with activity and engagement.

1.3. Statistical Challenges

The online gaming data, which is the application case described in detail in Section 2, and the particular business model of freemium, pose several statistical challenges and necessitates novel extensions of the joint modeling framework. We describe the details below.

(i) *Extreme zero-inflation*: Freemium behavior suggests that even if a player is active on a day, it very rarely leads to purchases or social media engagement on her part. Thus, though both activity and engagement are zero-inflated, engagement has an extremely zero-inflated distribution. Mixture distributions of which zero-inflated distributions are a special case are commonly used in this kind of data. While there are multiple models that have been developed to accommodate data with excess zeros (see, e.g., Olsen and Schafer 2001; Min and Agresti

2005; Han and Kronmal 2006; Alfò, Maruotti, and Trovato 2011; Greene 2009, and the references therein), there is not much attention on extreme zero-inflated data. Few recent works, for example, Hatfield et al. (2012) show promise though. We develop a joint modeling framework that can accommodate extreme zero-inflation. The proposed framework allows us to accommodate large incidences of no-engagement by active players, such as that observed in freemium markets and helps managers more accurately forecast sales potential for businesses with large active customer bases but small incidence of engagement by separating the confound between non-active and non-engaged. We highlight that this extreme zero inflated data are not only relevant to freemium markets but are also common in other businesses wherein a sizable portion of the active consumer base engages in very little purchase activity. For example, in the *online* setting, we may observe low incidences of online ratings (i.e., 1–5 star rating), user generated content creation, banner ad click-through, and search ad conversion (see, e.g., Urban et al. 2013; Haans, Raassens, and van Hout 2013). Likewise in the *offline* setting of purchase data, for example, most product categories comprise less than 5% planned or actual purchase for an individual's visit to the grocery store (Hui et al. 2013). Thus, if managers are interested in assessing promotion on sales or individual level purchase activity in these contexts, we may be confronted with data that contains an extreme number of zeros.

(ii) *Parametric constraints*: We develop a framework for incorporating domain specific structural constraints in our model for one may have prior knowledge that a vector of parameters lies on a simplex or follows a particular set of inequality constraints. It is quite common in gaming data to have prior information available on various activities of the player. For example, it is well-known that player characteristics will have a burgeoning weekend effect or marketers have prior knowledge on the comparative efficacies of the retention strategies particularly if they have known dosage demarcations. Using these side information is extremely important (James, Paulson, and Rusmevichientong 2013; Banerjee, Mukherjee, and Sun 2018) and the CEZIJ framework incorporates these domain expertise though convexity constraints in our model.

(iii) *Hierarchical variable selection*: In online gaming data one usually encounters numerous covariates related to both game specific and player specific variables and choosing the relevant set of covariates is highly desirable for improving predictability. It is also important that the inferential problems associated with these data properly account for the presence of a lot of possibly spurious covariates. The high-dimensionality of these datasets, however, renders classical variable selection techniques incompetent. We develop a novel algorithm for estimation in the CEZIJ framework that conducts variable selection from a large set of potential predictors in GLMM based joint model. To produce interpretable effects CEZIJ imposes a hierarchical structure on the selection mechanism and includes covariates either as fixed effects or composite effects where the latter are those covariates that have both fixed and random effects (Hui, Müller, and Welsh 2017a) (see Section 4 for details). Efficient selection of fixed and random effect components in a mixed model framework has received considerable attention in recent years (Bondell, Krishna, and Ghosh 2010; Fan and Li 2012;

Lin, Pang, and Jiang 2013; detailed background is provided in Section 4). Penalized quasi likelihood (PQL) approach has been used by Hui, Müller, and Welsh (2017b) to conduct simultaneous (but nonhierarchical) selection of mixed effects in a GLMM framework with adaptive lasso and adaptive group lasso regularization. The CREPE (composite random effects penalty) estimator of Hui, Müller, and Welsh (2017a) conducts hierarchical variable selection in a GLMM with a single longitudinal outcome and employs a Monte Carlo EM (MCEM) algorithm of Wei and Tanner (1990) to maximize the likelihood. The CREPE estimator ensures that variables are included in the final model either as fixed effects only or as composite effects. Our proposed CEZIJ framework is related to Hui, Müller, and Welsh (2017a) in its ability to conduct hierarchical variable selection in GLMMs. However, unlike Hui, Müller, and Welsh (2017a), CEZIJ performs hierarchical variable selection in a joint model of multiple correlated longitudinal outcomes. Additionally, it can also incorporate any convexity constraint on the fixed effects.

(iv) *Scalability*: For any mobile game app, gargantuan volumes of user activity data are automatically accumulated. Analyzing such big datasets not only involves inferential problems associated with high-dimensional data analysis but also the computational challenges of processing large-scale (sample) longitudinal data. To process large longitudinal datasets, CEZIJ leverages the benefits of distributed computing. Recently, algorithmic developments for increased scalability and reduced computational time without sacrificing the requisite level of statistical accuracy have received significant attention (see, e.g., Jordan 2013; Jordan, Lee, and Yang 2018; Lee et al. 2015, and references therein). A popular approach is to conduct inference independently and simultaneously on K subsets of the full dataset and then form a global estimator by combining the inferential results from the K nodes in a computation-efficient manner. We take a similar approach for the hierarchical selection of fixed and random effects by using the split-and-conquer approach of Chen and Xie (2014) that splits the original dataset into K nonoverlapping groups, conducts variable selection separately in each group and uses a majority voting scheme in assimilating the results from the splits.

(v) *Prediction and segmentation*: Predictive analysis of new player behavior is fundamental for the maintenance of existing as well as for the creation of new advertisement based monetization routes in these gaming platforms. Statistically, this necessitates construction of prognostic models that cannot only forecast new user activity, engagement and dropout behavior but also dynamically update such forecasts over time as new longitudinal information about them arrives. Based on our fitted joint model, we construct dropout probability profiles (over time) for an out-of-sample generic player population and use them for segmentation of idiosyncratic player behaviors. Segmentation is a key analytical tool for managers. Users in different segments respond differently to varied marketing promotions. This enables managers to use relevant marketing promotions that better match user responses in different segments and increase efficiency of their marketing campaign.

We develop our joint modeling framework which accommodates all of the above mentioned extensions through an efficient and scalable estimation procedure. To the best of our knowledge, we are the first to study constrained joint modeling of high-

dimensional data. Though we demonstrate the applicability of the CEZIJ inferential framework for the disciplined study of freemium behavior, it can be used in a wide range of other applications that needs analyzing multiple high-dimensional longitudinal outcomes along with a time-to-event analysis. To summarize, the key features of our CEZIJ framework are:

- Joint modeling of the highly related responses pertaining to daily player activity and engagement as well as the daily dropout probabilities using the freemium mobile game data described in [Section 2](#);
- The possibility of acute zero-inflation in the player engagement distribution is addressed by modeling the conditional probability of no engagement given that the player had used the app in the day (see [Figure 3](#));
- Convexity constraints pertinent to domain expertise and prior beliefs are incorporated in the modeling framework in [Section 3](#);
- A penalized EM algorithm (Wei and Tanner 1990) is used for simultaneous selection of fixed and random effects wherein data-driven weighted ℓ_1 penalties are imposed on the fixed effects as well as on the diagonal entries of the covariance matrix of the random effects while the common regularization parameter λ is chosen by a BIC-type criterion (see Equation (9));
- Hierarchical selection of the fixed and random effects is conducted in [Section 4](#) by using a reweighted ℓ_1 minimization algorithm that alternates between estimating the parameters and redefining the data-driven weights such that the weights used in any iteration are computed from the solutions of the previous iteration;
- The divide and conquer approach in [Section 5](#) distributes the problem into tractable parallel subgroups resulting in increased scalability;
- Prediction of the dropout probabilities as well as the activity and engagement characteristics of new players with the predictions being dynamically updated as additional longitudinal information is recorded (see [Section 6](#)). Based on these dynamic churn probability curves from our fitted joint model, we conduct segmentations of player profiles that can be used by game managers to develop improved promotion and retention policies specifically targeting different dominant player-types.

2. Motivating Data: Activity, Engagement, Churn and Promotion Effects in Freemium Mobile Games

We consider daily player level gaming information for a mobile app game where users use robot avatars to fight other robots till one is destroyed. There were 38,860 players in our database and we tracked daily player level activity and purchases for 60 consecutive days starting from the release date of the game. We use a part of the data (players) for estimation and the other part as the hold out set for prediction (see details in [Section 6](#)). There were three modes of the game and level progression can only be attained through the principal mode. However, the players get rewards (henceforth called *in-game rewards*) if they win games in all three modes. For the two

nonprincipal modes, collecting rewards is the main objective. The players can use these rewards for improving their fighting equipments through upgrades of their existing inventories or in getting access to powerful new robots or for acquiring fancy game themes and background changes. The player can also buy these facilities (add-ons) using real money through direct IAP. There were only 0.28% of the players who used real money for buying add-ons. The players are given premium rewards, which has much higher order of magnitudes than regular rewards, if they promote the game or the developers through social media (inviting friends on facebook for games, facebook likes, youtube likes, tweets) or through the app center or by downloading other related apps from the developer. Approximately 7.2% of the players in our data had premium rewards. We record daily engagement of a player by appropriately combining her real money purchases (direct source of revenue) with her varied involvement in promoting the game in media (indirect source of monetization) with the notion being that the highest form of engagement is the one leading to purchases. Daily engagement is an extremely zero-inflated variable. We assess player behavior in terms of her daily total playing time (activity), engagement value, and dropout probability. We say that a player has dropped out if she has not logged-in for a month consecutively. For each player we have a host of time-dependent covariates generated through the game-play which we model as composite effects. They include current level of the game, number of games played daily in the three different modes of the game, how are the in-game rewards spent, etc. (see [Table 3](#) and summary [Table 4](#) in [Section D](#) of the supplementary materials for details). From a gaming perspective, it is very interesting to study the effects on gaming time of the amount of in-game rewards that the players spend on either upgrading existing robots or purchasing new robots. Another interesting feature of the game, was the usage of “gacha” mechanism (Toto 2012) which allowed the players to gamble in-game currency through lottery draws. The “gacha” is a very popular feature in freemium games (Kanerva 2016). We use the currency employed by players in “gacha” as well as their gains, as covariates in modeling engagement. Also, several promotional and retention strategies were used by the developers, which encourage player activity. [Figure 1](#) contains a flowchart summarizing the key components of the game. The promotions intrinsically were of four different flavors: (a) award more reward percentages and battery life, (b) sale on robots, (c) thanksgiving holiday promotions, (d) E-mail and app-message based notifications for retention. Also, there were three different kinds of sales on robots. Thus, there were six different promotional strategies, with only one of them (if at all) being employed on a single day (see [Table 4](#) and [Figure 3](#) in [Section D](#) of the supplementary materials). In [Figure 2\(a\)](#) and (b), we present the activity, engagement, and churn profiles of the players in our data. Interestingly, the proportion of players with positive engagement is below 10% from day 3 onward and drops to less than 1% after the first 21 days. [Figure 2\(c\)](#) and (d), respectively, shows the 25th, 50th, and 75th percentiles of the distribution of total time played and the average engagement amount on each of the 60 days. Note that from day 20 onward the distribution of average engagement shows increased variability. This is not unexpected given the observation from [Figure 2\(a\)](#) which shows that the proportion of players with pos-

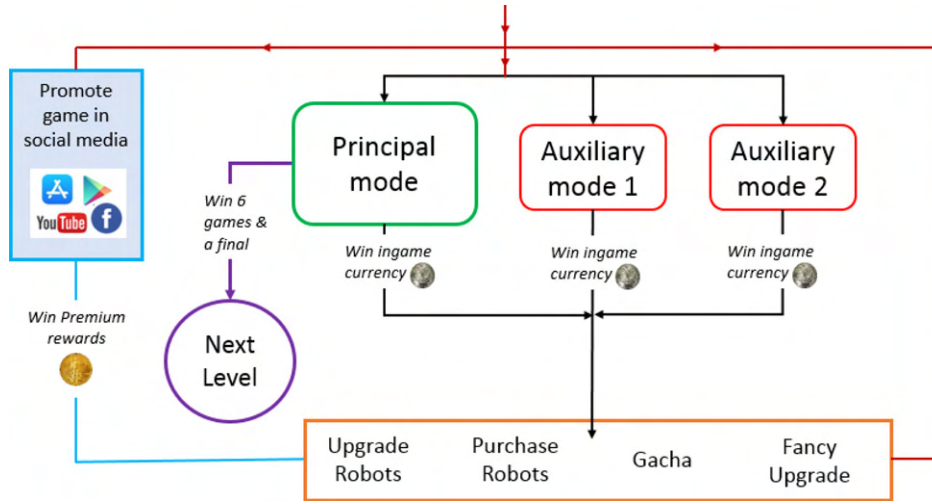


Figure 1. Game play flowchart.

itive engagement falls steadily. Also, note that the heavy tailed nature of the distributions of positive time played and positive engagement amount is evident from Figure 2 (Section D of the supplementary materials) which plots the empirical CDF of the two variables. So, in the following section we use log-normal distributions to model the nonzero activity and engagement values. Further details regarding the data are available in Section D of the supplementary materials.

3. CEZIJ Modeling Framework

Using the aforementioned motivation example, we now introduce our generic joint modeling framework. Consider data from n independent players where every player $i = 1, \dots, n$ is observed over m time points. Let \mathbb{A}_{ij} and \mathbb{E}_{ij} denote, respectively, the activity and engagement of player i at day j with $\mathbb{A}_i = (\mathbb{A}_{i1}, \dots, \mathbb{A}_{im})$ and $\mathbb{E}_i = (\mathbb{E}_{i1}, \dots, \mathbb{E}_{im})$ denoting the corresponding vector of longitudinal measurements taken on player i . Let \mathbb{D}_i denote the time of dropout for player i and \mathbb{C}_i the censoring time. We assume \mathbb{C}_i 's are independent of \mathbb{D}_i 's. Thus, $\mathbb{C}_i = m$ if player i never drops out. The observed time of dropout is $\mathbb{D}_i^* = \min(\mathbb{D}_i, \mathbb{C}_i)$, and the longitudinal measurements on any player i are available only over $m_i \leq \mathbb{D}_i^*$ time points. Suppose α_{ij} be the indicator of the event that player i is active ($\mathbb{A}_{ij} > 0$) on day j and ϵ_{ij} be the indicator that she positively engages ($\mathbb{E}_{ij} > 0$) on day j . Let $\pi_{ij} = \Pr(\alpha_{ij} = 1)$, $q_{ij} = \Pr(\epsilon_{ij} = 1 | \alpha_{ij} = 1)$, $\alpha_i = (\alpha_{i1}, \dots, \alpha_{im})$, and $\epsilon_i = (\epsilon_{i1}, \dots, \epsilon_{im})$. Note that, $\alpha_{ij} = 0$ implies $\mathbb{A}_{ij} = \mathbb{E}_{ij} = 0$ and also $\epsilon_{ij} = 0$. In these gaming apps, it is usually witnessed that any player's usage of the app always produces positive activity (however small). Thus, α_{ij} here corresponds to a player's daily activity indicator (AI). It forms the base (first level) of our joint model. The π_{ij} corresponds to daily usage probability whereas q_{ij} corresponds to the conditional probability of positive player engagement given that the player has used the app in the day. In Figure 3, we provide a schematic diagram of our joint model where we use two binary random variables: AI and engagement indicator (EI) to be, respectively, denoted α_{ij} and ϵ_{ij} . We jointly model the five components $[\alpha_i, \mathbb{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i] := [\mathbb{Y}_i^{(s)} : s \in$

$\{1, 2, 3, 4, 5\}]$ given the observations. Let \mathcal{I} be the full set of p predictors in the data with $\mathcal{I}_f \subset \mathcal{I}$ as the set of fixed effects (time invariant or not) and $\mathcal{I}_c = \mathcal{I} \setminus \mathcal{I}_f$ as the set of composite effect predictors, which are modeled by combination of fixed and random effects. Let $p_f = |\mathcal{I}_f|$ and $p_c = |\mathcal{I}_c|$ and so, $p_c + p_f = p$. For each of the first four submodels, $s = 1, \dots, 4$, we consider p fixed effects $\boldsymbol{\beta}^{(s)}$ (p_f of those are from the time invariant and the rest from the composite components) and p_c random effects $\mathbf{b}^{(s)}$ while for the dropout model, $s = 5$, we consider p new fixed effects $\boldsymbol{\beta}^{(5)}$ but share the random effects from the four submodels and calibrate their effects on dropouts through an association parameter vector $\boldsymbol{\eta}$. See Section 3.1 for further details.

Let $x_{ijk}^{(s)}$ denote the observed k th covariate value for the i th player on the j th day. Let $\mathbf{x}_{ij}^{(s)} = \{x_{ijk}^{(s)} | k \in \mathcal{I}\}$ and $\mathbf{z}_{ij}^{(s)} = \{z_{ijk}^{(s)} | k \in \mathcal{I}_c\}$ denote the set of covariate values pertaining to the in-model fixed and random effects; $\mathbb{X}^{(s)}$ and $\mathbb{Z}^{(s)}$, respectively, denote the data for these effects across all n players and $\boldsymbol{\beta} = \{\boldsymbol{\beta}^{(s)} : s \in \{1, 2, 3, 4\}\}$ and $\mathfrak{b} = \{\mathbf{b}^{(s)} : s \in \{1, 2, 3, 4\}\}$ be all the fixed and random effects across all players. To join the four models, we take a correlated random effects approach and assume that the random effects governing the four submodels have a multivariate Gaussian distribution. For player i , represent all her random effects by $\mathbf{b}_i = (\mathbf{b}_i^{(s)} : 1 \leq s \leq 4)$. We assume that $\{\mathbf{b}_i : 1 \leq i \leq n\}$ iid $N(\mathbf{0}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma}$ is the $4p_c \times 4p_c$ unknown covariance matrix. To model the dropouts, we again consider p new fixed effects $\boldsymbol{\beta}^{(5)}$ but share the random effects from the four submodels and calibrate their effects on dropouts through an association parameter vector $\boldsymbol{\eta}$. We model $[\mathbb{Y}^{(s)} : 1 \leq s \leq 5 | \mathbb{X}, \mathbb{Z}, \boldsymbol{\beta}, \boldsymbol{\Sigma}]$ as

$$\prod_{i=1}^n [\mathbf{b}_i | \alpha_i | \mathbb{X}_i^{(1)}, \boldsymbol{\beta}^{(1)}, \mathbb{Z}_i^{(1)}, \mathbf{b}_i^{(1)}] \\ [\mathbb{A}_i | \alpha_i, \mathbb{X}_i^{(2)}, \boldsymbol{\beta}^{(2)}, \mathbb{Z}_i^{(2)}, \mathbf{b}_i^{(2)}] \\ [\epsilon_i | \alpha_i, \mathbb{X}_i^{(3)}, \boldsymbol{\beta}^{(3)}, \mathbb{Z}_i^{(3)}, \mathbf{b}_i^{(3)}] \\ [\mathbb{E}_i | \alpha_i, \epsilon_i, \mathbb{X}_i^{(4)}, \boldsymbol{\beta}^{(4)}, \mathbb{Z}_i^{(4)}, \mathbf{b}_i^{(4)}] [\mathbb{D}_i | \mathbb{X}_i^{(5)}, \boldsymbol{\beta}^{(5)}, \mathbf{b}_i].$$

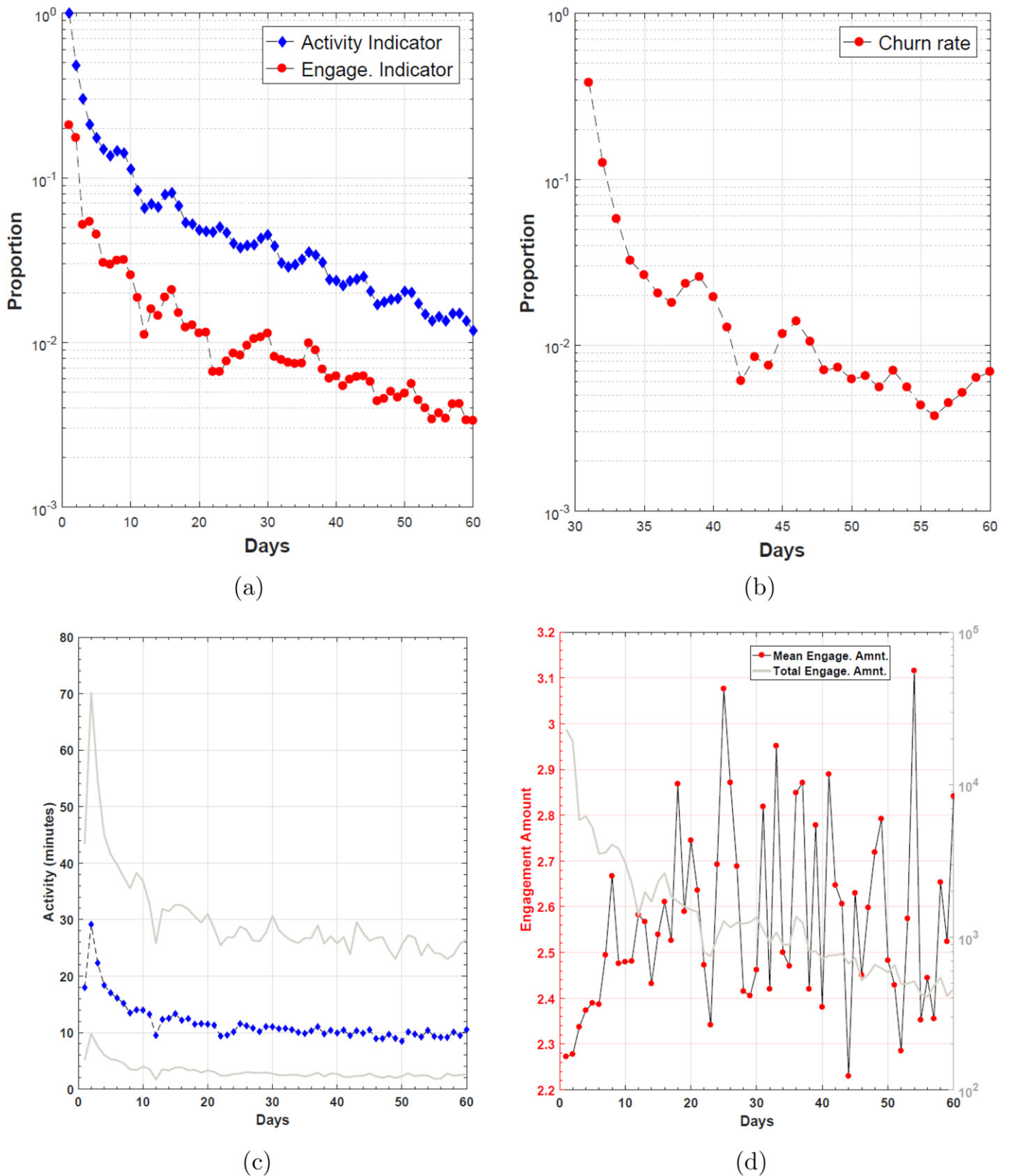


Figure 2. (a) Proportion of players active and proportion of players with positive engagement over 60 days. (b) Proportion of player churn from day 31 to day 60. (c) Median activity sandwiched between its 25th and 75th percentile. (d) Mean engagement amount and the total engagement amount over the 60 days.

Note that the dimension of each $\mathbf{b}_i^{(s)}$ in \mathbf{b}_i is p_c and that of $\mathbf{x}_{ij}^{(s)}$ is p . In the context of our mobile app game data, $p_c = 25$ and so Σ is 100×100 and $p = 31$ for each of the five submodels, thus making a set of 155 fixed effects (time invariant or not). See Section 6 for more details.

Remark 1. If we have data pertaining to social media interactions among players, it would be beneficial to include network or group effects among players. In the absence of such network information, we model $\mathbf{b}_i^{(s)}$ as iid across players.

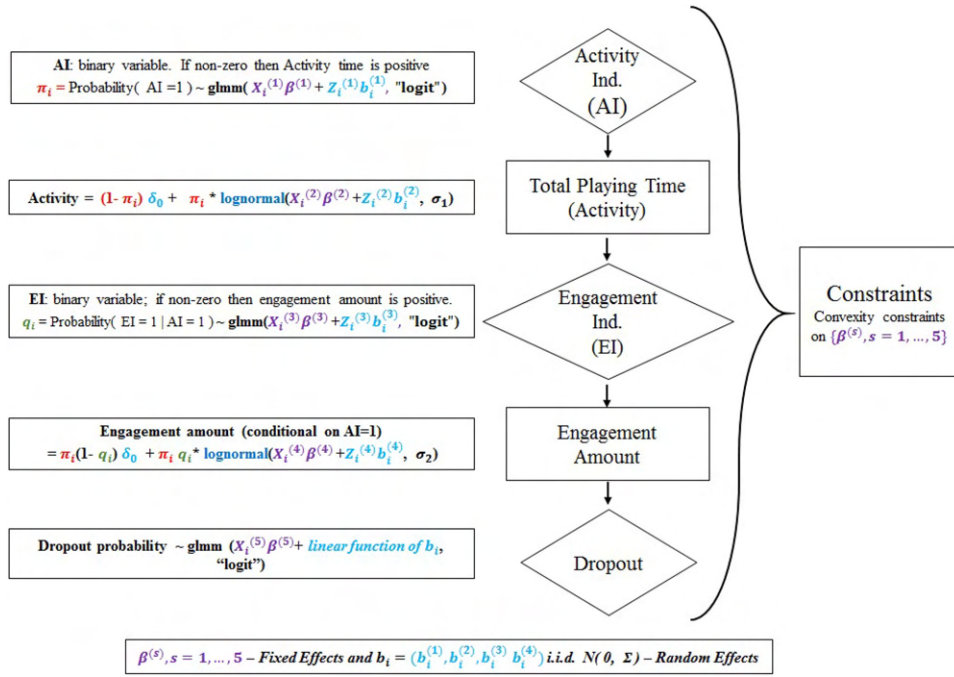


Figure 3. Schematic diagram of our joint model for player i . The suffix denoting the day number is dropped for presentational ease.

3.1. Longitudinal Submodels and Model for Dropouts

3.1.1. Zero Inflated Log-Normal for Modeling Activity

Since player i is active only at some time points j , the observed activity \mathbb{A}_i has a mix of many zeros and positive observations. In Equation (1), we consider a zero inflated (ZI) log normal model for \mathbb{A}_{ij} to capture both the prevalence of these excess zeros and possible large values observed in the support of \mathbb{A}_{ij} . Thus, the model for activity \mathbb{A}_{ij} has a mixture distribution with pdf

$$g_1(\alpha_{ij}, \mathbb{A}_{ij} | b_i^{(1)}, b_i^{(2)}) = (1 - \pi_{ij}) \mathbb{I}\{\alpha_{ij} = 0\} + \pi_{ij}(\sigma_1 \mathbb{A}_{ij})^{-1} \phi\left(\frac{\log \mathbb{A}_{ij} - \mu_{ij}}{\sigma_1}\right) \mathbb{I}\{\alpha_{ij} = 1\}, \quad (1)$$

where

$$\begin{aligned} \text{logit}(\pi_{ij}) &= \mathbf{x}_{ij}^{(1)T} \boldsymbol{\beta}^{(1)} + \mathbf{z}_{ij}^{(1)T} \mathbf{b}_i^{(1)} \rightarrow \text{Binary part}, \quad (2) \\ \text{and } \mu_{ij} &= \mathbf{x}_{ij}^{(2)T} \boldsymbol{\beta}^{(2)} + \mathbf{z}_{ij}^{(2)T} \mathbf{b}_i^{(2)} \rightarrow \text{Positive activity part}. \quad (3) \end{aligned}$$

The AI α_{ij} is modeled using a logistic regression model with random effects in Equation (2). In Equation (3), we use an identity link to connect the expected log activity with the covariates and the random effects. For convenience, hereon the dependence on the fixed effects and covariates are kept implicit in the notations and only the involved random effects are explicitly demonstrated.

3.1.2. Extreme ZI Log-Normal for Modeling Engagement

Note that, \mathbb{E}_i also has a mix of zeros and positive observations but the extreme zero events in the engagement variable are due to: (a) players are inactive on days and, (b) active players on a day may not exhibit engagement on the same day. To account for this excess prevalence of zeros, we use an extreme zero inflated

(EZI) log normal model that models $(\alpha_{ij}, \epsilon_{ij}, \mathbb{E}_{ij})$ as a flexible mixture distribution with joint pdf

$$g_2(\alpha_{ij}, \epsilon_{ij}, \mathbb{E}_{ij} | b_i^{(1)}, b_i^{(3)}, b_i^{(4)}) = (1 - \pi_{ij}) \mathbb{I}\{\alpha_{ij} = 0\} + \pi_{ij} g_3(\epsilon_{ij}, \mathbb{E}_{ij} | b_i^{(3)}, b_i^{(4)}) \mathbb{I}\{\alpha_{ij} = 1\}, \quad (4)$$

$$\text{where, } g_3(\epsilon_{ij}, \mathbb{E}_{ij} | b_i^{(3)}, b_i^{(4)}) = (1 - q_{ij}) \mathbb{I}\{\epsilon_{ij} = 0\} \quad (5)$$

$$+ q_{ij}(\sigma_2 \mathbb{E}_{ij})^{-1} \phi\left(\frac{\log \mathbb{E}_{ij} - \gamma_{ij}}{\sigma_2}\right) \mathbb{I}\{\epsilon_{ij} = 1\}$$

$$\text{logit}(q_{ij}) = \mathbf{x}_{ij}^{(3)T} \boldsymbol{\beta}^{(3)} + \mathbf{z}_{ij}^{(3)T} \mathbf{b}_i^{(3)} \rightarrow \text{Binary part}, \quad (6)$$

$$\gamma_{ij} = \mathbf{x}_{ij}^{(4)T} \boldsymbol{\beta}^{(4)} + \mathbf{z}_{ij}^{(4)T} \mathbf{b}_i^{(4)} \rightarrow \text{Positive engagement part}. \quad (7)$$

Note that a player can potentially engage ($\mathbb{E}_{ij} \geq 0$) only if she is active ($\alpha_{ij} = 1$) on that day. Thus, $g_3(\epsilon_{ij}, \mathbb{E}_{ij} | b_i^{(3)}, b_i^{(4)})$ in Equation (4) represents the joint pdf of $(\epsilon_{ij}, \mathbb{E}_{ij})$ conditional on the event that the player is active, that is, $\alpha_{ij} = 1$. However, even if the player is active, distribution of engagement again can have a mixture distribution, as the particular player may or may not exhibit positive engagement ($\mathbb{E}_{ij} > 0$). Thus, conditional on the player being active, we further model $(\epsilon_{ij}, \mathbb{E}_{ij})$ using another zero-inflated log normal model as shown in Equation (5). By combining Equations (4) and (5), intuitively, we use the EZI model to split the players into two groups: (1) who are not active and (2) who are active. Then conditional on being active, we further split the latter group of players into two additional segments: (1) who do not engage ($\epsilon_{ij} = 0$) and, (2) who engage ($\epsilon_{ij} = 1$) and thus demonstrate positive engagement ($\mathbb{E}_{ij} > 0$). Finally, we complete the specification of the EZI log normal model by connecting the binary response $\epsilon_{ij} | \alpha_{ij} = 1$ with the covariates and the random effects through a logit link in equation (6) and use an identity link for expected

log engagement γ_{ij} in Equation (7). Note that even though we model α_{ij} in Equations (1) and (4) using g_1 and g_2 , respectively, there is no discordance as $g_1(\alpha_{ij}) = g_2(\alpha_{ij})$ for all (i, j) .

3.1.3. Model for Dropouts

For the discrete time hazard of dropout, we model $\lambda_{ij} := P(\mathbb{D}_i = j | \mathbb{D}_i \geq j, \mathbf{b}_i)$ as

$$\text{logit}(\lambda_{ij}) = \mathbf{x}_{ij}^{(5)T} \boldsymbol{\beta}^{(5)} + \boldsymbol{\eta}^T \mathbf{b}_i, \tag{8}$$

and the pmf of \mathbb{D}_i^* is

$$g_4(\mathbb{D}_i^* = d | \mathbf{b}_i) = \left\{ \prod_{j=1}^{d-1} (1 - \lambda_{ij}) \right\} \lambda_{id}^{\delta_i^{\mathbb{D}}} (1 - \lambda_{id})^{1 - \delta_i^{\mathbb{D}}},$$

where $\delta_i^{\mathbb{D}} = I(\mathbb{D}_i \leq \mathbb{C}_i)$ is the indicator of dropout occurrence. Here $\boldsymbol{\eta}$ is a parameter vector that relates the longitudinal outcomes and the dropout time via the random effects \mathbf{b}_i . This approach to modeling the dropouts through Equation (8) is analogous to the shared parameter models in clinical trials that are used to account for potential not missing at random (NMAR) responses (see, e.g., Vonesh, Greene, and Schluchter 2006; Guo and Carlin 2004). If $\boldsymbol{\eta} = \mathbf{0}$ then the dropout is ignorable given the observed data. Figure 3 contains a schematic diagram of our joint model.

3.1.4. Correlating the Random Effects and Linking the Submodels

All the submodels described above carry information about the playing behavior of individuals and are therefore inter-related. To get the complete picture and to account for the heterogeneity across individual's outcomes, one must combine these effects by correlating the multiple outcomes. Without inter-relating or jointly considering these outcomes, it is not only hard to answer questions about how the evolution of one response (e.g., activity) is related to the evolution of another (e.g., engagement) or who is likely to dropout, but also problematic to model the heterogeneity. In such cases, it is natural to consider models where the dependency among the responses may be incorporated via the presence of one or more latent variables. A flexible solution is to model the association between different responses by correlating the random heterogeneous effects from each of the responses. In our joint modeling approach, random effects are assumed for each longitudinal response and they are associated by imposing a joint multivariate distribution on the random effects, that is, $\mathbf{b}_i = (\mathbf{b}_i^{(s)} : 1 \leq s \leq 4) \sim N(\mathbf{0}, \boldsymbol{\Sigma})$. Such a model borrows information across the various touch points and offers an intuitive way of describing the dependency between the responses. For example, questions such as, “is engagement related to activity for an individual?”, or “does higher activity increase the probability of engagement” can be answered using the estimated covariance structure of $\boldsymbol{\Sigma}$. Furthermore, we assume that the dependency between the longitudinal outcomes and the risk of dropout are described by the random effects \mathbf{b}_i and the covariates. In our context this is reasonable since, for instance, the longitudinal outcome AI may characterize player engagement, and player engagement can in turn influence the risk of dropout.

Table 1. Parameter constraints and their interpretation.

Constraint	Description
$\beta_{\text{weekend}}^{(s)} \geq 0, \forall s$	Expect increased player activity on weekends
$\beta_{\text{timesince}}^{(s)} \leq 0, \forall s$	Expect lower player activity as time since last login increases
$\beta_{\text{promV}}^{(1)}, \beta_{\text{promIV}}^{(1)} \geq 0$	Expect promotions IV, V to increase player activity
$\beta_{\text{promV}}^{(1)} \geq \beta_{\text{promIV}}^{(1)}$	Expect promotion V to have a higher positive impact on player activity than promotion IV
$\beta_{\text{prom}(i)}^{(2)} \geq 0$ for $i \neq \text{IV}$	All promotions other than IV to have a nonnegative impact on activity
$\beta_{\text{promIII}}^{(2)} \geq \beta_{\text{promV}}^{(2)}$	Promotions III leads to a higher increase in activity than promotion V
$\beta_{\text{promVI}}^{(2)} \geq \beta_{\text{promII}}^{(2)} \geq \beta_{\text{promV}}^{(2)}$	Promotions VI has the largest positive impact on activity followed by promotions II and V
$\beta_{\text{promV}}^{(2)} \geq \beta_{\text{promIV}}^{(2)}$	Promotion V leads to a higher increase in activity than promotion IV

NOTES: Here $\beta_{\text{prom}(i)}^{(s)}$ indicates the fixed effect coefficient for promotion $i = I, \dots, VI$ under model $s = 1, \dots, 5$.

3.2. Parametric Constraints

The CEZIJ framework can incorporate any convexity constraints on the fixed effects: $\mathfrak{f}^{(s)}(\boldsymbol{\beta}^{(s)}) \leq 0, s = 1, \dots, 5$, where \mathfrak{f} is any prespecified convex function. In the mobile game platform modeling application, domain expertise can be incorporated into our framework via these constraints. For example, industry belief dictates that all other factors remaining fixed, players have higher chance of being active in the game on weekends than on week days. Thus, a sign constraint on the unknown fixed effect coefficient for the variable ($\beta_{\text{weekend}}^{(s)} > 0$) that indicates whether day j is a weekend, is a simple yet effective way to include this additional information into our estimation framework. Also, different promotional and retention strategies used in these games are incorporated in the model as fixed effects through the binary variables demarcating the days they were applied (see Figure 3 in Section D of the supplementary materials for a distribution of the various promotion strategies across the $m = 60$ days). These strategies often have previously known efficacy levels which imply monotonicity constraints on their effects. For example, E-mail and app messaging based retention scheme should have at least a nonnegative increment effect on the daily usage probabilities π_{ij} 's; the engagement effect of a promotion that offers sale on only selected robots cannot exceed the increment effect of sale on all robots. As such, in our mobile game application, we assimilate these side information through structured affine constraints: $\mathbf{C}^{(s)} \boldsymbol{\beta}^{(s)} \leq \kappa^{(s)}, s = 1, \dots, 5$ where $\mathbf{C}^{(s)}$ and $\kappa^{(s)}$ are known. Details about these constraints are provided in Table 1 (and Table 5 in Section D of the supplementary materials), where we describe the six promotion strategies and the constraints that have been included in our estimation framework along with their business interpretation.

4. Variable Selection in CEZIJ

In the absence of any prior knowledge regarding variables that may appear in the true model, we conduct automated variable selection. Selection of fixed and random effect components in

a mixed model framework has received considerable attention. Under the special case of a linear mixed effect model, Bondell, Krishna, and Ghosh (2010) and Ibrahim et al. (2011) proposed penalized likelihood procedures to simultaneously select fixed and random effect components, while Fan and Li (2012), Peng and Lu (2012), and Lin, Pang, and Jiang (2013) conduct selection of fixed and random effects using a two stage approach. Procedures to select only the fixed effects or the random effects have also been proposed under a GLMM framework (see Pan and Huang 2014 and references therein). Simultaneous selection of fixed and random effect components in a GLMM framework is, however, computationally challenging. The high-dimensional integral with respect to the random effects in the marginal likelihood of GLMM often has no analytical form and several approaches have been proposed to tackle this computational hurdle: for example, Laplacian approximations (Tierney and Kadane 1986), adaptive quadrature approximations (Rabe-Hesketh, Skrondal, and Pickles 2002), PQL (Breslow and Clayton 1993), and EM algorithm (McCulloch 1997). We use a penalized EM algorithm and for proper interpretation of composite effects we conduct joint variable selection of fixed and random effects in a hierarchical manner, which ensures that nonzero random effects in the model are accompanied by their corresponding nonzero fixed effects. Let $\Theta = \{\beta^{(1)}, \dots, \beta^{(5)}, \sigma_1, \sigma_2, \eta, \text{vec}(\Sigma)\} := \{\theta, \text{vec}(\Sigma)\}$ denote the vector of all parameters to be estimated. The marginal log-likelihood of the observed data under the joint model is

$$\begin{aligned} \ell(\Theta) &= \sum_{i=1}^n \log \int p(\alpha_i, A_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i^* | \mathbf{b}_i, \theta) p(\mathbf{b}_i | \Sigma) d\mathbf{b}_i \\ &= \sum_{i=1}^n \ell_i(\Theta), \text{ where,} \\ \ell_i(\Theta) &= -\frac{1}{2} \log |\Sigma| \\ &\quad + \log \int \exp \left(\sum_{j=1}^{m_i} \log p(\alpha_{ij}, A_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_{ij}^* | \mathbf{b}_i, \theta) \right. \\ &\quad \left. - \frac{1}{2} \mathbf{b}_i^T \Sigma^{-1} \mathbf{b}_i \right) d\mathbf{b}_i. \end{aligned}$$

We estimate Θ using the EM algorithm for joint models (Rizopoulos 2012) where we treat the random effects \mathbf{b}_i as “missing data” and obtain $\hat{\Theta}$ by maximizing the expected value of the complete data likelihood $\ell^{\text{cl}}(\Theta, \mathbf{b})$ where

$$\begin{aligned} \ell^{\text{cl}}(\Theta, \mathbf{b}) &= -\frac{n}{2} \log |\Sigma| \\ &\quad + \sum_{i=1}^n \left(\sum_{j=1}^{m_i} \log p(\alpha_{ij}, A_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_{ij}^* | \mathbf{b}_i, \theta) \right. \\ &\quad \left. - \frac{1}{2} \mathbf{b}_i^T \Sigma^{-1} \mathbf{b}_i \right) = \sum_{i=1}^n \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i). \end{aligned}$$

Denote the Q-function $\ell^{\text{Q}}(\Theta) = \sum_{i=1}^n \mathbb{E} \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i)$ where the expectation is over the conditional distribution of \mathbf{b}_i given the observations at the current parameter estimates. We solve

the following maximization problem involving a penalized Q-function for variable selection

$$\begin{aligned} \max_{\theta, \Sigma > 0} \quad & \ell^{\text{Q}}(\Theta) - n\lambda \sum_{s=1}^5 \sum_{r=1}^p \left(c_{sr} |\beta_{sr}| + d_{sr} \Sigma_{rr}^{(s)} \mathbb{I}\{r \in \mathcal{I}_c\} \right) \\ \text{subject to} \quad & f^{(s)}(\beta^{(s)}) \leq 0, \quad s = 1, \dots, 5. \end{aligned} \quad (9)$$

Here, $\beta^{(s)} = \{\beta_{sr} : r \in \mathcal{I}\}$ and Σ is notationally generalized to include random effects corresponding to all p fixed effects—time invariant or not by introducing harmless zero rows and columns corresponding to time-invariant effects. This is done for presentational ease only to keep the indices same for the fixed and random effects and such degenerate large Σ matrix never crops in the computations. Also, $\Sigma_{rr}^{(s)}$ is the r th element of the vector $(\Sigma_{1+p_c(s-1), 1+p_c(s-1)}, \dots, \Sigma_{p_c s, p_c s})$ which represents the segmented covariance matrix corresponding to the s th model, \mathcal{I}_c is the index set of all composite effects and $\lambda > 0$ is the common regularization parameter which is chosen using a BIC-type criterion (Bondell, Krishna, and Ghosh 2010; Lin, Pang, and Jiang 2013; Hui, Müller, and Welsh 2017a) given by $\text{BIC}_\lambda = -2\ell^{\text{Q}}(\hat{\Theta}) + \log(n) \dim(\hat{\Theta})$ where $\dim(\hat{\Theta})$ is the number of nonzero components in $\hat{\Theta}$.

In many practical applications the composite effects impose the following hierarchy between fixed and random effects: a random component can have a nonzero coefficient only if its corresponding fixed effect is nonzero (Hui, Müller, and Welsh 2017a). To induce such hierarchical selection of fixed and random effects, we solve Equation (9) using a reweighted ℓ_1 minimization algorithm that alternates between estimating Θ and redefining the data-driven weights $(c_{sr}, d_{sr}) \in \mathbf{R}_+^2$ such that the weights used in any iteration are computed from the solutions of the previous iteration and are designed to maintain the hierarchy in selecting the fixed and random effects through their construction (see Candès, Wakin, and Boyd 2008; Zhao and Kočvara 2015; Lu, Lin, and Yan 2015, for details on these kind of approaches). Suppose $\Theta^{(t)}$ denote the solution to the maximization problem in Equation (9) at iteration t . Then we set $c_{sr}^{(t)} = \min(|\beta_{sr}^{(t)}|^{-\nu}, \epsilon_1^{-1})$ and $d_{sr}^{(t)} = \min(|\Sigma_{rr}^{(s,t)}|^{-\nu} |\beta_{sr}^{(t)}|^{-\nu}, \epsilon_2^{-1})$ for iteration $(t+1)$ with $\nu = 2$. We take $\epsilon_1 = 10^{-2}$ to provide numerical stability and to allow a nonzero estimate in the next iteration given a zero valued estimate in the current iteration (Candès, Wakin, and Boyd 2008) and fix $\epsilon_2 = 10^{-4}$ to enforce a large penalty on the corresponding diagonal element of Σ in iteration $(t+1)$ whenever $|\beta_{sr}^{(t)}| = 0$. Note that whenever $r \in \mathcal{I}_c$, the penalty d_{sr} on the diagonal elements of Σ encourages hierarchical selection of random effects. In Section C.2 of the supplementary materials, we conduct simulation experiments to demonstrate this property of our reweighted ℓ_1 procedure for solving Equation (9). We end this section with the observation that although the maximization problem based on criterion (9) does not conduct any selection on the association parameters η , it achieves that goal implicitly through the selection of the random effects.

5. Estimation Procedure

In this section, we discuss two key aspects of the estimation process.

5.1. Solving the Maximization Problem

We use an iterative algorithm to solve the maximization problem in Equation (9) which is analogous to the MCEM algorithm of Wei and Tanner (1990). Let $\Theta^{(t)}$ denote the parameter estimates at iteration t . In iteration $t + 1$, the MCEM algorithm performs the following two steps until convergence:

E-step: Recall $\mathbb{Y}_i = [\mathbb{Y}_i^{(s)}, s = 1, \dots, 5]$. Evaluate $\ell_{(t)}^Q(\Theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{b}_i | \Theta^{(t)}, \mathbb{Y}_i} \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i)$ where the expectation above is taken with respect to the conditional distribution of \mathbf{b}_i given the observations \mathbb{Y}_i at the current estimates $\Theta^{(t)}$. Thus,

$$\begin{aligned} & \mathbb{E}_{\mathbf{b}_i | \Theta^{(t)}, \mathbb{Y}_i} \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i) \\ &= \int \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i) p(\mathbf{b}_i | \alpha_i, \mathbb{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i^*, \Theta^{(t)}) d\mathbf{b}_i \\ &= \exp \{-\ell_i(\Theta^{(t)})\} \\ & \int \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i) p(\alpha_i, \mathbb{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i^* | \theta^{(t)}, \mathbf{b}_i) \phi_{p_c}(\mathbf{b}_i | \mathbf{0}, \Sigma^{(t)}) d\mathbf{b}_i, \end{aligned}$$

where, $\phi_{p_c}(\cdot | \mathbf{0}, \Sigma^{(t)})$ is the p_c dimensional normal density with mean $\mathbf{0}$ and variance $\Sigma^{(t)}$. Note that the expectation involves a multivariate integration with respect to the random effects \mathbf{b}_i which is evaluated by Monte Carlo integration. We approximate it as

$$\left(\sum_{d=1}^D \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i^d) p(\alpha_i, \mathbb{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i^* | \theta^{(t)}, \mathbf{b}_i^d) \right) / \left(\sum_{d=1}^D p(\alpha_i, \mathbb{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i^* | \theta^{(t)}, \mathbf{b}_i^d) \right),$$

where \mathbf{b}_i^d is a random sample from $\phi_{p_c}(\cdot | \mathbf{0}, \Sigma^{(t)})$ and $D = 2000$ is the number of Monte Carlo samples.

M-step: Solve the following maximization problem with data driven adaptive weights $(c_{sr}^{(t)}, d_{sr}^{(t)})$

$$\begin{aligned} \Theta^{(t+1)} &= \arg \max_{\theta, \Sigma > 0} \ell_{(t)}^Q(\Theta) - n\lambda \\ & \sum_{s=1}^5 \sum_{r=1}^p \left(c_{sr}^{(t)} |\beta_{sr}| + d_{sr}^{(t)} \Sigma_{rr}^{(s)} \mathbb{I}\{r \in \mathcal{I}_c\} \right) \\ & \text{subject to } \mathfrak{f}^{(s)}(\beta^{(s)}) \leq 0, s = 1, \dots, 5. \end{aligned}$$

The maximization problem above decouples into separate components that estimate $\beta^{(s)}$ as solutions to convex problems and Σ as a solution to a nonconvex problem. To solve the convex problems involving $\beta^{(s)}$, we use a proximal gradient descent algorithm after reducing the original problem to an ℓ_1 penalized least squares fit with convex constraints. See James, Paulson, and Rusmevichientong (2013) for related approaches of this kind. For estimating Σ , we use the coordinate descent algorithm of Wang (2014) that solves a lasso problem and updates Σ one column and row at a time while keeping the rest fixed. Further details regarding our estimation procedure is presented in Section A of the supplementary materials.

5.2. Split and Conquer

To enhance the computational efficiency of the estimation procedure, we use the split-and-conquer approach of Chen and Xie (2014) to split the full set of n players into K nonoverlapping groups and conduct variable selection separately in each group by solving K parallel maximization problems represented by Equation (9). Following Chen and Xie (2014), the selected fixed and random effects are then determined using a majority voting scheme across all the K groups as described below.

Let $\hat{\beta}^{(s)}[k]$ and $\hat{\Sigma}^{(s)}[k]$ denote, respectively, the estimate of the fixed effect coefficients for model s and the estimate of the p_c diagonal elements of Σ for model s on split k obtained by solving the maximization problem (9), where $k = 1, \dots, K$. We construct the set of selected effects as

Set of fixed effects:

$$\hat{\mathcal{I}}^{(s)} = \left\{ r: \sum_{k=1}^K \mathbb{I}(\hat{\beta}_{sr}[k] \neq 0) \geq w_0, r = 1, \dots, p \right\}$$

Set of random effects:

$$\hat{\mathcal{I}}_R^{(s)} = \left\{ r: \sum_{k=1}^K \mathbb{I}(\hat{\Sigma}_{r+p_c(s-1), r+p_c(s-1)}^{(k)} > 0) \geq w_1, r = 1, \dots, p_c \right\}.$$

Here w_0, w_1 are prespecified thresholds determining the severity of the majority voting scheme. For large datasets as in mobile apps application, a distributed computing framework utilizing the above scheme leads to substantial reduction in computation time. Section C of the supplementary material presents a discussion of the split-and-conquer approach along with numerical experiments that demonstrate the applicability of this method in our setting where data driven adaptive weights are used in the penalty and variable selection is conducted simultaneously across multiple models. Finally, based on the selected fixed and random effect components in $\hat{\mathcal{I}}^{(s)}$ and $\hat{\mathcal{I}}_R^{(s)}$, we use the entire data and estimate their effects more accurately by maximizing the likelihood based on only those components using the standard EM algorithm.

6. Analysis of Freemium Mobile Games Using CEZIJ

We apply our proposed CEZIJ methodology to the freemium mobile game data discussed in Section 2. This dataset holds player level gaming information for 38,860 players observed over a period of 60 days. The analyses presented here uses a sample of 33,860 players for estimation and the remaining 5000 players for out of sample validation. See Section D in the supplementary materials for a detailed description of the data. For submodels $s = 1, \dots, 4$, we consider a set of 30 predictors, of which 24 can have composite effects. The 24 composite effects are listed in Table 3 (Serial No. 1–24) of the supplementary materials. The remaining 6 predictors are the 6 promotion strategies summarized in Section D and Table 5 of the supplementary materials. We treat these promotion strategies as potential fixed effects with no corresponding random effect counterparts. For the dropout model, which shares its random effects with the four submodels, the entire list of 30 candidate predictors is taken as potential fixed effects. Overall, the selection mechanism must

Table 2. Selected fixed effect coefficients and their estimates under the submodels act. indicator, activity time, engag. indicator, engagement amount, and dropout.

Predictors	Act. indicator $\hat{\beta}^{(1)}$	Activity time $\hat{\beta}^{(2)}$	Engag. indicator $\hat{\beta}^{(3)}$	Engag. amount $\hat{\beta}^{(4)}$	Churn $\hat{\beta}^{(5)}$
intercept	-4.648*	0.932*	-1.560*	0.953*	-1.902
avg_session_length	-	0.269*	0.198*	-	-
p_fights	0.378*	0.169*	-0.126*	-	-
a1_fights	0.303*	0.379*	0.274*	-	-
a2_fights	0.334*	0.216*	-0.492*	-	-
level	0.084*	0.304*	0.282*	-	-
robot_played	-	-	-	-	-
gacha_sink	-	0.201*	0.509*	-	0.129
gacha_premium_sink	-	-	-	-	-
pfight_source	-	0.144*	-	-	-
a1fight_source	0.030	-0.239*	-0.727*	-	-
a2fight_source	-0.240*	-0.192*	0.482*	0.331*	-
gacha_source	0.182*	-0.212*	-	-	-
gacha_premium_source	-	-	0.133*	-	-
robot_purchase_count	0.134*	-	-	-	-
upgrade_count	0.093*	0.112*	0.404*	-	-
lucky_draw_wg	-	-	-0.240*	-	-
timesince	-2.065*	-0.641*	-0.229*	-	3.502
lucky_draw_og	-0.230*	-	-0.469*	-	-
fancy_sink	-	-	-0.110*	-	-
upgrade_sink	0.037*	-	-0.272*	-	-
robot_buy_sink	-	-	0.159	-	-
gain_gachaprem	-	-	-	-	-
gain_gachagrind	-0.127*	0.180*	-	-	-
weekend	0.302*	0.358*	-	-	-
promotion I	-	-	-	-1.153	-0.894
promotion II	0.178	0.134	-0.189	-	-0.934
promotion III	-0.129	-	-0.166	-1.791	-3.500
promotion IV	-	-	0.164	-3.345	-0.673
promotion V	-	-	-5.000	-	0.828
promotion VI	0.290	0.249	0.131	-2.389	-1.509

NOTES: The selected random effects are those variables that exhibit a (*) over their estimates. See Table 3 in Section D of the supplement for a detailed description of the covariates.

select random effects from a set of 100 potential random effects (24 for each of the four submodels and their 4 intercepts) and select fixed effects from a set of 155 potential fixed effects (30 for each of the five submodels and their 5 intercepts).

To model the responses at time point j , we consider time $j - 1$ values of the predictors that contain gaming characteristics of a player simply because at time j these characteristics are known only up to the previous time point $j - 1$. These gaming characteristics are marked with an (*) in Table 3 of the supplementary materials. All the remaining predictors like weekend indicator and the 6 indicator variables corresponding to the promotion strategies are applied at time j . We initialize the CEZIJ algorithm by fitting a saturated model on a subset of 200 players, which was also used to initialize the weights c_{sr} , d_{sr} in criterion (9). Finally, our algorithm is run on $K = 20$ splits where each split holds $n_k = 1693$ randomly selected players with the majority voting parameters w_0 , w_1 fixed at 12 and the regularization parameter λ_k is chosen as that value of $\lambda \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 0.25, 0.5, 1, 5\}$ which minimizes BIC_λ . Table 6 in Section E of the supplementary materials presents the voting results for each candidate predictor across the 20 splits.

6.1. The Fitted Joint Model and Its Interpretations

The final list of selected predictors and their estimated fixed effect coefficients for the submodels of AI, activity time (daily total time played), EI, engagement amount, and dropout is presented in Table 2. See Table 3 (Section D of the supplementary materials) for the description of the covariates. The selected

composite effects are those predictors that exhibit a (*) over their coefficient estimates in Table 2. All the selected fixed and random effects obey the hierarchical structure discussed in Section 4. We next discuss the fitted coefficients for each submodel.

Activity indicator: For modeling probability of AI, the CEZIJ methodology selects 18 fixed effects of which 14 are composite effects. As AI forms the base of our joint model, the fixed effects of its estimated marginal distribution have the least nuanced interpretation among the 5 submodels. All other things remaining constant, there is an overall increase in the odds for AI by 35% on the weekends and an 8% increase in odds for each level advancement in the game. Similarly, the conditional odds is boosted by 20%, 14%, and 9%, respectively, if the gacha game was played or robot purchases or upgrades were made in the previous days. Absence of log-in in the previous day adversely affects the odds with an average decrement of 88% for each absent day. Promotions II and VI, which provide sale on robots at different dosages are positively associated and increase odds of AI by approximately 20% and 30%, respectively.

Activity time: In this case the selection mechanism selects 17 fixed effects of which 15 are composite effects. The signs on the coefficients of `timesince` and `weekend` align with the constraints imposed on them and along with the game characteristics like number of primary and auxiliary fights played, level progressions, and robot upgrades, continue to provide a similar interpretation as with the AI model. This is the second layer of joint model which is conditioned on positive login occurrence. A key difference between these two models, however, lies in the inclusion of predictors `avg_session_length`,

`gacha_sink`, and `pfight_source`. They indicate that, keeping other things fixed, players interacting with the game through spending in-game currencies or winning the same through principal fights on the previous day have the natural incentive to spend more in-game time on the following day. In line with the monotonicity constraints imposed on the promotion strategies for this model, the coefficient for promotion VI is both positive and bigger than the coefficient for promotion II thus indicating that the strategy to promote sale on all robots has a higher impact on activity time than the strategy to offer the special particular “Boss” robots at a discount.

Engagement indicator and amount: Recall from Section 3 that we use an EZI log normal model for the engagement amount by first building a separate model for the probability of EI given activity. For the submodels that model EI and the engagement amount, the CEZIJ methodology selects, respectively, 22 fixed effects of which 16 are composite effects and 6 fixed effects of which 2 are composite effects. Direct interpretation of the fixed effect coefficients is difficult here, as this submodel is conditioned on the first two submodels. We see that some of the key player engagement characteristics like number of auxiliary fights played, level progression, in-app virtual currency spent and earned seem to positively impact the conditional likelihood of positive engagement at subsequent time points. A significant finding is that among the three different fight modes, only auxiliary fight second mode which involve time restricted fights seems to lead to substantially higher player engagement implying that all other variables remaining constant, player engagement in game promotion through social media is more while playing time attack fights.

Dropout: In this case, the selection mechanism selects 9 fixed effects. The sign on the coefficient for `timesince` is positive, which is natural, and indicates that players who do not frequent the game often (low frequency of AI) exhibit a high likelihood of dropping out at subsequent time points. It is also interesting to see, through `gacha_sink`, that all else being equal, players who spend more of their virtual currencies on gacha exhibit a high likelihood of dropping out at subsequent time points. This can potentially be explained through a “make-gacha-work-for-all-players” (Agelle 2016) phenomenon where the player spends a major portion of her virtual currency on gacha however the value of the items won is largely worthless when compared to the amount of currency spent, thus inducing a lack of interest in the game at future time points. All the promotions with exception of promotion V, reduce the odds of dropouts validating their usage as retention schemes.

From the heatmap in Figure 4, the random effects of the selected composite effect predictors demonstrate correlations within the four submodels that were modeled jointly, indicating that players exhibit idiosyncratic profiles over time. Moreover, we notice several instances of cross correlations across the four submodels. For example, from Figure 5, the random effect associated with the number of championship fights played (predictor `p_fights`) in the AI model has a positive correlation with the amount of virtual currency earned through auxiliary fights (predictor `a2_fights_source`) played in the model for activity time which suggests that the modeled responses are correlated

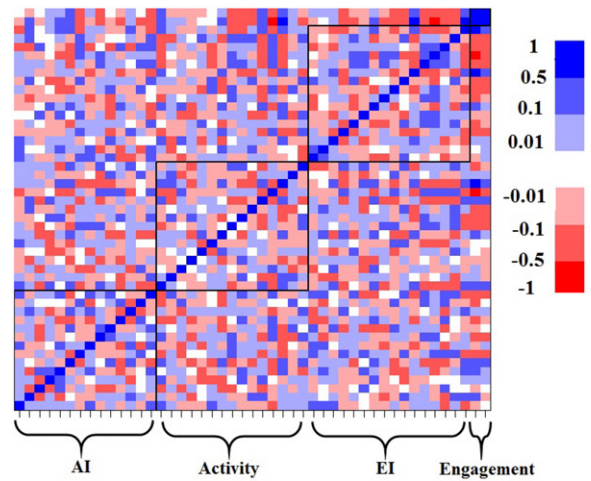


Figure 4. Heatmap of the 47×47 correlation matrix obtained from $\hat{\Sigma}$. On the horizontal axis are the selected composite effects of the four submodels: AI, activity time, EI, and engagement amount. The horizontal axis begins with the intercept from the AI model and ends with `a2fight_source` from the engagement amount model.

for a player. Our joint model allows us to borrow information across these related responses and may aid game managers and marketers in understanding how the outcomes depend on each other.

6.2. Out-of-Sample Validation

We use the hold-out sample of 5000 players from the original data for assessing the predictive accuracy of our model. Our scheme consists of predicting the four outcomes—AI, activity time, EI, and engagement amount, dynamically over the next 29 days using the fitted model discussed in Section 6.1. Note that the time frame of prediction covers the first 30 days of game usage for each player, and so by definition, no player drops out which leaves us with the aforementioned four outcomes to predict. As benchmarks to our fitted model, we consider four competing models—Benchmark I to Benchmark IV which we describe below.

For Benchmark I we consider a setup where there are no random effects, the outcomes are not modeled jointly and variable selection is conducted using the R-package `glmLasso` (Schelldorfer, Meier, and Bühlmann 2014) that uses an ℓ_1 -penalized algorithm for fitting high-dimensional GLMMs with logit links for AI, EI and identity link for the two continuous outcomes of positive activity time and engagement amount. In case of Benchmark II, we continue to model the outcomes separately and use the R-package `rpq1` (Hui, Müller, and Welsh 2017b) that performs joint selection of fixed and random effects in GLMMs using a regularized PQL (Breslow and Clayton 1993) with similar link functions as used in Benchmark I. The remaining two Benchmark models rely on the selected variables from the CEZIJ model itself and do not conduct their respective variable selection. In particular, Benchmark III uses the selected predictors from the CEZIJ methodology and models the outcomes via generalized linear models with logit links for AI, EI and identity link for the two continuous outcomes of positive activity time and engagement amount. Thus Benchmark III, like

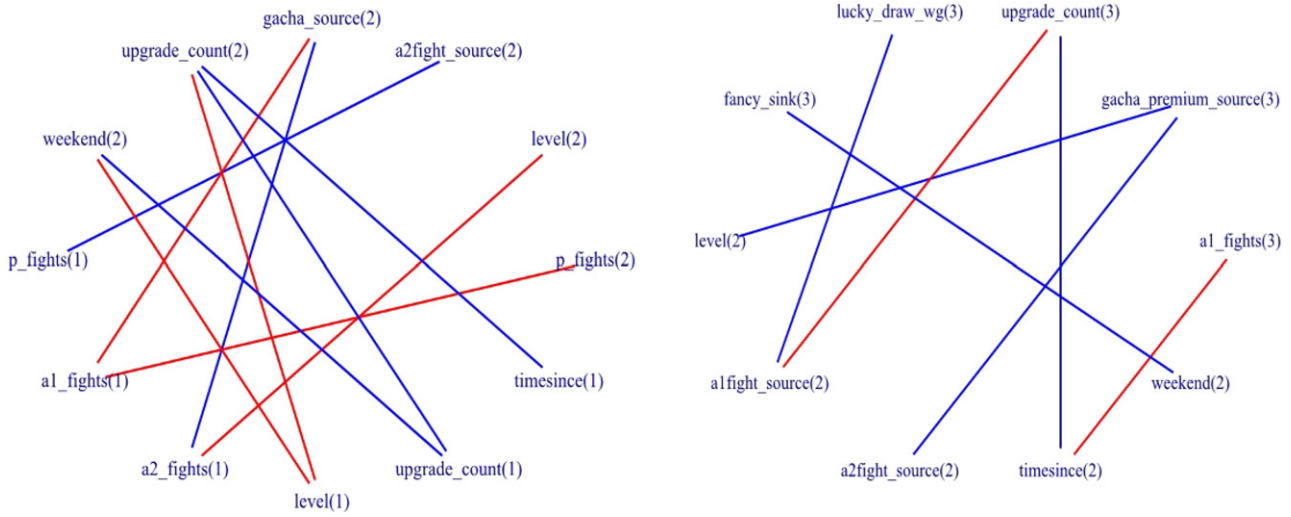


Figure 5. Two networks that demonstrate several cross correlations across the models. Blue line represents positive correlation and red line represents negative correlation. The model numbers are inside the parenthesis next to the predictor names. Left: Key cross correlations between the submodels AI and activity time. Right: Key cross correlations between the submodels activity time and EI.

Benchmark I, represents a setup where there are no random effects and the outcomes are not modeled jointly. Benchmark IV, on the other hand, represents a more sophisticated setup wherein it resembles the fitted CEZIJ model in every aspect except that the random effects across the four submodels are not correlated. It achieves this by using the selected fixed and composite effects from CEZIJ model but employs a slightly modified covariance matrix $\hat{\Sigma}$ where the covariances between random effects originating from the different submodels are set to 0, thus representing a setup where the outcomes are not modeled jointly.

The out-of-sample validation requires predicting the responses dynamically over time. For Benchmarks I and III this step is easily carried out by running the fitted model on the validation data. However, for Benchmark II, IV and CEZIJ model the prediction mechanism must, respectively, estimate the latent random effects and appropriately account for the endogenous nature of the responses. To do that we utilize the simulation scheme discussed in Rizopoulos (2012, sec. 7.2) and Rizopoulos (2011, sec. 3), and calculate the expected time j responses given the observed responses until time $j-1$, the estimated parameters and the event that the player has not churned until time $j-1$ (details provided in Section B of the supplementary materials). Table 3 summarizes the results of predictive performance of CEZIJ and the benchmark models. For AI and EI, Table 3 presents, for each model, the false positive (FP) rate and the false negative (FN) rate, respectively, averaged over the 29 time points. The FP rate measures the percentage of cases where the model incorrectly predicted activity (or engagement) whereas the FN rate measures the percentage of cases where the model incorrectly predicted no activity (or no engagement). Benchmark II, for example, exhibits the lowest FP rate and has the highest FN rate followed by Benchmark III. The low FP rate of Benchmark II, however, belies the relatively poor performance of this model in predicting zero inflated responses which becomes apparent in the higher FN rates especially for the EI model. The CEZIJ model along with Benchmark IV, on the other hand, have the lowest FN rates

demonstrating their relatively superior ability in predicting the zero inflated responses of AI and EI. For positive activity times and positive engagement values we take a slightly different approach and first calculate the time j prediction errors PE_j for the Benchmark models and CEZIJ as follows. For any model $\mathcal{M} \in \{\text{Benchmark I}, \dots, \text{Benchmark IV}, \text{CEZIJ}\}$, we define $PE_j^{\mathcal{M}}$ for submodel $s = 2$ at time $j = 1, \dots, 29$ as

$$PE_j^{\mathcal{M}}(\mathbb{Y}^{*(s)}, \hat{\mathbb{Y}}^{*(s)}) = \sum_{i=1}^n \left| \log \mathbb{Y}_{ij}^{*(s)} - \log \hat{\mathbb{Y}}_{ij}^{*(s)} \right|, \quad (10)$$

where $\mathbb{Y}_{ij}^{*(s)} = \mathbb{Y}_{ij}^{(s)}$ if $\alpha_{ij} = 1$ and 1 otherwise, and $\hat{\mathbb{Y}}_{ij}^{*(s)} = \hat{\mathbb{Y}}_{ij}^{(s)}$ if $\hat{\alpha}_{ij} = 1$ and 1 otherwise with $\hat{\mathbb{Y}}_{ij}^{(s)}, \hat{\alpha}_{ij}$ being model \mathcal{M} predictions of activity time, AI, respectively, for player i at time j . The time j prediction error for submodel $s = 4$ is also defined in a similar fashion with $\alpha_{ij}, \hat{\alpha}_{ij}$ replaced with $\epsilon_{ij}, \hat{\epsilon}_{ij}$, respectively, and measures the total absolute deviation of the prediction from the truth at any time j . For notational convenience the dependence of $PE_j^{\mathcal{M}}$ on $\alpha_{ij}, \hat{\alpha}_{ij}$ (or $\epsilon_{ij}, \hat{\epsilon}_{ij}$) have been suppressed but the inclusion of these predicted and observed indicators in Equation (10) is aimed at exploiting the dependencies between the responses, if any. For the two submodels ($s = 2, 4$), Table 3 presents the ratio of the prediction errors of the Benchmarks to the CEZIJ model averaged over the 29 time points where a ratio in excess of 1 indicates a larger absolute deviation of the prediction from the truth when compared to CEZIJ model. All Benchmark models exhibit prediction error ratios bigger than 1 with Benchmarks II and III being the worse for engagement amount and activity time models, respectively. Benchmark IV, on the other hand, profits from the structure of the various components of CEZIJ model but is unable to account for the dependencies between the responses which is reflected in its prediction error ratios being slightly bigger than 1 but along with the CEZIJ model, it continues to demonstrate superior prediction error ratios across the two submodels.

Table 3. Results of predictive performance of CEZIJ model and Benchmarks I to IV.

Submodel	Benchmark I	Benchmark II	Benchmark III	Benchmark IV	CEZIJ
Activity indicator	1.32%/6.71%	0.27%/7.83%	1.19%/6.33%	5.92%/4.15%	5.86%/4.12%
Total time played	1.742	1.961	4.662	1.041	1
Engagement indicator	0.09%/1.87%	0%/1.89%	0.05%/1.89%	3.56%/1.48%	3.54%/1.47%
Engagement amount	1.408	8.619	1.217	1.067	1

NOTES: For activity and engagement indicators, the false positive (FP) rate/the false negative (FN) rate averaged over the 29 time points are reported. For nonzero activity time and engagement amounts, the ratio of prediction errors (10) of Benchmarks I to IV to CEZIJ model averaged over the 29 time points are reported.

6.3. Player Segmentation Using Predicted Churn Probabilities

Player subpopulations with similar churn characteristics over time provide valuable insights into user profiles that are more likely to dropout and can be used to design future retention policies specifically targeting those characteristics. In this section, we use the fitted churn model of Section 6.1 to predict the temporal trajectories of churn probabilities on a sample of 1000 players who are 30 days into the game and use the predicted probabilities over the next 25 days to cluster the players into homogeneous subgroups. The churn probabilities are predicted in a similar fashion as discussed in Section 6.2 and Section B of the supplementary materials where the churn probability at time j is predicted conditional on the estimated parameters, the observed responses until time $j - 1$ and the event that the player has not churned until time $j - 1$. To determine the player subgroups, we use R package `fda.usc` to cluster the rows of the 1000×25 predicted churn probability matrix using functional K-means clustering. We use the prediction strength algorithm of Tibshirani and Walther (2005) to determine the number of clusters.

In Figure 6, the three cluster centroids segment the sample into groups which demonstrate distinct temporal churn profiles. For instance, cluster 3, which holds almost 48% of the players, exhibits rising churn probabilities until day 5 but tapers down

under the influence promotions I, II, and VI. Cluster 2, with 34% of the players, has a different trajectory than cluster 3 and appears to respond favorably to promotion VI. Of particular importance are those players that belong to cluster 1 which holds 18% of the players and is characterized by rising churn probabilities over time. The churn profile of this cluster represents players who have been relatively inactive in the game and continue to do so even under the influence of various promotion strategies. During days 17 to 19, their churn probabilities are predicted to diminish under the effect of promotion VI however subsequent promotions do not appear to have any favorable impact. These segment curves suggest that there are some key differences in customer attrition patterns. For example, Cluster 1 shows increasing attrition rates over time, which suggests that the game is not able to retain these players. Cluster 2 shows increasing attrition initially, but then the attrition rate starts to decline significantly after 45 days. This segment is potentially beneficial to the platform as it demonstrates that there is a core set of players who are loyal to the game. Players in Cluster 3 on average start with a much higher attrition rate than the other two segments, but their attrition rate tapers down significantly after five days and then stays at a very low level over time. Interestingly, Cluster 3 seems to be responding to promotions I, II, and IV. These differences in user behavior across the segments can be leveraged to increase the efficiency of player

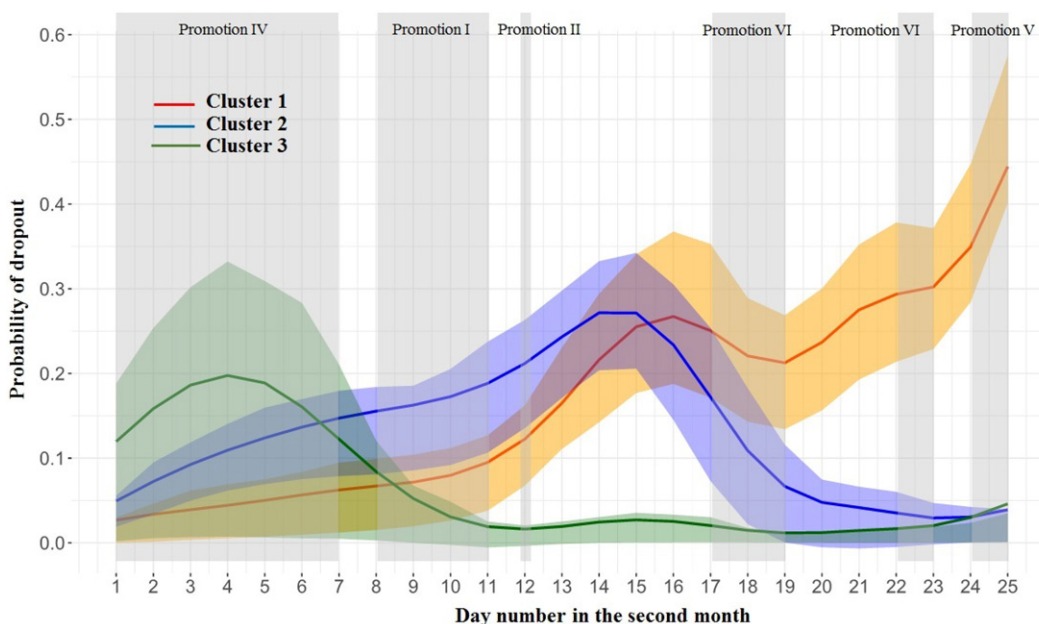


Figure 6. Functional cluster analysis of predicted dropout probabilities over time. The plot presents three cluster centroids. The shaded band around the centroids are the 25th and 75th percentiles of the churn probabilities. The vertical shaded regions in the graph correspond to the days on which different promotion strategies were in effect. The number of clusters were identified using prediction strength (Tibshirani and Walther 2005).

retention policies. They also suggest that the platform should adopt different business strategies. For instance, in Cluster 3 many players have been weeded out early. This indicates that short term visitors to the gaming portal have left the platform more quickly in Cluster 3 compared to the other two segments. So it is important for the platform to emphasize promotional activities that increase player engagement. On the contrary, for Clusters 1 and 2, it is important for the platform to emphasize promotional activities that increase player log-in or activity. This relative emphasis across the segments can increase the efficiency of marketing campaigns.

7. Discussion

We propose a very scalable joint modeling framework CEZIJ for unified inference and prediction of player activity and engagement in freemium mobile games. The rapid growth of mobile games globally has generated significant research interest in different business areas such as marketing, management, and information sciences. Our proposed algorithm conducts variable selection by maintaining the hierarchical congruity of the fixed and random effects and produces models with interpretable composite effects. A key feature of our framework is that it allows incorporation of side information and domain expertise through convexity constraints. We exhibit the superior performance of CEZIJ in producing dynamic predictions. It is also used to segment players based on their churn rates, with the analysis revealing several idiosyncratic player behaviors that can be used for targeted marketing of players in future freemium games. The segmentation findings have important business implications for monetization of the platforms. They can be used to enhance the effectiveness and efficiency of promotional activities and also future user acquisition and retention strategies.

Our inferential framework is based on modern optimization techniques and is very flexible. It can be used in a wide range of big-data applications that need analyzing multiple high-dimensional longitudinal outcomes along with a time-to-event analysis. In future, we would like to extend our joint modeling program for providing comprehensive statistical guidance regarding the growth, development and optimal pricing of generic digital products that use the freemium model. For that purpose, it will be interesting to investigate extensions of our CEZIJ modeling framework, in particular, the possibility of incorporating nonparametric components for modeling the nonlinear time effects since player behavior may change over time. Furthermore, the current dropout model in Equation (8) may be enhanced to include more sophisticated structures involving cumulative effects parameterization and conduct variable selection on the high-dimensional vector of association parameter η , which the current CEZIJ framework implicitly achieves through the selection of the random effects. An alternative and computationally less demanding approach may be to consider the following low dimensional representation wherein the dropout model is of the form $\text{logit}(\lambda_{ij}) = \mathbf{x}_{ij}^{(s)T} \boldsymbol{\beta}^{(s)} + \sum_{s=1}^4 \eta_s \mathbf{z}_{ij}^{(s)T} \mathbf{b}_i^{(s)}$ so that η is then only a 4×1 vector. Finally, while the focus of this article is the CEZIJ modeling framework

and its applicability in the disciplined study of freemium behavior and other applications that needs analyzing multiple high-dimensional longitudinal outcomes along with a time-to-event analysis, a natural extension of our work, as future research, will be targeted toward estimating standard errors of the estimated coefficients and confidence intervals under the CEZIJ framework using ideas from recent developments in post-selection inference (see, e.g., Javanmard and Montanari 2014; Lee et al. 2016).

Of the thousands of freemium games that are developed every month, very few of them go on to make adequate amount through IAP. Most games resemble our data where a significant part of the revenue is earned through in-game ads and social media usages. In these games, such low incidence of real money purchases present a challenge in model development as the robustness of the estimated model coefficients will be significantly impacted in case real money purchases are modeled as a separate response variable. Thus, in very low IAP incidence games it is useful to model the combined revenue using game specific weights to blend direct and indirect engagement as is done in this article. For games with significant amount of IAP, we envision modeling direct and indirect engagement separately and study their interactions.

Supplementary Materials

The supplementary materials contain the following items: details around the maximization problem in equation (9), the prediction equations used in section 6.3, discussion around the split-and-conquer approach, data description and variable selection voting results.

Acknowledgments

We thank Professors Gil Appel and Milan Miric for helpful discussions and the editor, the associate editor, and two anonymous referees for several constructive suggestions that have greatly improved the presentation of the article.

Funding

The research here was partially supported by NSF DMS-1811866, Indian Institute of Management-Bangalore Challenge grant and the Zumberge individual award from the University of Southern California.

References

- Agelle, P. (2016), "Getting Gacha Right: Tips for Creating Successful In-Game Lotteries," *Pocket Gamer*, available at <http://www.pocketgamer.biz/comment-and-opinion/63620/getting-gacha-right-tips-for-creating-successful-in-game-lotteries/>. [549]
- Alfo, M., Maruotti, A., and Trovato, G. (2011), "A Finite Mixture Model for Multivariate Counts Under Endogenous Selectivity," *Statistics and Computing*, 21, 185–202. [540]
- AppBrain (2017), "Free vs. Paid Android Apps," *AppBrain*, available at <http://www.appbrain.com/stats/free-and-paid-android-applications>. [538]
- Appel, Gil and Libai, Barak and Muller, Eitan and Shachar, Roni, On the Monetization of Mobile Apps (January 20, 2019). Available at SSRN: <https://ssrn.com/abstract=2652213> or <http://dx.doi.org/10.2139/ssrn.2652213> [538]
- Banerjee, S., Carlin, B. P., and Gelfand, A. E. (2014), *Hierarchical Modeling and Analysis for Spatial Data*, Boca Raton, FL: CRC Press. [539]

- Banerjee, T., Mukherjee, G., and Sun, W. (2018), "Adaptive Sparse Estimation With Side Information," arXiv no. 1811.11930. [540]
- Bondell, H. D., Krishna, A., and Ghosh, S. K. (2010), "Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effects Models," *Biometrics*, 66, 1069–1077. [540,546]
- Bordeaux, Kevin and Jeppesen, Lars Bo and Miric, Milan, Competing on Free(mium): Digital Competition with Network Effects (February 22, 2019). Available at SSRN: <https://ssrn.com/abstract=2984546> or <http://dx.doi.org/10.2139/ssrn.2984546> [538]
- Breslow, N. E., and Clayton, D. G. (1993), "Approximate Inference in Generalized Linear Mixed Models," *Journal of the American Statistical Association*, 88, 9–25. [546,549]
- Candes, E. J., Wakin, M. B., and Boyd, S. P. (2008), "Enhancing Sparsity by Reweighted l_1 Minimization," *Journal of Fourier Analysis and Applications*, 14, 877–905. [546]
- Chen, X., and Xie, M.-G. (2014), "A Split-and-Conquer Approach for Analysis of Extraordinarily Large Data," *Statistica Sinica*, 24, 1655–1684. [540,547]
- Diele, O. (2013), "State of Online Gaming Report," *Spil Games*, available at http://auth-83051f68-ec6c-44e0-afe5-bd8902acff57.cdn.spilcloud.com/v1/archives/1384952861.25_State_of_Gaming_2013_US_FINAL.pdf. [538]
- Fan, Y., and Li, R. (2012), "Variable Selection in Linear Mixed Effects Models," *Annals of Statistics*, 40, 2043. [540,546]
- Garg, R., and Telang, R. (2013), "Inferring app Demand from Publicly Available Data." *MIS Quarterly*, 37(4), 1253–1264. [538]
- Greene, W. (2009), "Models for Count Data With Endogenous Participation," *Empirical Economics*, 36, 133–173. [540]
- Guo, X., and Carlin, B. P. (2004), "Separate and Joint Modeling of Longitudinal and Event Time Data Using Standard Computer Packages," *The American Statistician*, 58, 16–24. [545]
- Gustafsson, A., Johnson, M. D., and Roos, I. (2005), "The Effects of Customer Satisfaction, Relationship Commitment Dimensions, and Triggers on Customer Retention," *Journal of Marketing*, 69, 210–218. [539]
- Haans, H., Raassens, N., and van Hout, R. (2013), "Search Engine Advertisements: The Impact of Advertising Statements on Click-Through and Conversion Rates," *Marketing Letters*, 24, 151–163. [540]
- Han, C., and Kronmal, R. (2006), "Two-Part Models for Analysis of Agatston Scores With Possible Proportionality Constraints," *Communications in Statistics—Theory and Methods*, 35, 99–111. [540]
- Hatfield, L. A., Boye, M. E., Hackshaw, M. D., and Carlin, B. P. (2012), "Multilevel Bayesian Models for Survival Times and Longitudinal Patient-Reported Outcomes With Many Zeros," *Journal of the American Statistical Association*, 107, 875–885. [540]
- Hui, F. K., Müller, S., and Welsh, A. (2017a), "Hierarchical Selection of Fixed and Random Effects in Generalized Linear Mixed Models," *Statistica Sinica*, 27, 501–518. [540,546]
- (2017b), "Joint Selection in Mixed Models Using Regularized PQL," *Journal of the American Statistical Association*, 112, 1323–1333. [540,549]
- Hui, S. K., Inman, J. J., Huang, Y., and Suher, J. (2013), "The Effect of In-Store Travel Distance on Unplanned Spending: Applications to Mobile Promotion Strategies," *Journal of Marketing*, 77, 1–16. [540]
- Hwong, C. (2016), "Using Audience Measurement Data to Boost User Acquisition and Engagement," *Verto Analytics*, available at <http://www.vertoanalytics.com/report-leveling-mobile-game-using-audience-measurement-data-boost-user-acquisition-engagement/>. [538]
- Ibrahim, J. G., Zhu, H., Garcia, R. I., and Guo, R. (2011), "Fixed and Random Effects Selection in Mixed Effects Models," *Biometrics*, 67, 495–503. [546]
- James, G. M., Paulson, C., and Rusmevichientong, P. (2013), "Penalized and Constrained Regression," Technical Report. [540,547]
- Javanmard, A., and Montanari, A. (2014), "Confidence Intervals and Hypothesis Testing for High-Dimensional Regression," *The Journal of Machine Learning Research*, 15, 2869–2909. [552]
- Jerath, K., Fader, P. S., and Hardie, B. G. (2011), "New Perspectives on Customer "Death" Using a Generalization of the Pareto/NBD Model," *Marketing Science*, 30, 866–880. [539]
- Jiang, J. (2007), *Linear and Generalized Linear Mixed Models and Their Applications*, New York: Springer Science & Business Media. [539]
- Jordan, M. I. (2013), "On Statistics, Computation and Scalability," *Bernoulli*, 19, 1378–1390. [540]
- Kanerva, T. (2016), "Cultures Combined: Japanese Gachas Are Sweeping F2P Mobile Games in the West," *GameRefinery*, available at <http://www.gamerefinery.com/japanese-gachas-sweeping-f2p-games-west/>. [541]
- Koetsier, J. (2015), "Why 2016 Is the Global Tipping Point for the Mobile Economy," *Tune*, available at <https://www.tune.com/blog/global-mobile-why-2016-is-the-global-tipping-point-for-the-mobile-economy/>. [538]
- Kumar, V. (2014), "Making 'freemium' Work," *Harvard Business Review*, 92, 27–29. [539]
- Lee, J. D., Sun, D. L., Sun, Y., and Taylor, J. E. (2016), "Exact Post-selection Inference, With Application to the Lasso," *The Annals of Statistics*, 44, 907–927. [552]
- Lee, J. D., Sun, Y., Liu, Q., and Taylor, J. E. (2015), "Communication-Efficient Sparse Regression: A One-Shot Approach," arXiv no. 1503.04337. [540]
- Lin, B., Pang, Z., and Jiang, J. (2013), "Fixed and Random Effects Selection by REML and Pathwise Coordinate Optimization," *Journal of Computational and Graphical Statistics*, 22, 341–355. [540,546]
- Liu, C. Z., Au, Y. A., and Choi, H. S. (2014), "Effects of Freemium Strategy in the Mobile App Market: An Empirical Study of Google Play," *Journal of Management Information Systems*, 31, 326–354. [538]
- Lu, C., Lin, Z., and Yan, S. (2015), "Smoothed Low Rank and Sparse Matrix Recovery by Iteratively Reweighted Least Squares Minimization," *IEEE Transactions on Image Processing*, 24, 646–654. [546]
- MarketingCharts (2017), "App Retention Rates Still Low, but Improving," *Marketing Charts*, available at <http://www.marketingcharts.com/online/app-retention-rates-still-low-but-improving-75135/>. [538]
- McCulloch, C. E. (1997), "Maximum Likelihood Algorithms for Generalized Linear Mixed Models," *Journal of the American Statistical Association*, 92, 162–170. [546]
- (2008), "Joint Modelling of Mixed Outcome Types Using Latent Variables," *Statistical Methods in Medical Research*, 17, 53–73. [539]
- McDonald, E. (2017), "The Global Games Market," *Newzoo*, available at <https://newzoo.com/insights/articles/the-global-%20games-market-%20will-reach-108-%20billion-%20in-2017-%20with-mobile-%20taking-42>. [538]
- Michael I. Jordan, Jason D. Lee & Yun Yang (2018) "Communication-Efficient Distributed Statistical Inference," *Journal of the American Statistical Association*, <https://dx.doi.org/10.1080/01621459.2018.1429274> [540]
- Min, Y., and Agresti, A. (2005), "Random Effect Models for Repeated Measures of Zero-Inflated Count Data," *Statistical Modelling*, 5, 1–19. [540]
- Needleman, S. E., and Loten, A. (2012), "When Freemium Fails," *WSJ*, available at <https://www.wsj.com/articles/SB10000872396390443713704577603782317318996>. [539]
- Niculescu, M. F., and Wu, D. J. (2011), "When Should Software Firms Commercialize New Products via Freemium Business Models" (working paper) available at <https://www.misrc.umn.edu/wise/papers/3b-1.pdf> [538]
- Olsen, M. K., and Schafer, J. L. (2001), "A Two-Part Random-Effects Model for Semicontinuous Longitudinal Data," *Journal of the American Statistical Association*, 96, 730–745. [539]
- Pan, J., and Huang, C. (2014), "Random Effects Selection in Generalized Linear Mixed Models via Shrinkage Penalty Function," *Statistics and Computing*, 24, 725–738. [546]
- Peng, H., and Lu, Y. (2012), "Model Selection in Linear Mixed Effect Models," *Journal of Multivariate Analysis*, 109, 109–129. [546]
- Perro, J. (2016), "Mobile Apps: What's a Good Retention Rate?," *Localytics*, available at <http://Info.Localytics.Com/Blog/Mobile-apps-Whats-A-Good-Retention-Rate>. [538]
- Pocket Gamer (2018), "Number of Applications Submitted per Month to the iTunes App Store," *Pocket Gamer*, available at <http://www.pocketgamer.biz/metrics/app-store/submissions/>. [538]
- Rabe-Hesketh, S., Skrondal, A., and Pickles, A. (2002), "Reliable Estimation of Generalized Linear Mixed Models Using Adaptive Quadrature," *The Stata Journal*, 2, 1–21. [546]

- Rizopoulos, D. (2011), "Dynamic Predictions and Prospective Accuracy in Joint Models for Longitudinal and Time-to-Event Data," *Biometrics*, 67, 819–829. [550]
- (2012), *Joint Models for Longitudinal and Time-to-Event Data: With Applications in R*, Boca Raton, FL: CRC Press. [539,546,550]
- Rizopoulos, D., and Lesaffre, E. (2014), "Introduction to the Special Issue on Joint Modelling Techniques," *Statistical Methods in Medical Research*, 23, 3–10. [539]
- Rizopoulos, D., Verbeke, G., and Lesaffre, E. (2009), "Fully Exponential Laplace Approximations for the Joint Modelling of Survival and Longitudinal Data," *Journal of the Royal Statistical Society, Series B*, 71, 637–654. [539]
- Rizopoulos, D., Verbeke, G., and Molenberghs, G. (2010), "Multiple-Imputation-Based Residuals and Diagnostic Plots for Joint Models of Longitudinal and Survival Outcomes," *Biometrics*, 66, 20–29. [539]
- Schellendorfer, J., Meier, L., and Bühlmann, P. (2014), "Glmmlasso: An Algorithm for High-Dimensional Generalized Linear Mixed Models Using 1-Penalization," *Journal of Computational and Graphical Statistics*, 23, 460–477. [549]
- Statista (2018), "How Many Hours in a Typical Week Would You Say You Play Games?," *Statista*, available at <https://www.statista.com/statistics/273311/time-spent-gaming-weekly-in-the-uk-by-age/>. [538]
- Swrve (2016), "Monetization Report 2016," *swrve*, available at <https://www.swrve.com/images/uploads/whitepapers/swrve-monetization-report-2016.pdf>. [539]
- Taube, A. (2013), "People Spend Way More on Purchases in Free Apps Than They Do Downloading Paid Apps," *Business Insider*, available at <http://www.businessinsider.com/inapp-purchases-dominate-revenue-share-2013-12>. [538]
- Tibshirani, R., and Walther, G. (2005), "Cluster Validation by Prediction Strength," *Journal of Computational and Graphical Statistics*, 14, 511–528. [551]
- Tierney, L., and Kadane, J. B. (1986), "Accurate Approximations for Posterior Moments and Marginal Densities," *Journal of the American Statistical Association*, 81, 82–86. [546]
- Toto, S. (2012), "Gacha: Explaining Japan's Top Money-Making Social Game Mechanism," *Kantan Games*, available at <https://www.serkantoto.com/2012/02/21/gacha-social-games/>. [541]
- Urban, G. L., Liberali, G., MacDonald, E., Bordley, R., and Hauser, J. R. (2013), "Morphing Banner Advertising," *Marketing Science*, 33, 27–46. [540]
- Vonesh, E. F., Greene, T., and Schluchter, M. D. (2006), "Shared Parameter Models for the Joint Analysis of Longitudinal Data and Event Times," *Statistics in Medicine*, 25, 143–163. [545]
- Wang, H. (2014), "Coordinate Descent Algorithm for Covariance Graphical Lasso," *Statistics and Computing*, 24, 521–529. [547]
- Wei, G. C., and Tanner, M. A. (1990), "A Monte Carlo Implementation of the EM Algorithm and the Poor Man's Data Augmentation Algorithms," *Journal of the American Statistical Association*, 85, 699–704. [540,541,547]
- Yang, Z., and Peterson, R. T. (2004), "Customer Perceived Value, Satisfaction, and Loyalty: The Role of Switching Costs," *Psychology & Marketing*, 21, 799–822. [539]
- Zhao, Y.-B., and Kočvara, M. (2015), "A New Computational Method for the Sparsest Solutions to Systems of Linear Equations," *SIAM Journal on Optimization*, 25, 1110–1134. [546]

Supplementary Material to “A Large-scale Constrained Joint Modeling Approach For Predicting User Activity, Engagement And Churn With Application To Freemium Mobile Games”

Trambak Banerjee*, Gourab Mukherjee*, Shantanu Dutta[†] and Pulak Ghosh[‡]

January 15, 2019[§]

This supplementary material holds the following items: details around the maximization problem in equation (9) (section A), the prediction equations used in section 6.3 of the main paper (section B), discussion around the split-and-conquer approach (section C), data description (section D) and variable selection voting results (section E).

A Technical details around the maximization problem in equation (9)

In this section, we will first show that the maximization problem in equation (9) decouples into separate components that estimate $\beta^{(s)}$ (and σ_1, σ_2 for the activity and engagement models) and Σ as solutions to independent optimization problems (section A.1). Thereafter, we show that the optimization problems involving $\beta^{(s)}$ are convex and can be solved after reducing the original problem to an ℓ_1 penalized least squares fit with convex constraints (section A.2), while the coordinate descent algorithm of Wang (2014) provides a

[§]The research here was partially supported by NSF DMS-1811866 and by the Zumberge individual award from the University of Southern California’s James H. Zumberge Faculty Research and Innovation Fund. Corresponding author: gmukherj@marshall.usc.edu

*Department of Data Sciences and Operations, University of Southern California, Los Angeles

[†]Department of Marketing, University of Southern California, Los Angeles

[‡]Department of Decision Sciences and Information Systems, Indian Institute of Management, Bangalore

solution to the non-convex problem involving Σ (section A.3).

A.1 Simplifying equation (9)

Note that in the E-step of section 5, $\ell_i^Q(\Theta)$ is approximated by $\sum_{i=1}^n \sum_{d=1}^D \ell_i^{\text{cl}}(\Theta, \mathbf{b}_i^d) w_{id}^{(t)}$ where

$$w_{id}^{(t)} = p(\alpha_{ij}, \mathbb{A}_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_i^* | \mathbf{b}_i^d, \theta^{(t)}) / \sum_{d=1}^D p(\alpha_{ij}, \mathbb{A}_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_i^* | \mathbf{b}_i^d, \theta^{(t)})$$

is a known constant at iteration (t) and

$$\ell_i^{\text{cl}}(\Theta, \mathbf{b}_i^d) = -\frac{1}{2} \log \Sigma - \frac{1}{2} \mathbf{b}_i^{dT} \Sigma^{-1} \mathbf{b}_i^d + \sum_{j=1}^{m_i} \log p(\alpha_{ij}, \mathbb{A}_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_i^* | \mathbf{b}_i^d, \theta).$$

Moreover given the random effects \mathbf{b}_i^d , $\log p(\alpha_{ij}, \mathbb{A}_{ij}, \epsilon_{ij}, \mathbb{E}_{ij}, \mathbb{D}_i^* | \mathbf{b}_i^d, \theta)$ factorizes into

$$\begin{aligned} & \log p(\alpha_{ij} | \mathbf{b}_i^{d(1)}, \beta^{(1)}) + \log p(\mathbb{A}_{ij} | \alpha_{ij}, \mathbf{b}_i^{d(2)}, \beta^{(2)}, \sigma_1) + \log p(\epsilon_{ij} | \alpha_{ij}, \mathbf{b}_i^{d(3)}, \beta^{(3)}) + \\ & \log p(\mathbb{E}_{ij} | \epsilon_{ij}, \alpha_{ij}, \mathbf{b}_i^{d(4)}, \beta^{(4)}, \sigma_2) + \log p(\mathbb{D}_i^* | \mathbf{b}_i^d, \beta^{(5)}, \eta) \end{aligned}$$

wherein the s^{th} term, for $s = 1, \dots, 5$, in the display above is solely a function of the unknown parameter $\beta^{(s)}$ (and σ_1, σ_2, η for $s = 2, 4, 5$ respectively). This suffices to show that the maximization problem in equation (9) decouples into six separate problems for estimating $\beta^{(1)}$, $(\beta^{(2)}, \sigma_1)$, $\beta^{(3)}$, $(\beta^{(4)}, \sigma_2)$, $(\beta^{(5)}, \eta)$ and Σ .

A.2 Estimating $\beta^{(s)}$

In what follows, we will show that the optimization problems involving $\beta^{(s)}$ (and σ_1, σ_2 for the activity and engagement models) are convex and can be solved after reducing the original problem to an ℓ_1 penalized least squares fit with convex constraints.

AI model - First note that $\sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log p(\alpha_{ij} | \mathbf{b}_i^{d(1)}, \beta^{(1)}) w_{id}^{(t)}$ can be written as $f_1(\beta^{(1)}) + f_2(\beta^{(1)}) +$ terms independent of $\beta^{(1)}$ where

$$f_1(\beta^{(1)}) = - \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log \left[1 + \exp \left(\mathbf{x}_{ij}^{(1)T} \beta^{(1)} + \mathbf{z}_{ij}^{(1)T} \mathbf{b}_i^{d(1)} \right) \right] w_{id}^{(t)}$$

is concave in $\boldsymbol{\beta}^{(1)}$ and

$$f_2(\boldsymbol{\beta}^{(1)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \mathbf{x}_{ij}^{(1)T} \boldsymbol{\beta}^{(1)} w_{id}^{(t)}$$

is affine in $\boldsymbol{\beta}^{(1)}$. Now from equation (9), the minimization problem for $\boldsymbol{\beta}^{(1)}$ is

$$\min_{\boldsymbol{\beta}^{(1)}} f(\boldsymbol{\beta}^{(1)}) + h(\boldsymbol{\beta}^{(1)}) \quad (1)$$

where $f(\boldsymbol{\beta}^{(1)}) = -f_1(\boldsymbol{\beta}^{(1)}) - f_2(\boldsymbol{\beta}^{(1)})$ is convex and differentiable with respect to $\boldsymbol{\beta}^{(1)}$, and $h(\boldsymbol{\beta}^{(1)}) = n\lambda \sum_{r=1}^p c_{1r} |\beta_{1r}| + \mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(1)})$ is convex but non-differentiable, with $\mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(1)})$ as the indicator function of the closed, convex set $\mathcal{C} = \{\boldsymbol{\beta}^{(1)} : \mathbf{f}^{(1)}(\boldsymbol{\beta}^{(1)}) \leq 0\}$. To solve equation (1), we use the proximal gradient method that updates $\boldsymbol{\beta}^{(1)}$ in iteration $k = 1, 2, 3, \dots$ as

$$\boldsymbol{\beta}_{(k)}^{(1)} = \text{prox}_{t_k, h} \left(\boldsymbol{\beta}_{(k-1)}^{(1)} - t_k \nabla f(\boldsymbol{\beta}_{(k-1)}^{(1)}) \right) \quad (2)$$

where $t_k > 0$ is the step size determined by backtracking line search and

$$\text{prox}_{t_k, h}(\mathbf{u}) = \arg \min_{\boldsymbol{\beta}} \left(h(\boldsymbol{\beta}) + \frac{1}{2t_k} \|\boldsymbol{\beta} - \mathbf{u}\|_2^2 \right) \quad (3)$$

is the proximal mapping of h with $\mathbf{u} = \boldsymbol{\beta}_{(k-1)}^{(1)} - t_k \nabla f(\boldsymbol{\beta}_{(k-1)}^{(1)})$ and

$$\nabla f(\boldsymbol{\beta}_{(k-1)}^{(1)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D w_{id}^{(t)} \left\{ \left[1 + \exp \left(-\mathbf{x}_{ij}^{(1)T} \boldsymbol{\beta}_{(k-1)}^{(1)} - \mathbf{z}_{ij}^{(1)T} \mathbf{b}_i^{d(1)} \right) \right]^{-1} - \alpha_{ij} \right\} \mathbf{x}_{ij}^{(1)}$$

being the derivative of $f(\boldsymbol{\beta}^{(1)})$ with respect to $\boldsymbol{\beta}^{(1)}$ evaluated at $\boldsymbol{\beta}_{(k-1)}^{(1)}$. The proximal mapping in equation (3) for our specific application is, unfortunately, not available in an analytical form. We resort to computing the proximal mappings numerically by re-writing the minimization problem in equation (3) as an ℓ_1 penalized least squares fit with convex constraints as follows:

$$\begin{aligned} \min_{\tilde{\boldsymbol{\beta}}^{(1)}} \frac{1}{2t} \|\mathbf{u} - \mathbf{A}^{(1)} \tilde{\boldsymbol{\beta}}^{(1)}\|_2^2 + n\lambda \|\tilde{\boldsymbol{\beta}}^{(1)}\|_1 \\ \text{subject to } \tilde{\mathbf{f}}^{(1)}(\tilde{\boldsymbol{\beta}}^{(1)}) \leq 0 \end{aligned} \quad (4)$$

where $t = t_k$, $\tilde{\beta}_{1r} = c_{1r} \beta_{1r}$, $\mathbf{A}^{(1)}$ is a $p \times p$ diagonal matrix with $A_{r,r}^{(1)} = 1/c_{1r}$ and $\tilde{\mathbf{f}}^{(1)}$ are the transformed convexity constraints on $\tilde{\boldsymbol{\beta}}^{(1)}$. For instance, if $\mathbf{f}^{(1)}(\boldsymbol{\beta}^{(1)}) = \mathbf{C}^{(1)} \boldsymbol{\beta}^{(1)}$ for some

matrix $\mathbf{C}^{(1)}$ with p columns then $\tilde{\mathbf{f}}^{(1)}(\tilde{\boldsymbol{\beta}}^{(1)}) = \mathbf{C}^{(1)} \mathbf{A}^{(1)} \tilde{\boldsymbol{\beta}}^{(1)}$. Finally, we solve (4) using CVX (Grant et al., 2008).

Activity Time model - Define $\tau_1 = \sigma_1^{-1}$, $\bar{\boldsymbol{\beta}}^{(2)} = \tau_1 \boldsymbol{\beta}^{(2)}$ and re-write

$\sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log p(\mathbf{A}_{ij} | \alpha_{ij}, \mathbf{b}_i^{d(2)}, \boldsymbol{\beta}^{(2)}, \sigma_1) w_{id}^{(t)}$ as $-f(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) + \text{constant terms}$, where

$$f(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} w_{id}^{(t)} \left[\frac{1}{2} \left(\tau_1 \log \mathbf{A}_{ij} - \mathbf{x}_{ij}^{(2)T} \bar{\boldsymbol{\beta}}^{(2)} \right)^2 - \log \tau_1 \right]$$

is convex in $(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$. Thus from equation (9), the minimization problem for $(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ is

$$\min_{\tau_1, \bar{\boldsymbol{\beta}}^{(2)}} f(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) + h(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) \quad (5)$$

where $f(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ is convex and differentiable with respect to $(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$, and $h(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) = n\lambda \sum_{r=1}^p c_{2r} |\bar{\beta}_{2r}| + \mathbb{I}_{\mathcal{C}}(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ is convex but non-differentiable, with $\mathbb{I}_{\mathcal{C}}(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ as the indicator function of the closed, convex set $\mathcal{C} = \{(\tau_1, \bar{\boldsymbol{\beta}}^{(2)}) : \mathbf{f}^{(2)}(\bar{\boldsymbol{\beta}}^{(2)}) \leq 0, \tau_1 \geq v\}$. Here v is a small positive number used to enforce $\tau_1 > 0$. To solve (5) we use the proximal gradient method discussed in equations (2) and (3) wherein the proximal mapping of h is given by

$$\text{prox}_{t_k, h}(\mathbf{u}) = \arg \min_{\tau_1, \bar{\boldsymbol{\beta}}} \left(h(\tau_1, \bar{\boldsymbol{\beta}}) + \frac{1}{2t_k} \|(\tau_1, \bar{\boldsymbol{\beta}})^T - \mathbf{u}\|_2^2 \right)$$

where $\mathbf{u} = (\tau_1^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)})^T - t_k \nabla f(\tau_1^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)})$ and $\nabla f(\tau_1^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)})$

$$= \left[\begin{array}{c} \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} w_{id}^{(t)} \left\{ -1/\tau_1^{(k-1)} + \log \mathbf{A}_{ij} \left(\tau_1^{(k-1)} \log \mathbf{A}_{ij} - \mathbf{x}_{ij}^{(2)T} \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)} \right) \right\} \\ - \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} w_{id}^{(t)} \left(\tau_1^{(k-1)} \log \mathbf{A}_{ij} - \mathbf{x}_{ij}^{(2)T} \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)} \right) \mathbf{x}_{ij}^{(2)} \end{array} \right]$$

being the derivative of $f(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ with respect to $(\tau_1, \bar{\boldsymbol{\beta}}^{(2)})$ evaluated at $(\tau_1^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(2)})$.

The above proximal mapping is computed in CVX by solving an ℓ_1 penalized least squares fit with convex constraints as shown in equation (4).

EI model - Like the AI model, $\sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log p(\epsilon_{ij} | \alpha_{ij}, \mathbf{b}_i^{d(3)}, \boldsymbol{\beta}^{(3)}) w_{id}^{(t)}$ can be written as $f_1(\boldsymbol{\beta}^{(3)}) + f_2(\boldsymbol{\beta}^{(3)}) + \text{terms independent of } \boldsymbol{\beta}^{(3)}$ where

$$f_1(\boldsymbol{\beta}^{(3)}) = - \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \log \left[1 + \exp \left(\mathbf{x}_{ij}^{(3)T} \boldsymbol{\beta}^{(3)} + \mathbf{z}_{ij}^{(3)T} \mathbf{b}_i^{d(3)} \right) \right] w_{id}^{(t)}$$

is concave in $\boldsymbol{\beta}^{(3)}$ and

$$f_2(\boldsymbol{\beta}^{(3)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \epsilon_{ij} \mathbf{x}_{ij}^{(3)T} \boldsymbol{\beta}^{(3)} w_{id}^{(t)}$$

is affine in $\boldsymbol{\beta}^{(3)}$. So the minimization problem for $\boldsymbol{\beta}^{(3)}$ in (9) is

$$\min_{\boldsymbol{\beta}^{(3)}} f(\boldsymbol{\beta}^{(3)}) + h(\boldsymbol{\beta}^{(3)}) \quad (6)$$

where $f(\boldsymbol{\beta}^{(3)}) = -f_1(\boldsymbol{\beta}^{(3)}) - f_2(\boldsymbol{\beta}^{(3)})$ is convex and differentiable with respect to $\boldsymbol{\beta}^{(3)}$, and $h(\boldsymbol{\beta}^{(3)}) = n\lambda \sum_{r=1}^p c_{3r} |\beta_{3r}| + \mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(3)})$ is convex but non-differentiable, with $\mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(3)})$ as the indicator function of the closed, convex set $\mathcal{C} = \{\boldsymbol{\beta}^{(3)} : \mathfrak{f}^{(3)}(\boldsymbol{\beta}^{(3)}) \leq 0\}$. To solve equation (6), we use the proximal gradient method discussed in equations (2) and (3) wherein the proximal mapping of h is given by

$$\text{prox}_{t_k, h}(\mathbf{u}) = \arg \min_{\boldsymbol{\beta}} \left(h(\boldsymbol{\beta}) + \frac{1}{2t_k} \|\boldsymbol{\beta} - \mathbf{u}\|_2^2 \right)$$

where $\mathbf{u} = \boldsymbol{\beta}_{(k-1)}^{(3)} - t_k \nabla f(\boldsymbol{\beta}_{(k-1)}^{(3)})$ and

$$\nabla f(\boldsymbol{\beta}_{(k-1)}^{(3)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} w_{id}^{(t)} \left\{ \left[1 + \exp \left(-\mathbf{x}_{ij}^{(3)T} \boldsymbol{\beta}_{(k-1)}^{(3)} - \mathbf{z}_{ij}^{(3)T} \mathbf{b}_i^{d(3)} \right) \right]^{-1} - \epsilon_{ij} \right\} \mathbf{x}_{ij}^{(3)}$$

being the derivative of $f(\boldsymbol{\beta}^{(3)})$ with respect to $\boldsymbol{\beta}^{(3)}$ evaluated at $\boldsymbol{\beta}_{(k-1)}^{(3)}$. The above proximal mapping is finally computed in CVX by solving an ℓ_1 penalized least squares fit with convex constraints as shown in equation (4).

Engag. Amount model - Like the Activity time model, define $\tau_2 = \sigma_2^{-1}$, $\bar{\boldsymbol{\beta}}^{(4)} = \tau_2 \boldsymbol{\beta}^{(4)}$ and re-write

$\sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log p(\mathbb{E}_{ij} | \alpha_{ij}, \epsilon_{ij}, \mathbf{b}_i^{d(4)}, \boldsymbol{\beta}^{(4)}, \sigma_2) w_{id}^{(t)}$ as $-f(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) + \text{constant terms}$, where

$$f(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \epsilon_{ij} w_{id}^{(t)} \left[\frac{1}{2} \left(\tau_2 \log \mathbb{E}_{ij} - \mathbf{x}_{ij}^{(4)T} \bar{\boldsymbol{\beta}}^{(4)} \right)^2 - \log \tau_2 \right]$$

is convex in $(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$. Thus from equation (9), the minimization problem for $(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ is

$$\min_{\tau_2, \bar{\boldsymbol{\beta}}^{(4)}} f(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) + h(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) \quad (7)$$

where $f(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ is convex and differentiable with respect to $(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$, and $h(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) = n\lambda \sum_{r=1}^p c_{4r} |\bar{\beta}_{4r}| + \mathbb{I}_{\mathcal{C}}(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ is convex but non-differentiable, with $\mathbb{I}_{\mathcal{C}}(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ as the indicator function of the closed, convex set $\mathcal{C} = \{(\tau_2, \bar{\boldsymbol{\beta}}^{(4)}) : \mathbf{f}^{(4)}(\bar{\boldsymbol{\beta}}^{(4)}) \leq 0, \tau_2 \geq v\}$. Here v is a small positive number used to enforce $\tau_2 > 0$. To solve (7) we use the proximal gradient method discussed in equations (2) and (3) wherein the proximal mapping of h is given by

$$\text{prox}_{t_k, h}(\mathbf{u}) = \arg \min_{\tau_2, \bar{\boldsymbol{\beta}}} \left(h(\tau_2, \bar{\boldsymbol{\beta}}) + \frac{1}{2t_k} \|(\tau_2, \bar{\boldsymbol{\beta}})^T - \mathbf{u}\|_2^2 \right)$$

where $\mathbf{u} = (\tau_2^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)})^T - t_k \nabla f(\tau_2^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)})$ and $\nabla f(\tau_2^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)})$

$$= \left[\begin{array}{c} \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \epsilon_{ij} w_{id}^{(t)} \left\{ -1/\tau_2^{(k-1)} + \log \mathbb{E}_{ij} \left(\tau_2^{(k-1)} \log \mathbb{E}_{ij} - \mathbf{x}_{ij}^{(4)T} \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)} \right) \right\} \\ - \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \alpha_{ij} \epsilon_{ij} w_{id}^{(t)} \left(\tau_2^{(k-1)} \log \mathbb{E}_{ij} - \mathbf{x}_{ij}^{(4)T} \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)} \right) \mathbf{x}_{ij}^{(4)} \end{array} \right]$$

being the derivative of $f(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ with respect to $(\tau_2, \bar{\boldsymbol{\beta}}^{(4)})$ evaluated at $(\tau_2^{(k-1)}, \bar{\boldsymbol{\beta}}_{(k-1)}^{(4)})$. Finally, CVX is used to compute the above proximal mapping by solving an ℓ_1 penalized least squares fit with convex constraints as shown in equation (4).

Dropout model - For the dropout model, we re-write

$\sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log p(\mathbb{D}_i^* | \mathbf{b}_i^d, \boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) w_{id}^{(t)}$ as $f_1(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) + f_2(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) + \text{constant terms}$ where

$$f_1(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) = - \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \log \left[1 + \exp \left(\mathbf{x}_{ij}^{(5)T} \boldsymbol{\beta}^{(5)} + \boldsymbol{\eta}^T \mathbf{b}_i^d \right) \right] w_{id}^{(t)}$$

is concave in $(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$ and

$$f_2(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D \delta_i^{\mathbb{D}} \left(\mathbf{x}_{ij}^{(5)T} \boldsymbol{\beta}^{(5)} + \boldsymbol{\eta}^T \mathbf{b}_i^d \right) w_{id}^{(t)}$$

is affine in $(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$. Now from equation (9), the minimization problem for $(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$ is

$$\min_{\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}} f(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) + h(\boldsymbol{\beta}^{(5)}) \quad (8)$$

where $f(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) = -f_1(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta}) - f_2(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$ is convex and differentiable with respect to $(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$, and $h(\boldsymbol{\beta}^{(5)}) = n\lambda \sum_{r=1}^p c_{5r} |\beta_{5r}| + \mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(5)})$ is convex but non-differentiable, with $\mathbb{I}_{\mathcal{C}}(\boldsymbol{\beta}^{(5)})$ as the indicator function of the closed, convex set $\mathcal{C} = \{\boldsymbol{\beta}^{(5)} : \mathbf{f}^{(5)}(\boldsymbol{\beta}^{(5)}) \leq 0\}$. To

solve equation (8), we use the proximal gradient method discussed in equations (2) and (3) wherein the proximal mapping of h is given by

$$\text{prox}_{t_k, h}(\mathbf{u}) = \arg \min_{\boldsymbol{\beta}, \boldsymbol{\eta}} \left(h(\boldsymbol{\beta}) + \frac{1}{2t_k} \|(\boldsymbol{\beta}, \boldsymbol{\eta})^T - \mathbf{u}\|_2^2 \right)$$

where $\mathbf{u} = (\boldsymbol{\beta}_{(k-1)}^{(5)}, \boldsymbol{\eta}_{(k-1)})^T - t_k \nabla f(\boldsymbol{\beta}_{(k-1)}^{(5)}, \boldsymbol{\eta}_{(k-1)})$ and $\nabla f(\boldsymbol{\beta}_{(k-1)}^{(5)}, \boldsymbol{\eta}_{(k-1)})$

$$= \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{d=1}^D w_{id}^{(t)} \left\{ \left[1 + \exp \left(-\mathbf{x}_{ij}^{(5)T} \boldsymbol{\beta}_{(k-1)}^{(5)} - \boldsymbol{\eta}_{(k-1)} \mathbf{b}_i^d \right) \right]^{-1} - \delta_i^{\mathbb{D}} \right\} \begin{bmatrix} \mathbf{x}_{ij}^{(5)} \\ \mathbf{b}_i^d \end{bmatrix}$$

being the derivative of $f(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$ with respect to $(\boldsymbol{\beta}^{(5)}, \boldsymbol{\eta})$ evaluated at $(\boldsymbol{\beta}_{(k-1)}^{(5)}, \boldsymbol{\eta}_{(k-1)})$. We use CVX to compute the above proximal mapping by solving an ℓ_1 penalized least squares fit with convex constraints as shown in equation (4).

A.3 Estimating Σ

From equation (9), the optimization problem for estimating Σ in iteration (t) can be expressed as

$$\min_{\Sigma > 0} \log |\Sigma| + \text{trace}(\mathbf{Q}\Sigma^{-1}) + 2\lambda \|\mathbf{P} * \Sigma\|_1 \quad (9)$$

where $\mathbf{Q}_{4p_c \times 4p_c} = n^{-1} \sum_{i=1}^n \sum_{d=1}^D \mathbf{b}_i^d \mathbf{b}_i^{dT} w_{id}^{(t)}$, $\mathbf{P}_{4p_c \times 4p_c} = \text{diag}(d_{s1}^{(t)}, \dots, d_{s4p_c}^{(t)})$. Here $*$ denotes elementwise multiplication and for any matrix \mathbf{A} , $\|\mathbf{A}\|_1 = \|\text{vec}(\mathbf{A})\|_1 = \sum_{i,j} |A_{ij}|$. The above minimization problem in Σ is non-convex (Bien and Tibshirani, 2011) and we use the coordinate descent based algorithm of Wang (2014) that updates Σ one row and one column at a time while keeping the remaining elements fixed to obtain a solution. In particular, given inputs $(\mathbf{Q}, \mathbf{P}, \lambda)$ and iteration $(k+1)$, the aforementioned algorithm first partitions

$$\Sigma^{(k+1)} = \begin{pmatrix} \Sigma_{11}^{(k)} & \sigma_{12} \\ \sigma_{12}^T & \sigma_{22} \end{pmatrix}, \quad \mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{q}_{12} \\ \mathbf{q}_{12}^T & q_{22} \end{pmatrix}.$$

where $\Sigma_{11}^{(k)}$ and \mathbf{Q}_{11} are the sub-matrices obtained from the first $4p_c - 1$ columns. Then with $\boldsymbol{\beta} = \sigma_{12}$ and $\gamma = \sigma_{22} - \sigma_{12}^T \Sigma_{11}^{-1(k)} \sigma_{12}$, it uses coordinate descent algorithms (Friedman et al., 2007) to obtain the estimates $(\hat{\boldsymbol{\beta}}, \hat{\gamma})$ (see equations (5)-(7) in Wang (2014)) and

finally updates $\sigma_{12}^{(k+1)} = \hat{\beta}$ and $\sigma_{22}^{(k+1)} = \hat{\gamma} + \hat{\beta}^T \Sigma_{11}^{-1(k)} \hat{\beta}$. This procedure is repeated for every row and column (keeping others fixed) until convergence.

B Prediction equations

We first focus on the prediction problem discussed in section 6.3 of the main paper. For player i , let $\mathcal{Y}_i(t) = \{\alpha_{ij}, \mathbf{A}_{ij}, \epsilon_{ij}, \mathbb{E}_{ij} : 0 \leq j \leq t\}$ denote the observed responses until time t and $\rho_i(u | t)$ be the conditional probability of drop-out at time $u > t > 0$ given no drop-out until time t . Then,

$$\begin{aligned} \rho_i(u | t) &= \Pr(\mathbb{D}_i^* = u | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta) \\ &= \int \Pr(\mathbb{D}_i^* = u | \mathbb{D}_i^* > t, \mathcal{Y}_i(t), \mathbf{b}_i; \Theta) p(\mathbf{b}_i | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta) d\mathbf{b}_i \\ &= \int \Pr(\mathbb{D}_i^* = u | \mathbb{D}_i^* > t, \mathbf{b}_i; \Theta) p(\mathbf{b}_i | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta) d\mathbf{b}_i \end{aligned}$$

Following section 3 of Rizopoulos (2011) and the fitted dropout model in equation (8), an estimate of $\rho_i(u | t)$ is

$$\hat{\rho}_i(u | t) = \Pr(\mathbb{D}_i^* = u | \mathbb{D}_i^* > t, \hat{\mathbf{b}}_i; \hat{\Theta})$$

where $\hat{\mathbf{b}}_i = \arg \max_{\mathbf{b}} \log p(\mathbf{b} | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \hat{\Theta})$.

In section 6.2, we are interested in predicting the time $u > t$ expected longitudinal outcomes of AI, Activity, EI and Engagement given the observed responses $\mathcal{Y}_i(t)$ for player i who has not dropped-out at time t . We consider the case of predicting $w_i(u | t) := \mathbb{E}\{\mathbf{A}_{iu} | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta\}$ as an example as the rest follow along similar lines. Let $\hat{\alpha}_{iu}$ be the predicted AI at time u conditional on $\mathcal{Y}_i(t)$ and no dropout until time t . Then note that

$$\mathbb{E}\{\mathbf{A}_{iu} | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta\} = \int \mathbb{E}\{\mathbf{A}_{iu} | \mathbf{b}_i; \Theta\} p(\mathbf{b}_i | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \Theta) d\mathbf{b}_i$$

and from section 7.2 of Rizopoulos (2012) an estimate of $w_i(u | t)$ is given by

$$\hat{w}_i(u | t) = \begin{cases} 0, & \text{if } \hat{\alpha}_{iu} = 0 \\ \exp\left(\mathbf{x}_{iu}^{(2)T} \hat{\beta}^{(2)} + \mathbf{z}_{iu}^{(2)T} \hat{\mathbf{b}}_i^{(2)} + \frac{\hat{\sigma}_1^2}{2}\right), & \text{otherwise} \end{cases}$$

where $\hat{\mathbf{b}}_i = (\hat{\mathbf{b}}_i^{(s)} : 1 \leq s \leq 4) = \arg \max_{\mathbf{b}} \log p(\mathbf{b} | \mathbb{D}_i^* > t, \mathcal{Y}_i(t); \hat{\Theta})$.

C Split-and-Conquer Approach and Numerical Experiments

In this section, we first discuss the split-and-conquer approach of [Chen and Xie \(2014\)](#) (section [C.1](#)) and thereafter conduct numerical experiments to demonstrate the applicability of this approach in our GLMM setup (section [C.2](#)).

C.1 Split-and-Conquer approach

To enhance the computational efficiency of the estimation procedure, CEZIJ uses the split-and-conquer approach of [Chen and Xie \(2014\)](#) to split the full set of n players into K non-overlapping groups and conducts variable selection separately in each group by solving K parallel maximization problems represented by equation (9). In the process, our methodology uses data-driven adaptive weights $(c_{sr}, d_{sr}) \in \mathbf{R}_+^2$ in the penalty with weights in any iteration being computed from the solutions of the previous iteration. The selected fixed and random effects are then determined using a majority voting scheme across all the K groups as described in Section 5 of the main paper. In their original article however, [Chen and Xie \(2014\)](#) use this approach in a GLM setup, and conduct selection and thereafter estimation of the selected coefficients, by first solving K penalized likelihood problems (with fixed penalty λ that may vary with K) across the K splits of the data and then averaging across the selected coefficients in each split. Theorem 1 in their paper demonstrates that the estimator so obtained is sign consistent under some regularity conditions and as long as $\log(Kp) = o(n/K)$ where p denotes the number of candidate predictors and n the sample size. Moreover along with Theorem 1, Theorem 2 establishes that this averaged estimator is asymptotically equivalent to the estimator obtained by solving the penalized likelihood problem on the entire data.

While these theoretical results do not directly extend to a GLMM setting, in section [C.2](#) we empirically demonstrate the applicability of the above scheme in selecting fixed and random effects in our setting where data-driven adaptive weights are used in the penalty and variable selection is conducted simultaneously across multiple models. In terms of

computational efficiency, the split-and-conquer approach is efficient in the sense that if an estimation procedure requires $O(n^a p^b)$ computing steps for some $a > 1, b \geq 0$, then the split-and-conquer approach results in an efficiency gain of $O(K^{a-1})$ in computing steps (see theorem 5 in [Chen and Xie \(2014\)](#)). Figure 1 presents a comparison of the computing time for the two simulation settings considered in in section C.2 and demonstrates that in both these settings CEZIJ, through its split-and-conquer approach for variable selection, offers a potential gain in computational efficiency against the conventional and memory intensive approach of running the selection algorithm on the undivided data.

C.2 Numerical Experiments

Here we present numerical experiments that assess the model selection performance of CEZIJ under the longitudinal and Dropout models discussed in Section 3.1 of the main paper. The MATLAB code for these simulation experiments is available at <https://github.com/trambakbanerjee/cezij#what-is-cezij>. We consider two simulation settings as follows:

Simulation setting I - We consider a sample of $n = 500$ players and for each player i , let $\mathbb{X}_i = (\mathbf{X}_{i,1}, \dots, \mathbf{X}_{i,p})$ denote the $m \times p$ matrix of candidate predictors where $\mathbf{X}_{i,k} = (x_{i1k}, \dots, x_{imk})^T$. We fix $m = 30$, $p = 10$ and take $\mathcal{I}_f = \{1, \dots, 8\}$, $\mathcal{I}_c = \{9, 10\}$ so that $p_f = 8$, $p_c = 2$. Thus, the first 8 columns of \mathbb{X}_i represent fixed effects while the last 2 represent composite effects.

The five responses $[\alpha_i, \mathbf{A}_i, \epsilon_i, \mathbf{E}_i, \mathbf{D}_i]$ corresponding to equations (2), (3), (6), (7) and (8) of section 3.1 are generated from the following models: $\text{logit}(\pi_{ij}) = \beta_0^{(1)} + x_{ij1}\beta_1^{(1)} + b_{i1}$, $\mu_{ij} = \beta_0^{(2)} + x_{ij2}\beta_1^{(2)} + b_{i2}$ with $\sigma_1 = 0.5$, $\text{logit}(q_{ij}) = \beta_0^{(3)} + x_{ij3}\beta_1^{(3)} + b_{i3}$, $\gamma_{ij} = \beta_0^{(4)} + x_{ij4}\beta_1^{(4)} + b_{i4}$ with $\sigma_2 = 0.5$ and $\text{logit}(\lambda_{ij}) = \beta_0^{(5)} + x_{ij5}\beta_1^{(5)} + \eta_1 b_{i1} + \dots + \eta_4 b_{i4}$ where the true values of the fixed effect coefficients are: $\boldsymbol{\beta}^{(1)} = (1, -1.5)$, $\boldsymbol{\beta}^{(2)} = (3.5, -2)$, $\boldsymbol{\beta}^{(3)} = (1, -1)$, $\boldsymbol{\beta}^{(4)} = (3, -3)$, $\boldsymbol{\beta}^{(5)} = (-1, 2)$ and, $\boldsymbol{\eta} = (\eta_1, \dots, \eta_4) = (-0.1, 0.2, 0.1, -0.2)$. Thus setting I presents a relatively simple scenario wherein there are no composite effects in the true model. The random effects $\mathbf{b}_i = (b_{i1}, \dots, b_{i4})$ are sampled from $N_4(\mathbf{0}, \boldsymbol{\Sigma})$, independently

for each i , where

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.2 & 0.4 & 0.5 \\ 0.2 & 3 & 0.9 & 0.7 \\ 0.4 & 0.9 & 0.8 & 0.5 \\ 0.5 & 0.7 & 0.5 & 4 \end{pmatrix}$$

Since the CEZIJ framework can incorporate convexity constraints on the fixed effect coefficients, we impose the following sign constraints:

$$\beta_0^{(1)} > 0, \beta_1^{(1)} < 0; \beta_1^{(3)} < 0; \beta_0^{(4)} > 0, \beta_1^{(4)} < 0.$$

Finally, to complete the specification, we sample $(x_{ij1}, \dots, x_{ij4})$ from $N_4(\mathbf{0}, 4\mathbf{I}_4)$, independently for each $i = 1, \dots, n, j = 1, \dots, m$. To ensure that the generated sample contains players that have not churned for at least the first 7 to 10 days, we let $\mathbf{X}_{i.5}$ to be an m dimensional ordered sample from $\text{Unif}(-1, 1)$ so that $\mathbf{X}_{i.5} = (x_{i15} \leq \dots \leq x_{im5})$ and, generate the remaining predictors independently from $\text{Unif}(-1, 1)$. In this respect, $\mathbf{X}_{i.5}$ mimics the variable `timesince` (see section 6.1 and table 3) that gradually increases with m and appears in the fitted Dropout model in table 2 of section 6.1.

Simulation setting II - In this setting, we consider a larger design and fix $n = 2000, m = 30, p = 20$ and, take $\mathcal{I}_f = \{1, 3, 5, \dots, 8, 11, \dots, p\}$, $\mathcal{I}_c = \{2, 4, 9, 10\}$ so that $p_f = 16$ and $p_c = 4$. The five responses $[\alpha_i, \mathbf{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i]$ are generated from the following models: $\text{logit}(\pi_{ij}) = \beta_0^{(1)} + x_{ij1}\beta_1^{(1)} + b_{i1}$, $\mu_{ij} = \beta_0^{(2)} + b_{i2} + x_{ij2}(\beta_1^{(2)} + b_{i3})$ with $\sigma_1 = 0.5$, $\text{logit}(q_{ij}) = \beta_0^{(3)} + x_{ij3}\beta_1^{(3)} + b_{i4}$, $\gamma_{ij} = \beta_0^{(4)} + b_{i5} + x_{ij4}(\beta_1^{(4)} + b_{i6})$ with $\sigma_2 = 0.5$ and $\text{logit}(\lambda_{ij}) = \beta_0^{(5)} + x_{ij5}\beta_1^{(5)} + \eta_1 b_{i1} + \dots + \eta_6 b_{i6}$ where the true values of the fixed effect coefficients are identical to setting I and, $\boldsymbol{\eta} = (\eta_1, \dots, \eta_6) \stackrel{i.i.d}{\sim} \text{Unif}(-0.3, 0.3)$. The random effects $\mathbf{b}_i = (b_{i1}, \dots, b_{i6})$ are sampled from $N_6(\mathbf{0}, \boldsymbol{\Sigma})$, independently for each i , where

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.2 & -0.3 & 0.4 & 0.5 & 0.3 \\ 0.2 & 3 & -0.2 & 0.9 & 0.7 & 0.1 \\ -0.3 & -0.2 & 1 & 0.2 & 0.3 & 0.2 \\ 0.4 & 0.9 & 0.2 & 0.8 & 0.5 & 0.4 \\ 0.5 & 0.7 & 0.3 & 0.5 & 4 & 0.3 \\ 0.3 & 0.1 & 0.2 & 0.4 & 0.3 & 1 \end{pmatrix}$$

Table 1: Simulation setting I ($n = 500, m = 30, p = 10, K = 5$) - average False Negatives(FN), average False Positives (FP) for fixed (composite or not) and random effects and, % datasets with non-hierarchical selection.

Model	Fixed Effects		Random Effects		% Non Hier. Selec.
	FN	FP	FN	FP	
AI	0	2.76	0.16	0.36	0
Activity Time	0	1.44	0	0.04	0
EI	0	4.48	0	1.12	0
Engage. Time	0	1.52	0	0.04	0
Dropout	0	1.52	-	-	-

and the convexity constraints on the fixed effect coefficients continue to resemble that of setting I. Finally, we continue to let $\mathbf{X}_{i,5}$ be an m dimensional ordered sample from $\text{Unif}(-1, 1)$ and sample the remaining $p - 1$ predictors from a multivariate Gaussian distribution with mean $\mathbf{0}$ and covariance matrix $\text{Cov}(x_{ijr}, x_{ijs}) = 0.5^{|r-s|}$ independently for each $i = 1, \dots, n, j = 1, \dots, m$.

Recall that CEZIJ uses the split-and-conquer approach of [Chen and Xie \(2014\)](#) to split the full set of n players into K non-overlapping groups and conducts variable selection separately in each group by solving K parallel maximization problems represented by equation (9). The selected fixed and random effects are then determined using a majority voting scheme across all the K groups as described in Section 5 of the main paper. For settings I and II, we fix $(K, \omega_1 \omega_2)$ at $(5, 3, 3)$ and $(10, 6, 6)$, respectively, so that n/K is 100 in setting I and 200 in setting II. For each setting, we generate 50 datasets and assess the model selection performance in terms of the average False Negatives (in-model predictors falsely identified as being out of model) and average False Positives (out of model predictors falsely identified as being in-model) for the fixed effects (composite or not) and the random effects. To evaluate the hierarchical selection property of our framework, we also report the percentage of datasets where our method conducted non-hierarchical selection and chose predictors with random effects only.

Table 2: Simulation setting II ($n = 2000, m = 30, p = 20, K = 10$) - average False Negatives(FN), average False Positives (FP) for fixed (composite or not) and random effects and, % datasets with non-hierarchical selection.

Model	Fixed Effects		Random Effects		% Non Hier. Selec.
	FN	FP	FN	FP	
AI	0	4.13	0.07	1	0
Activity Time	0	0.47	0	0	0
EI	0	5.80	0	0.8	0
Engage. Time	0	1.07	0	0	0
Dropout	0	1.13	-	-	-

Tables 1 and 2 report the results of these simulation experiments. We see that across both simulation settings, CEZIJ selects the correct in-model predictors for the five models. The relatively higher fixed effects False positives for the AI and EI models possibly indicate some over-fitting due to the prevalence of large number of zeros in these models. However, CEZIJ selects the fixed and random effects in a hierarchical fashion such that no random effect predictor appears in any of the four models without their fixed effect counterparts. This is not surprising given the way CEZIJ updates the adaptive weights $(c_{sr}^{(t)}, d_{sr}^{(t)})$ are after each iteration. Figure 1 presents a comparison of the computing time for the two simulation settings considered here. In particular, it demonstrates that in both these settings CEZIJ, through its split-and-conquer approach for variable selection, offers a potential gain in computational efficiency against the conventional and memory intensive approach of running the selection algorithm on the undivided data. The efficiency gain reported in these figures, however, rely on the specific system configuration which in our case was Windows 7, 64 bit, 32GB RAM on an Intel i7-5820K CPU with 12 cores.

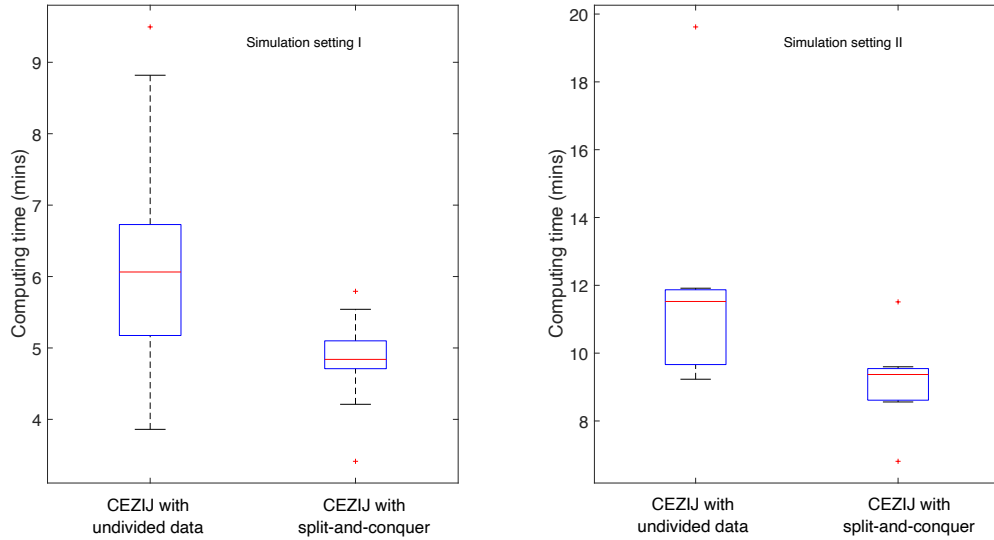


Figure 1: Computing time comparison for a fixed regularization parameter λ . Left: Simulation setting I with $n = 500, m = 30, p = 10, K = 5$. Right: Simulation setting II with $n = 2000, m = 30, p = 20, K = 10$.

D Data Description

In this section we describe the data that holds player level gaming information for a free-to-play Robot versus Robot Fighting game based on the movie Real Steel for Windows, iOS and Android devices. The primary game-play revolves around fighting and upgrading the robots while the secondary goals are to own as many robots as possible and collect rewards. A key feature of the game is a Lucky Draw which is a card game where players bet on their earnings to earn exciting in-app consumables, virtual currencies for robot upgrades or even robots! There are 38,860 players with first activity date 24-Oct-2014 and the analyses presented in section 6 uses the cohort of 33,860 players for estimation and the remaining 5,000 players for prediction. In table 3, we list the raw covariates along with their description that were available in the data and table 4 presents a descriptive summary of the raw covariates.

Promotion strategies - As discussed in Section 2, we also have side information about the different retention and promotion strategies that were used across the 60 days. These

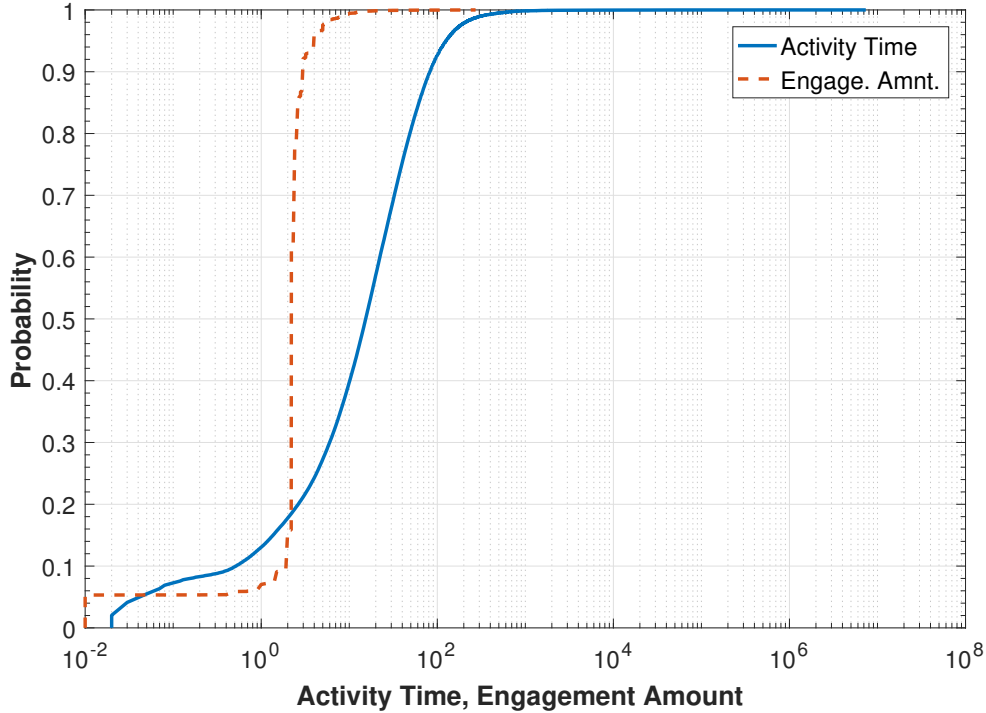


Figure 2: Empirical CDF of Activity Time and Engagement Amount.

strategies were carefully designed by the game marketers to induce player activity, boost engagement and in-app purchases at different points in time. Table 5 and figure 3 provide a summary of the 6 different promotion strategies that were used during the 60 days that the players were observed. In what follows, we provide a short description of the 6 promotion strategies.

- *Promotion strategy I* - awards extra energy points or rewards during fights when the player wins a combat.
- *Promotion strategy II* - constitutes the sale of ‘boss’ robots that possess special combat moves not available in other robots and can only be acquired by defeating the boss robot itself.
- *Promotion strategy III*- provides discounts on the purchase of powerful robots that are usually available in higher levels of the game.
- *Promotion strategy IV*- offered discounts on in-app purchases during the Black-Friday

Table 3: List of covariates and the five responses. The gaming characteristics are marked with an (*).

SI No	Covariates	Description
1	avg_session_length*	Average Session Length in Minutes
2	p_fights*	Total No. Of Principal Fights Played
3	a1_fights*	Auxiliary 1 Fights Played
4	a2_fight*	Auxiliary 2 Fights Played
5	level*	Last Principal Level Fight Played
6	robot_played*	Total No. Of Robots Played with
7	gacha_sink*	Amount of In-Game Currency Spent for Gacha
8	gacha_premium_sink*	Amount of Premium In-Game Currency Spent on Gacha
9	pfight_source*	Amount of In-Game Currency Earned by Playing Principal Fights
10	a1fight_source*	Amount of In-Game Currency Earned by Playing Auxiliary 1 Fights
11	a2fight_source*	Amount of In-Game Currency Earned by Playing Auxiliary 2 Fights
12	gacha_source*	Amount of In-Game Currency Earned by Playing Lucky Draw
13	gacha_premium_source*	Amount of Premium In-Game Currency Earned by Playing the Lucky Draw
14	robot_purchase_count*	No. Of Robot Purchased per Day
15	upgrade_count*	No. Of Robot Upgrades Done per Day
16	lucky_draw_ig*	No. Of Lucky Draw played per Day Inside Game
17	timesince*	Time Since Last Login in Days
18	lucky_draw_og*	No. Of Lucky Draw played per Day Outside Game
19	fancy_sink*	Amount of In-Game Currency Spent on Buying Accessories
20	upgrade_sink*	Amount of In-Game Currency Spent for Robot Upgrade
21	robot_buy_sink*	Amount of In-Game Currency Spent for Robot Purchase
22	gain_gachaprem*	% gain over gacha_premium_sink
23	gain_gachagrind*	% gain over gacha_sink
24	weekend	Weekend Indicator (0 - No, 1 - Yes)

SI No	Response	Description
1	AI	Whether active in a day (0 - No, 1 - Yes)
2	activity time	Total Time Played in a day in Minutes
3	EI	Whether positive engagement from the player in a day (0 - No, 1 - Yes)
4	engagement amount	Total positive engagement amount from the player in a day in dollars
5	dropout	Whether dropped out on that day (0 - No, 1 - Yes)

and Thanksgiving holiday week.

- *Promotion strategy V*- designed to promote different robots and their combat skills through emails and notifications
- *Promotion strategy VI*- provides discounts on the purchase of all robots.

Table 4: Summary statistics of the covariates reporting % of 0, mean, the 25th, 50th, 75th, 95th percentiles and the standard deviation of all active players ($\alpha_{ij} = 1$) across all $m = 60$ days. For timesince, however, the statistics are reported for all players and not just active.

Covariates	% of 0	Mean	25 th	50 th	75 th	95 th	Std.
avg_session_length	0.01	32.55	7.32	2.66	13.63	30.97	5925.89
p_fights	43.69	2.83	1.00	0.00	4.00	12.00	5.07
a1_fights	57.11	1.59	0.00	0.00	1.00	8.00	3.57
a2_fight	84.83	0.64	0.00	0.00	0.00	4.00	2.35
level	43.69	3.70	1.00	0.00	5.00	15.00	5.06
robot_played	30.68	1.15	1.00	0.00	2.00	3.00	1.16
gacha_sink	73.06	6.19	0.00	0.00	1.00	31.50	27.65
gacha_premium_sink	96.80	0.30	0.00	0.00	0.00	0.00	3.76
pfight_source	43.69	27.18	0.98	0.00	10.24	116.82	108.44
a1fight_source	57.12	3.07	0.00	0.00	1.56	15.82	10.98
a2fight_source	84.83	1.47	0.00	0.00	0.00	5.12	12.27
gacha_source	71.81	1.67	0.00	0.00	0.80	8.50	6.56
gacha_premium_source	88.93	0.85	0.00	0.00	0.00	5.00	3.81
robot_purchase_count	91.62	0.10	0.00	0.00	0.00	1.00	0.38
upgrade_count	55.26	3.30	0.00	0.00	3.00	17.00	6.74
lucky_draw_wg	55.62	1.18	0.00	0.00	1.00	4.00	3.56
timesince	7.91	22.93	21.00	7.00	37.00	54.00	17.57
lucky_draw_log	77.97	1.92	0.00	0.00	0.00	10.00	8.56
fancy_sink	87.57	0.69	0.00	0.00	0.00	1.60	10.51
upgrade_sink	55.50	18.23	0.00	0.00	10.29	85.20	74.21
robot_buy_sink	91.64	8.68	0.00	0.00	0.00	35.00	56.10
gain_gachaprem	98.36	0.04	0.00	0.00	0.00	0.00	0.45
gain_gachagrind	77.98	0.13	0.00	0.00	0.00	0.47	1.30
weekend	63.33	0.37	0.00	0.00	1.00	1.00	0.48

Table 5: Summary of the promotion strategies

Strategy	Description	No. of days	%
-	No strategy	20	33.33
I	More energy or rewards	8	13.33
II	Sale of boss robots	4	6.67
III	Discounts on powerful robots	8	13.33
IV	Holiday sale	7	11.67
V	Promotion via emailing and messaging	5	8.33
VI	Sale on all robots	8	13.33

In table 1 of the main paper, we provide the list of convex constraints imposed on the fixed effects coefficients while solving the maximization problem in equation (9).

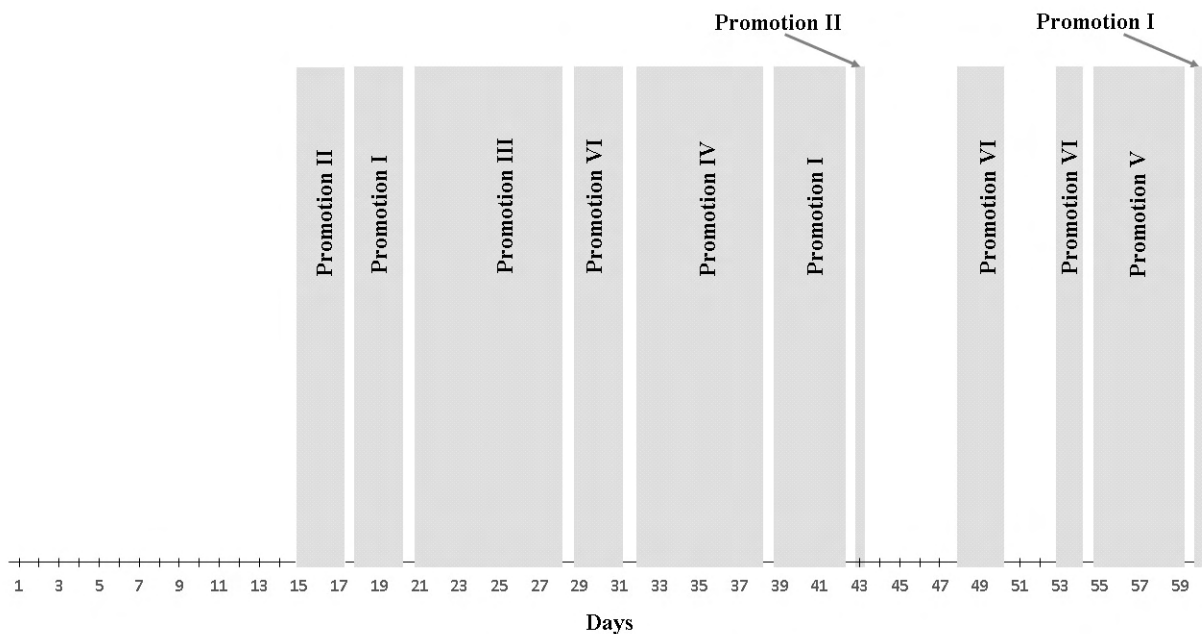


Figure 3: Distribution of the six promotion strategies over 60 days

E Variable selection by split and conquer : Voting results

Table 6 provides the voting results of the variable selection by split and conquer approach of section 5.

Table 6: The number of times each candidate predictor is selected as fixed effect and random effect across the $K = 20$ splits for the five sub-models. For each sub-model, the predictors with atleast 12 occurrences across 20 splits were selected.

Predictors	AI		Act. Time		EI		Engage. Amnt		Dropout
	Fixed Eff.	Random Eff.	Fixed Eff.	Random Eff.	Fixed Eff.	Random Eff.	Fixed Eff.	Random Eff.	Fixed Eff.
Intercept	20	20	20	20	14	14	20	20	17
avg_session_length	5	5	18	18	14	14	5	3	3
p_fights	20	20	18	18	14	14	11	11	3
a1_fights	20	20	18	18	14	14	9	9	3
a2_fights	20	20	18	18	14	14	11	11	3
level	15	14	18	18	14	14	3	0	5
robot_played	5	5	11	11	0	0	8	8	3
gacha_sink	3	0	18	18	14	14	11	0	14
gacha_premium_sink	0	0	0	0	0	0	5	5	2
pfight_source	3	0	18	18	0	0	9	8	3
a1fight_source	12	11	18	18	14	14	8	8	5
a2fight_source	17	17	18	18	14	14	12	12	11
gacha_source	20	20	18	18	0	0	11	11	3
gacha_premium_source	11	11	0	0	14	14	5	5	2
robot_purchase_count	17	17	0	0	0	0	3	0	0
upgrade_count	20	20	12	12	14	14	5	5	2
lucky_draw_wg	2	0	2	2	14	14	9	9	3
timesince	20	20	20	20	14	14	2	0	20
lucky_draw_og	18	18	0	0	14	14	6	6	3
fancy_sink	2	2	0	0	14	14	5	5	0
upgrade_sink	17	15	0	0	14	14	6	6	3
robot_buy_sink	5	5	0	0	14	6	3	2	2
gain_gachaprem	3	3	0	0	0	0	2	2	0
gain_gachagrind	18	18	18	18	0	0	11	5	3
weekend	20	12	18	18	0	0	2	0	0
promotion I	0		0		0		12		17
promotion II	20		18		14		3		17
promotion III	14		0		14		14		17
promotion IV	0		0		14		14		15
promotion V	0		0		14		0		14
promotion VI	20		18		14		14		17
# selected	18	14	17	15	22	16	6	2	9

References

- Bien, J. and Tibshirani, R. J. (2011). Sparse estimation of a covariance matrix. *Biometrika*, 98(4):807–820.
- Chen, X. and Xie, M.-g. (2014). A split-and-conquer approach for analysis of extraordinarily large data. *Statistica Sinica*, pages 1655–1684.
- Friedman, J., Hastie, T., Höfling, H., Tibshirani, R., et al. (2007). Pathwise coordinate optimization. *The Annals of Applied Statistics*, 1(2):302–332.
- Grant, M., Boyd, S., and Ye, Y. (2008). Cvx: Matlab software for disciplined convex programming.
- Rizopoulos, D. (2011). Dynamic predictions and prospective accuracy in joint models for longitudinal and time-to-event data. *Biometrics*, 67(3):819–829.
- Rizopoulos, D. (2012). *Joint models for longitudinal and time-to-event data: With applications in R*. CRC Press.
- Wang, H. (2014). Coordinate descent algorithm for covariance graphical lasso. *Statistics and Computing*, 24(4):521–529.

MATLAB implementation of the Constrained Zero Inflated Joint Modeling (CEZIJ) framework of [Banerjee et al. \(2019\)](#).

March 18, 2019

1 Introduction

CEZIJ ([Banerjee et al., 2019](#)) is a novel framework for parameter estimation in joint models with multiple longitudinal outcomes along with a time-to-event analysis. Longitudinal data from modern datasets usually exhibit a large set of potential predictors and choosing the relevant set of predictors is highly desirable for various purposes including improved predictability. To achieve this goal, CEZIJ conducts simultaneous selection of fixed and random effects in high-dimensional penalized generalized linear mixed models and maintains the hierarchical congruity of the fixed and random effects, thus producing models with interpretable composite effects. It not only accommodates extreme zero-inflation in the responses in a joint model setting but also incorporates domain-specific, convex structural constraints on the model parameters. For analyzing such large-scale datasets, variable selection and estimation is conducted via a distributed computing based split-and-conquer approach ([Chen and Xie, 2014](#)) that massively increases scalability.

2 Installation requirements

The GitHub repository holds the MATLAB toolbox `cezij.mltbx` that provides an implementation of the CEZIJ procedure developed in [Banerjee et al. \(2019\)](#). To install this

toolbox, simply download `cezij.mltbx` in your computer and double click to install it. For a successful installation, please make sure that the following system requirements are met.

- Access to 32GB RAM and at least 8 CPU cores for parallel computing
- MATLAB 2016b or higher with the following toolboxes (and their dependencies):
 - Statistics toolbox
 - Optimization toolbox
 - Parallel Computing toolbox
 - Data Acquisition toolbox
- [CVX for MATLAB \(version 2.1 or higher\)](#)

3 A numerical example

In this section, we will use a numerical example to illustrate the use of the `cezij` toolbox. Our goal is to assess the model selection performance of CEZIJ under the longitudinal and Dropout models discussed in Section 3.1 of [Banerjee et al. \(2019\)](#). To do that, we use the simulation example of setting I discussed in section C.2 of the supplementary materials and indicate which files in the toolbox should be edited to test a different dataset. Figure 1 presents three MATLAB scripts that should be edited to prepare the `cezij` toolbox for analyses. In what follows, we discuss these scripts.

3.1 Generating simulated data - `simulate_data.m`

To run the `cezij` variable selection algorithm on a dataset of your own choice, please edit the m file `simulate_data.m` and ensure that your data is in the same format as the output of this m file. In the default setting, `simulate_data.m` generates a simulated dataset that we describe below.

Consider a sample of $n = 500$ players and for each player i , let $\mathbb{X}_i = (\mathbf{X}_{i.1}, \dots, \mathbf{X}_{i.p})$ denote the $m \times p$ matrix of candidate predictors where $\mathbf{X}_{i.k} = (x_{i1k}, \dots, x_{imk})^T$, $m = 30$

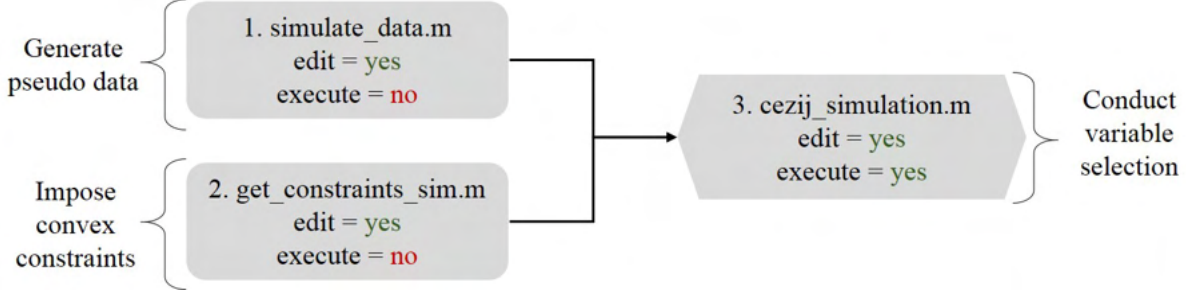


Figure 1: There are three MATLAB scripts, numbered 1-3, that must be edited to prepare the `cezij` toolbox for analyses. To reproduce table 4 in the supplementary file or table 1, script 3 must be executed without making any changes to the default parameters in scripts 1-3.

denotes the number of time points for which we observe each player and $p = 10$ is number of candidate predictors. We take $\mathcal{I}_f = \{1, \dots, 8\}$ as the indices of the fixed effects and, $\mathcal{I}_c = \{9, 10\}$ as the indices of the composite effects so that $p_f = |\mathcal{I}_f| = 8$, $p_c = |\mathcal{I}_c| = 2$. Thus, the first 8 columns of \mathbb{X}_i represent fixed effects while the last 2 represent composite effects. The five responses $[\alpha_i, \mathbf{A}_i, \epsilon_i, \mathbb{E}_i, \mathbb{D}_i]$ corresponding to AI, positive Activity, EI, Positive Engagement and Dropout are generated from the following models: $\text{logit}(\pi_{ij}) = \beta_0^{(1)} + x_{ij1}\beta_1^{(1)} + b_{i1}$, $\mu_{ij} = \beta_0^{(2)} + x_{ij2}\beta_1^{(2)} + b_{i2}$ with $\sigma_1 = 0.5$, $\text{logit}(q_{ij}) = \beta_0^{(3)} + x_{ij3}\beta_1^{(3)} + b_{i3}$, $\gamma_{ij} = \beta_0^{(4)} + x_{ij4}\beta_1^{(4)} + b_{i4}$ with $\sigma_2 = 0.5$ and $\text{logit}(\lambda_{ij}) = \beta_0^{(5)} + x_{ij5}\beta_1^{(5)} + \eta_1 b_{i1} + \dots + \eta_4 b_{i4}$ where the true values of the fixed effect coefficients are: $\boldsymbol{\beta}^{(1)} = (1, -1.5)$, $\boldsymbol{\beta}^{(2)} = (3.5, -2)$, $\boldsymbol{\beta}^{(3)} = (1, -1)$, $\boldsymbol{\beta}^{(4)} = (3, -3)$, $\boldsymbol{\beta}^{(5)} = (-1, 2)$ and, $\boldsymbol{\eta} = (\eta_1, \dots, \eta_4) = (-0.1, 0.2, 0.1, -0.2)$. Thus this setting presents a scenario wherein there are no composite effects in the true model. The random effects $\mathbf{b}_i = (b_{i1}, \dots, b_{i4})$ are sampled from $N_4(\mathbf{0}, \boldsymbol{\Sigma})$, independently for each i , where

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.2 & 0.4 & 0.5 \\ 0.2 & 3 & 0.9 & 0.7 \\ 0.4 & 0.9 & 0.8 & 0.5 \\ 0.5 & 0.7 & 0.5 & 4 \end{pmatrix}$$

Finally, to complete the specification, we sample $(x_{ij1}, \dots, x_{ij4})$ from $N_4(\mathbf{0}, 4\mathbf{I}_4)$, independently for each $i = 1, \dots, n$, $j = 1, \dots, m$. To ensure that the generated sample contains

players that have not churned for at least the first 7 to 10 days, we let $\mathbf{X}_{i,5}$ to be an m dimensional ordered sample from $\text{Unif}(-1, 1)$ so that $\mathbf{X}_{i,5} = (x_{i15} \leq \dots \leq x_{im5})$ and, generate the remaining predictors independently from $\text{Unif}(-1, 1)$. Although the MATLAB file `simulate_data.m` stores the above simulation setting, it can easily be modified to test different settings.

3.2 Imposing convex constraints - `get_constraints_sim.m`

The CEZIJ framework can incorporate convexity constraints on the fixed effect coefficients and this MATLAB file stores the following default constrains:

$$\beta_0^{(1)} > 0, \beta_1^{(1)} < 0; \beta_1^{(3)} < 0; \beta_0^{(4)} > 0, \beta_1^{(4)} < 0.$$

Please modify this file to enforce constraints specific to your application or leave this file unchanged to reproduce table 4 of the supplementary materials or table 1 in `cezij_help.pdf`.

3.3 Running the joint model - `cezij_simulation.m`

Recall that CEZIJ uses the split-and-conquer approach of [Chen and Xie \(2014\)](#) to split the full set of n players into K non-overlapping groups and conducts variable selection separately in each group by solving K parallel maximization problems represented by equation (9) of [Banerjee et al. \(2019\)](#). The selected fixed and random effects are then determined using a majority voting scheme across all the K groups as described in Section 5 of the above paper.

In the MATLAB file `cezij_simulation.m`, lines 15-27 can be used to specify a number of user defined parameters. For this example, we fix $K = 5$ so that n/K is 100 while $q = 3$ indicates the number of random effects including a random intercept. We generate `nsets = 25` datasets and run `cezij_simulation.m` to assess the model selection performance in terms of the average False Negatives (in-model predictors falsely identified as being out of model) and average False Positives (out of model predictors falsely identified as being in-model) for the fixed effects (composite or not) and the random effects. To evaluate

Table 1: ($n = 500, m = 30, p = 10, K = 5$) - average False Negatives(FN), average False Positives (FP) for fixed (composite or not) and random effects and, % datasets with non-hierarchical selection.

Model	Fixed Effects		Random Effects		% Non Hier. Selec.
	FN	FP	FN	FP	
AI	0	2.76	0.16	0.36	0
Activity Time	0	1.44	0	0.04	0
EI	0	4.48	0	1.12	0
Engage. Time	0	1.52	0	0.04	0
Dropout	0	1.52	-	-	-

the hierarchical selection property of our framework, the code also reports the percentage of datasets where `cezij` conducted non-hierarchical selection and chose predictors with random effects only.

Running `cezij_simulation.m` with the default parameters generates Table 1 that reports the results of the simulation exercise under setting I that is discussed in section C.2 of the supplementary materials. We see that CEZIJ selects the correct in-model predictors for the five models. The relatively higher fixed effects False positives for the AI and EI models possibly indicate some over-fitting due to the prevalence of large number of zeros in these models. However, CEZIJ selects the fixed and random effects in a hierarchical fashion such that no random effect predictor appears in any of the four models without their fixed effect counterparts. This is not surprising given the way CEZIJ updates the adaptive weights ($c_{sr}^{(t)}, d_{sr}^{(t)}$) are after each iteration.

4 Simulation flow

In figure 2, we present a simulation flow diagram that depicts the main scripts that are called when `cezij_simulation.m` is executed. The scripts highlighted in blue are editable and can be used to run the analyses on a different data set as described in section 3. The

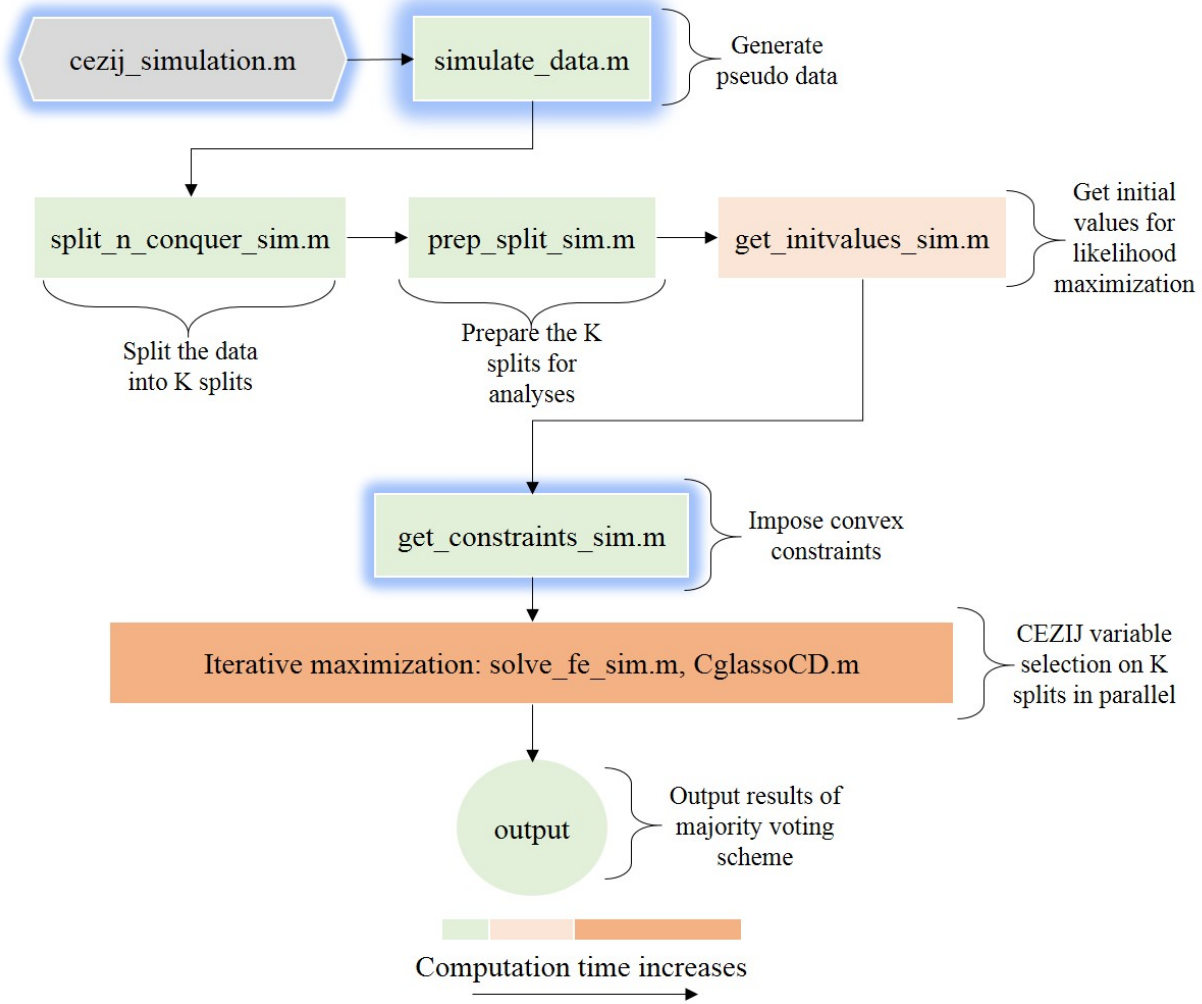


Figure 2: Simulation flow diagram that depicts the main scripts that are called when `cezij_simulation.m` is executed. The scripts highlighted in blue are editable while the color coding of the scripts indicate their relative contribution to total computation time.

color coding of the scripts indicate their relative contribution to total computation time. For instance, the iterative maximization step that is executed in parallel across the K splits is the most computationally intensive step of the `cezij` algorithm and, as discussed in section C.2 of the supplementary material, relies on the specific system configuration and the number of computation cores available. In the default setting that is used to reproduce table 1, this step takes approximately 5 minutes to execute (see figure 7 of the supplementary materials). Depending on the number of splits K , availability of additional

computational cores may further reduce the overall computation time.

References

Banerjee, T., Mukherjee, G., Dutta, S., and Ghosh, P. (2019). A large-scale constrained joint modeling approach for predicting user activity, engagement and churn with application to freemium mobile games. *under review*.

Chen, X. and Xie, M.-g. (2014). A split-and-conquer approach for analysis of extraordinarily large data. *Statistica Sinica*, pages 1655–1684.