



This is a self-archived – parallel-published version of an original article. This version may differ from the original in pagination and typographic details. When using please cite the original.

AUTHOR	R. Numminen, M. J. Viljanen and T. Pahikkala
TITLE	Bayesian inference for predicting the monetization percentage in free-to-play games
YEAR	2020.
DOI	http://www.doi.org/10.1109/TG.2020.3014660
CITATION	R. Numminen, M. J. Viljanen and T. Pahikkala, "Bayesian inference for predicting the monetization percentage in free-to-play games," in <i>IEEE Transactions on Games</i> , doi: 10.1109/TG.2020.3014660.

Bayesian Inference for Predicting the Monetization Percentage in Free-to-Play Games

Riikka Numminen

Department of Future Technologies
University of Turku
Turku, Finland
rimanu@utu.fi

Markus Viljanen

Department of Future Technologies
University of Turku
Turku, Finland
majuvi@utu.fi

Tapio Pahikkala

Department of Future Technologies
University of Turku
Turku, Finland
aatapa@utu.fi

Abstract—Free-to-play has become one of the most popular monetization models, and as a consequence game developers need to get the players to purchase in the game instead of getting players to buy the game. Game analytics and player monetization prediction are important parts in estimating the profitability of a free-to-play game. In this paper, we concentrate on predicting the fraction of monetizing players among all players. Our method is based on a survival analysis mixture cure model, and can be applied to unlabeled data collected from any free-to-play game. We formulate a statistical model and use the Expectation Maximization algorithm to solve the latent monetization percentage and the monetization rate. The original method is modified by using Bayesian inference, and the results of the versions are compared. The method can be applied as a preliminary profitability study in situations where there is no extensive historical game data available, such as game and business development scenarios that need to utilize real time analytics.

Index Terms—Bayesian Inference, Free-to-play, Monetization, Survival Analysis

I. INTRODUCTION

Total revenue from video games reached approximately 110 billion dollars in 2018 [1]. As the gaming industry has grown, the revenue models have also evolved. Most revenue in the past was from game purchases and subscriptions whereas today free-to-play games account for the majority of all game titles and revenues [1]. Money is made through advertisements, premium upgrades and in-app purchases. However, only around 5% of players in a successful free-to-play game can be expected to monetize [2]. These developments have made it important to understand exactly what percentage of players monetize and why they do so.

The data sets collected from games can be divided into a scale between two extremes:

- 1) Extensive historical game data spanning maybe years.
- 2) A completely new game data set with short follow-ups.

The first setting allows the usage of supervised machine learning models, since it is known which users made the purchases and can thus be generalized to the same game. However, in the second setting a new game with only few observed purchases may be dealt with and the data set is semisupervised because the correct answers are not known for most of the players.

The first setting is well suited for academic research. The second setting occurs when game developers want to use real-time analytics to understand their current game. They want to know as soon as possible how profitable a game is expected to be, as they do not want to expend finances on advertising if the game is expected to be unprofitable. Game literature has demonstrated that it is possible to train machine learning models with good predictive performance in the first setting, many of which were featured in a recent competition [3]. Research has been more limited in the second setting. The method presented in this article investigates the second setting, for a new game with only a short period of data collection.

The goal of our method is to predict the proportion of all players that will monetize over time. The method belongs to the field of game analytics, which is concerned with understanding player behavior. Game developers use real-time analytics when they are developing or planning to launch a game. A data set is generated by tracking players for a certain duration, which we call the follow-up. Some players monetize during the follow-up, some will monetize later but have not yet done so, and some players will not monetize. We describe the percentage of monetizing players as an unknown latent variable and solve it using the Expectation Maximization algorithm. This article is an extended version of the conference presentation [4]: prior information of the monetization percentage is taken into account and the method is tested with a new data set. The results, however, demonstrate that using unlabelled data is a more challenging task and requires further research. Even though our focus is on the monetization percentage and rate, in principle the latent variable formulation could be used together with many supervised machine learning models to train them on data that is semisupervised because of a limited follow-up. Many real world data sets are somewhere between these two extremes and the approaches could complement each other.

There is increasing awareness about the ethical issues of data collection and model based prediction [5]. Many of these issues are also present in games, especially in the retention and monetization aspects of player modeling [6]. The free-to-play monetization model is also under increased scrutiny, as there can be problems about disadvantaging certain players and encouraging problematic financial behavior [7]. Game

developers that plan to use micro-transactions in their game should also consider these aspects. In our research, the most relevant ethical issues are obtaining informed consent and anonymizing the collected data. We have taken into account the ethical aspects as far as they are related to our research, although we present the research from a mathematical point of view.

II. RELATED WORK

Regression and machine learning have been used in studies to predict player purchases in various problem formulations for labeled historical data. Random Forest, linear SVM, and Decision Tree were used first in [8] to classify whether a player would buy an in-game item after a match. Both general in-game items purchases and hard currency purchases were predicted. Similarly, given purchasing and non-purchasing players with two weeks of history before the purchase, Decision Tree, Logistic Regression, and SVM were used in [9] to predict which group the player belonged to, with a focus on the game agnostic features. In [10], both classification and regression were used to predict whether the user would make a first purchase and how many purchases would occur. Decision Trees, Random Forests, and Support Vector Machines were used with data set balancing methods. Finally, in [11] two linear regressions were used to model the number of purchases and the number of coins purchased at a given level. The focus was on understanding the impact of gating mechanisms on retention and monetization.

Survival analysis has been used in gaming for various other tasks, see [12] for a review. Noncontractual probability models used in marketing [13] are closest to our approach. These models predict player purchase counts over time, given a data set in the second setting. However, one of the most popular models (BG/NBD) was tested in free-to-play games and the authors found that the model struggled with covering real data [14]. It has been suggested that further research should be conducted to redesign or adjust existing models, in order to examine better assumptions for free-to-play games. In a recent approach [15], a model-free method was developed to measure the mean customer lifetime value (LTV) in the second setting, but this approach did not predict into the future. Studies have also investigated how the first purchase predicts overall LTV [16], which is very useful when used together with our model.

III. METHOD

A. Monetization as a survival analysis model

Monetization can be expressed as a survival analysis model. Survival analysis is used for analysing time-to-event data with limited follow-ups and the data gathered from a new game with short follow-ups is exactly this type of data. In survival analysis there is a sample for which it is observed which individuals have a specific event during the follow-up [17]. In our case the event is making the first purchase. The follow-up means the time an individual is in the study. We will use the calendar time between a player starting to play the game and the player making the first purchase. However, not every player

makes the first purchase during the data collection interval. Those that do not purchase anything during the follow-up, are censored, which means that their event times are not known because the follow-up ended before the events occurred.

It can be known which players made a purchase during the interval and which players were followed until the censoring time. Mathematically said, there are two variables that are observed in time-to-event data with limited follow-up: One of them is the time $T = \min(T^*, C)$ which is either the purchase time T^* or the censoring time C whichever is smaller. The other variable is the censoring indicator $\delta = \mathbb{I}(C \leq T^*)$ which is a binary variable demonstrating whether a player made a purchase before censoring or was censored by the follow-up. Realizations of these random variables are denoted with $t_i = \min(t_i^*, c_i)$ and $\delta_i = \mathbb{I}(c_i \leq t_i^*)$ for the i^{th} player.

In this paper, we assume that the purchase time $T^* \sim \text{Exp}(\lambda)$ and the censoring time C is implied by the data collection time. The assumption of purchase time following the exponential distribution is a special case of the playtime principle introduced in [18] and has been used to model player survival [19].

The distribution of event time T^* is defined by a survival function S . In general, the survival function $S(t) = 1 - F(t)$, where F is the cumulative distribution function [20]. Thus S describes the probability that the event time is greater than the observation time t and for the exponential model it is

$$S(t) = P(T^* > t) = e^{-\lambda t}. \quad (1)$$

Now that the probability that a player makes the first purchase later than the time t is known, the second thing to consider is the risk for a player to purchase at time t . The instantaneous risk that a player purchases at time t given that he/she did not do so until that time, is described by a hazard function h [21]:

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T^* < t + \Delta t | T^* > t)}{\Delta t} = \lambda. \quad (2)$$

A probability density function can be calculated from survival and hazard functions. It is

$$f(t) = h(t) S(t) \quad (3)$$

and means the density of purchases at time t .

B. Monetizing players as a mixture cure model

In standard survival analysis, each individual is assumed to eventually have the event [17]. This assumption does not apply in free-to-play games since many players seem to never buy anything. This means that there are two kinds of players, monetizing and unmonetizing, and the whole population of the players is a mixture of the two sub-populations. Hence a mixture cure model [22] is required in order to properly model this situation.

All unmonetizing players are always censored because they never purchase anything. However, also some of the monetizing players might be censored if they were not followed for long enough. Thus the division into censored and purchased players does not tell the whole truth about the number of monetizing players. A third variable, a monetizing indicator ζ , is used for describing which sub-population the player belongs

to: $\zeta = 0$ for the monetizing population and $\zeta = 1$ for the unmonetizing. Monetizing indicator $\zeta \sim \text{Bern}(1 - \pi)$, and the probabilities that a player is a monetizing or an unmonetizing player are $P(\zeta = 0) = \pi$ and $P(\zeta = 1) = 1 - \pi$, respectively. This variable is partly latent because the value of it is known only for those players that made a purchase before the censoring.

In a mixture cure model, the survival function is a weighted sum of the survival functions of the sub-populations: $S(t) = \pi S_m(t) + (1 - \pi) S_u(t)$, where the weight π is the percentage of monetizing individuals in the whole sample, $S_m(t) = e^{-\lambda t}$ and $S_u(t) \equiv 1$. Given that the purchase time follows the exponential distribution, there are now two parameters that describe the model: $\Psi = (\pi, \lambda)$. They are monetization percentage and conversion rate.

An example of this kind of data is shown in Fig. 1. In the example there are 50 players, of whom nine are monetizing. It can be seen that two monetizing players have not purchased before censoring. The players have started within three units of the calendar time and censoring occurs five time units after the first player arrived as can be seen in the upper subplot. This results in different follow-ups for the players as is demonstrated in the lower subplot. In reality the monetization status of players, i.e. the colors in the figure, are not known, and it is necessary to infer from the data set how many players are going to make a purchase.

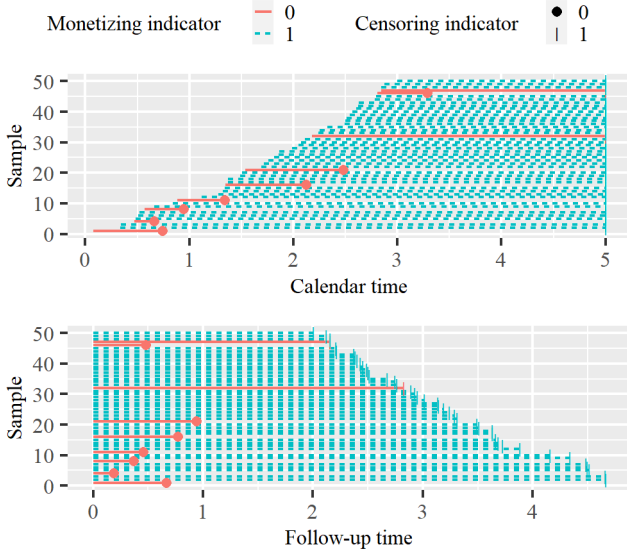


Fig. 1. Simulated data example with 50 players: 9 monetizing and 41 unmonetizing. The data were generated with parameters $(\pi, \lambda) = (0.1, 1.0)$.

C. Fitting the mixture cure model

We can infer the monetization percentage and the conversion rate by finding the model parameter vector Ψ . The likelihood function L shows how likely the probability distribution samples are, given values for the parameters. The maximum likelihood (ML) estimate $\hat{\Psi}_{\text{ML}}$ is the parameter vector that maximizes the likelihood function, i.e. parameter values that

make the given data most likely. In survival analysis the likelihood function

$$L(\Psi) = \prod_{i=1}^n f(t_i|\Psi)^{1-\delta_i} S(t_i|\Psi)^{\delta_i}, \quad (4)$$

where n is the sample size. However, the logarithm of it,

$$l(\Psi) = \sum_{i=1}^n \{(1 - \delta_i) \log f(t_i|\Psi) + \delta_i \log S(t_i|\Psi)\}, \quad (5)$$

is often used instead [23]. Since the logarithm is a strictly increasing function, the maximum likelihood estimate $\hat{\Psi}_{\text{ML}}$ is the same for both (4) and (5). In the maximum likelihood estimation the parameter values converge in probability to the true parameter values as $n \rightarrow \infty$ [24].

If the latent monetizing status ζ_i is somehow known for every player, the data is said to be complete and the solution is both simple and intuitive. The total number of players is denoted with $n = n_0 + n_1$, where n_0 stands for the number of monetizing players and n_1 is the number of unmonetizing players. In order to find the maximum likelihood estimate, the roots of the partial derivatives of (4) or (5) with respect to π and λ are found separately. Then the parameters simply are

$$\hat{\pi} = \frac{n_0}{n} \text{ and } \hat{\lambda} = \frac{n_0}{\sum_{i:\zeta_i=0} t_i}. \quad (6)$$

In other words, the monetization percentage is the fraction of players that purchase something. The conversion rate is the number of monetized players divided by their total exposure time. When considering the fact that $\mathbb{E}[T^*] = 1/\lambda$, it can be seen that the expected purchase time is the average of exposure times in the monetizing population.

However, there is the latent variable in the mixture cure model since it is not known which players are monetizing. That variable makes it impossible to find an explicit equation for the maximum likelihood estimate. Expectation Maximization (EM) algorithm [25] is an iterative algorithm that is suitable for maximum likelihood estimation in situations where there are missing data or latent variables. In a situation like this, the observed data is said to be incomplete. As shown in the appendix, this results in an iterative algorithm, which updates the current value of the parameter $\pi^{(k)}$ by

$$\pi^{(k)} = \frac{1}{n} \left[\sum_{i:\delta_i=1} \frac{\pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}}{1 - \pi^{(k-1)} + \pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}} + \sum_{i:\delta_i=0} 1 \right] \quad (7)$$

and the current value of $\lambda^{(k)}$ is calculated with

$$\lambda^{(k)} = \frac{\sum_{i:\delta_i=0} 1}{\sum_{i:\delta_i=1} \frac{\pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}}{1 - \pi^{(k-1)} + \pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}} t_i + \sum_{i:\delta_i=0} t_i}. \quad (8)$$

These estimates converge to the global maximum of the log-likelihood function when $k \rightarrow \infty$. The values of the log-likelihood function (5), and the method iterations (7) and (8) are illustrated with a black line in Fig. 2 for the data represented in Fig. 1.

There are also some problematic cases when the model does not work correctly. A zero-frequency problem is present when all players are censored and there is no information showing that some of the players monetize. In such a situation the model is not able to distinguish between the survival

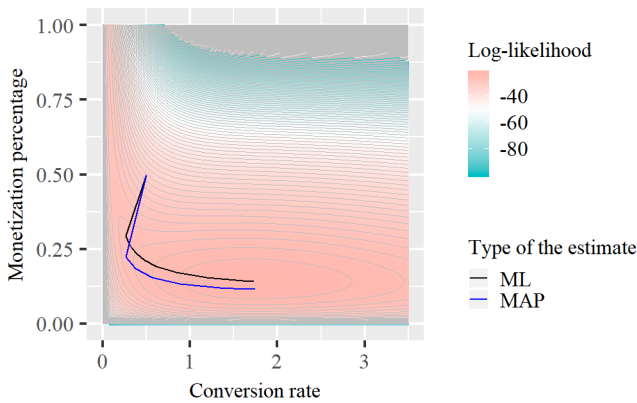


Fig. 2. Values of the log-likelihood as a function of π and λ . The paths of the EM-algorithm iterations are presented as curves from the initial guess $\Psi^{(0)} = (0.5, 0.5)$ to the maximum likelihood estimate $(\hat{\pi}_{ML}, \hat{\lambda}_{ML}) = (0.14, 1.73)$ and to the maximum a posteriori estimate $(\hat{\pi}_{MAP}, \hat{\lambda}_{MAP}) = (0.12, 1.75)$.

analysis model and the mixture cure model, and the method predicts what the survival analysis assumes, i.e. that eventually every player monetizes. This problem can be avoided by using Laplace smoothing [26] which is a method that makes all classes (unmonetizing and monetizing populations in our case) existent by adding pseudo-observations to the data. It is enough to add one unmonetizing pseudo-observation with infinite follow-up to the data to avoid the zero-frequency problem and obtain credible results in reasonable computation time.

D. Prior information

In addition to Laplace smoothing, the problematic situations can be avoided by using Bayesian analysis. In that both prior information and sample information (likelihood function) are used to obtain a posterior distribution of the considered parameters [27]. Prior information means that something is known about the values of the parameters before seeing any data. For example, it is known that even successful free-to-play games have monetization percentages in the single digit range. This means we have prior information about the monetization percentage π . A natural prior for π is the Beta (α, β) distribution because its parameters are related to the numbers of monetizing and unmonetizing players. In fact, the Laplace smoothing that was used in [4] corresponds to a Beta $(1, 2)$ distribution. Four examples of Beta distribution are shown in Fig. 3. By changing the values of the parameters, the density can be concentrated around a certain monetization percentage, e.g. with Beta $(2, 20)$ distribution it is around 5% and with Beta $(10, 82)$ distribution it is around 10%.

In Bayesian analysis the optimal values for parameters are called maximum a posteriori (MAP) estimates because they maximize the posterior density function. EM-algorithm can be used also for finding maximum a posteriori estimates. The deriving of the formulas is shown in the appendix. Now that

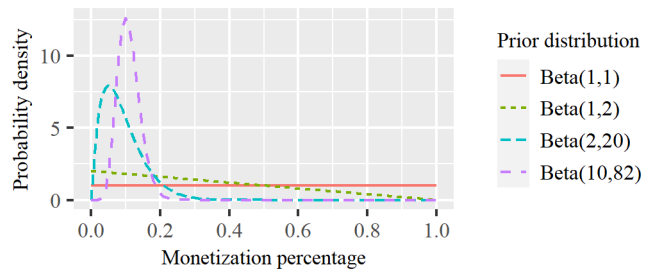


Fig. 3. Four examples of the Beta (α, β) prior distribution for the monetization percentage.

we use a prior only for the monetization percentage π , only the formula (7) needs to be modified from the maximum likelihood case. The current value of the parameter $\pi^{(k)}$ is now updated by

$$\pi^{(k)} = \left[\sum_{i:\delta_i=1} \frac{\pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}}{1 - \pi^{(k-1)} + \pi^{(k-1)} e^{-\lambda^{(k-1)} t_i}} + \sum_{i:\delta_i=0} 1 + \alpha - 1 \right] \cdot \frac{1}{n + \alpha + \beta - 2}. \quad (9)$$

As can be seen, there are only minor changes in the formula. They seem very obvious when thinking about the meaning of the parameters α and β . The parameters can be thought as the numbers of pseudo-observations that are added in the data: $\alpha - 1$ monetizing players and $\beta - 1$ unmonetizing players. It can also be seen that (9) is equal to (7) if a noninformative prior Beta $(1, 1)$ is used. The EM-algorithm iterations obtained by using (9) instead of (7) are shown in Fig. 2 as a blue curve. The MAP estimate is almost the same as the ML estimate, but it avoids the zero-frequency problem and formalizes the idea of having prior assumptions about the data.

IV. DATA SETS

Based on the formulas we derived, the model was implemented with the R programming language [28]. First we tested the variability and bias of the model on generated data. After that we used a game data set and CDNOW data set to see how well the model works with real data. Now we introduce the data sets and verify that the real data sets follow the assumptions of the model.

A. Generated data

To our knowledge, there is no publicly available person level data set with many free-to-play games. One of our primary motivations is to apply the algorithm to completely new games with limited follow-ups, therefore we simulated different sample sizes n , censoring times c , and monetization percentages π and then calculated the estimate $\hat{\pi}$ for each data set.

In the first experiment, we calculated the monetization percentage estimates for different sample sizes and follow-up times, given a true value of $\pi = 0.10$ and Beta $(2, 10)$ and Beta $(10, 82)$ priors. The effect of sample size n was tested

with values 100, 500, 1000 and 5000. The censoring time c was tested with values defined in a way that there is 25%, 50%, 75% or 100% probability for the monetizing players to purchase before censoring. For each combination, we conduct a thousand experiments $r = 1, \dots, 1000$. In each experiment, we sampled a player data set of size n , where the observed time $t_i = \min(t_i^*, c)$, and calculated the predicted value of $\hat{\pi}_r$ using the EM-algorithm.

In the second experiment, we tested the effect of the true monetization percentage on the estimate, and aimed to find appropriate prior distributions for them. The effect of the true monetization percentage was tested with values 0.01, 0.05, 0.1, 0.5 and 0.75. Censoring times were defined the same way as in the first simulation but with 10–100% probabilities to purchase before censoring. The effect of the sample size was tested with 10 different values varying from 100 to 1000 by 100. For each (n, π, c) triplet, 1000 estimates were computed and averaged.

B. Game data

The real game data were collected from a free-to-play mobile game by the company whose game it is. The game was in-development during data collection, and many developments were made to each version over the game development cycle. Periodical user acquisition tests were used to evaluate the current performance. In these tests, a group of players was obtained by using paid advertisement in Facebook targeted by country and device operating system. Each player has a random hexadecimal ID and hence the data sets are anonymous. The behaviour of the players was recorded only inside the game.

The number of players varies between different versions and only a small percentage of the players purchased an item as can be seen in Table I. The version subsets are considered as independent data sets and it is not taken into account that there may be some players that are included in multiple of them. The monetization in the game was likely improved from version 1.18 to version 1.21, but the subsequent development in the 1.3x series had no large effect on the percentage of purchasing players.

TABLE I
NUMBERS OF PLAYERS AND MONETIZED PLAYERS IN THE REAL DATA.

Version	# of players	# of monetized players	π
1.18	1604	6	0.004
1.21	309	6	0.019
1.31	1691	24	0.014
1.32	1582	21	0.013
1.33	1211	18	0.015
1.35	2364	35	0.015

We used these data to create censored data sets that replicate the actual version user test. The event time was the calendar time from the beginning of the first session to the first purchase. We defined censored data sets by varying the data collection date as $1, \dots, D$ days of calendar time from the beginning of the first session of the first player. Maximal

follow-up D denotes the actual data collection date, after which we have no data. This created censored players with different follow-up times, exactly the same way the data set would be obtained if it was updated at the end of each day after the test began.

C. CDNOW data

CDNOW data set consists of a purchase history of 23 570 individuals in an online retail shop. Every individual in this data set made the first purchase at CDNOW in 1997 during the first quarter. The data were collected until the end of June 1998. This data set was first used in [29].

We are interested in predicting how many of the customers will return to the online shop and buy again. The beginning of the follow-up is now the date of the first purchase and the event time is the date of the second purchase. Now that we only know the date instead of knowing also the time of the purchase, there are some observed times that are equal to zero, i.e. the second purchase was made the same day as the first one.

We used the same method as with the game data also with these data. We created censored data sets that replicate the situation that the number of returning customers was estimated every day since the beginning of the data collection. Again the maximal follow-up is the date when the data collection ended, which is 30.6.1998.

D. Model assumption verification

There are two assumptions in the model about the data which need to be verified for the real data. The assumptions are:

- 1) the event times are exponentially distributed and
- 2) there are individuals that do not have the event: $\pi < 1$.

Assumption 1 is verified with Q–Q-plots in which the quantiles of observed event times are compared to the quantiles of an exponential distribution. The Q–Q-plots for the game data are shown in Fig. 4. There are some exceptions with late purchase times, but most of the points are along a straight line. Most of the points are along the straight line also for CDNOW data which can be seen in Fig. 5. These suggest that for the purposes of estimating the monetization fraction, we may assume that the event times follow exponential distribution.

Akaike’s information criterion [30] is used to verify assumption 2. This method consists of calculating an AIC value for each compared model and the smaller the value, the better the model describes the data. The value is calculated with

$$\text{AIC} = 2n_p - 2 \log(\hat{L}), \quad (10)$$

where n_p is the number of parameters in a model and $\hat{L} = L(\hat{\Psi})$. A maximum likelihood estimate is calculated for the incomplete data likelihood function because the values of the monetizing indicator ζ are not known and the value of the complete data likelihood function cannot be calculated. AIC values in Table II show that the mixture cure model is better than the regular survival model at explaining the data for every

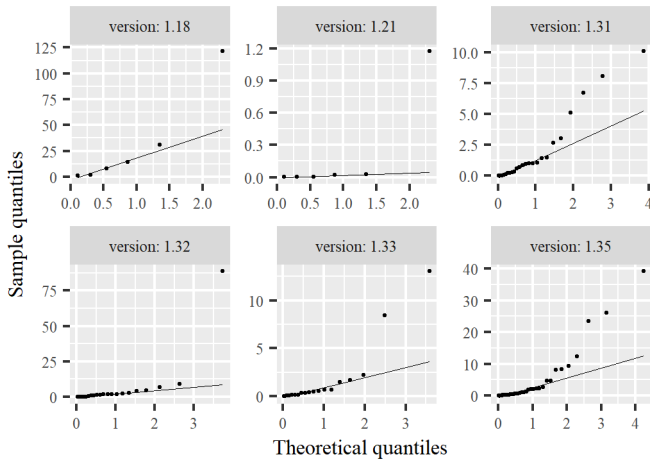


Fig. 4. Q-Q-plots of the event times of each version.

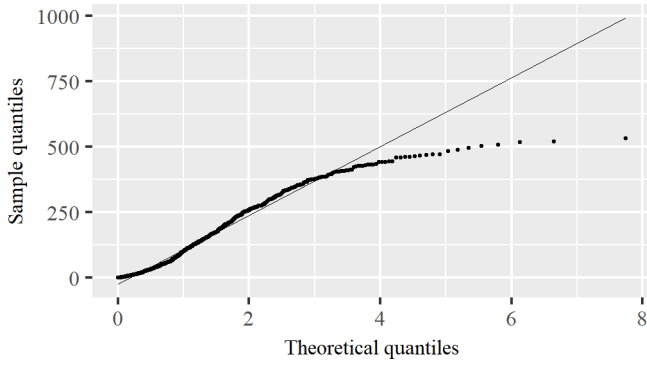


Fig. 5. Q-Q-plot of the event times of CDNOW data.

game version and for the CDNOW data. This is not surprising, given our prior knowledge about free-to-play monetization.

TABLE II
AIC VALUES FOR MIXTURE CURE AND REGULAR SURVIVAL MODELS.

Version	Model	
	$\pi < 1$	$\pi = 1$
1.18	135.6720	154.0050
1.21	114.0415	132.1888
1.31	450.3302	501.7088
1.32	401.4259	436.8337
1.33	338.0034	364.8098
1.35	659.8481	693.9284
CDNOW	16288.19	17165.52

V. RESULTS AND DISCUSSION

The simulation studies were run several times with different prior distributions. The prior distributions used with real data sets are chosen based on the results of the simulation studies.

A. Simulation

The results of the first simulation are shown in Fig. 6. The distributions of the estimates are visualised with boxplots for

every sample size and censoring time. When the probability to purchase before censoring is one, we know which players are monetizing and which are not. Thus that case is the same as if the data were complete. The estimates were calculated by using two different prior distributions with mode equal to the true value $\pi = 0.1$.

With this experiment we demonstrate the effect of the prior on the variance of the estimate. It can be seen that the stronger the used prior is, the smaller the variance of the estimate is. When the sample size is small and the follow-up is short, there is hardly any information in the data, the estimate is mostly based on the prior, and the variance of the estimate depends on the variance of the prior. When there is more information in the data (greater sample size or longer follow-up), the effect of the prior decreases, and the variability of data implies larger variance of the estimate. When the follow-up time approaches infinity, there is sufficient information in the data and the variance decreases again.

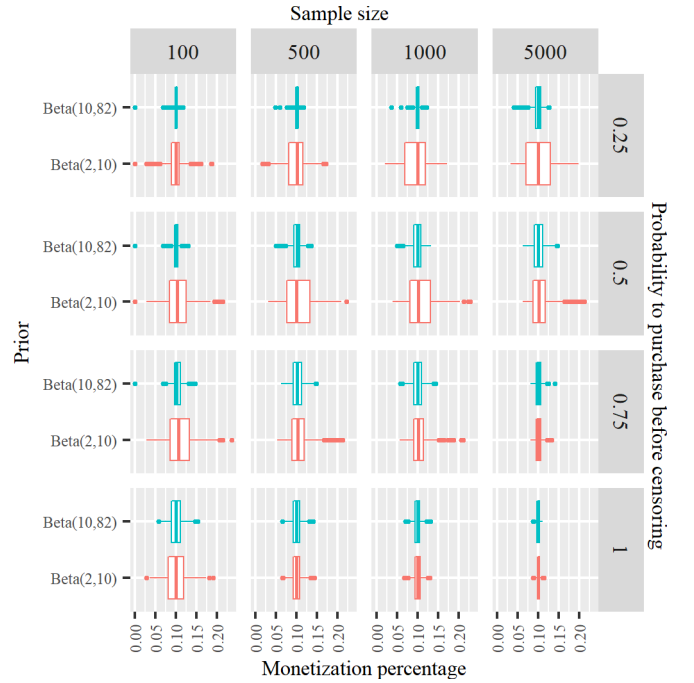


Fig. 6. Boxplots of the estimates of monetization percentage π .

The effect of the prior on the bias can be seen in Fig. 7, where the results of the second simulation study are represented. First the values of the relative bias were shifted by 1 so that all of them were positive and then logarithm was taken. Again it can be seen that the stronger the prior and the less there are monetized individuals in the data, the more the prior affects the estimate. We can see in the subplots on the diagonal of Fig. 7 that the most suitable prior distribution is the one with mode equal to the real monetization percentage.

Based on the simulation studies we can say that the additional prior information improves the method in that it decreases the variance and if it is appropriately chosen, it also helps to decrease the bias. Strong inappropriate prior might

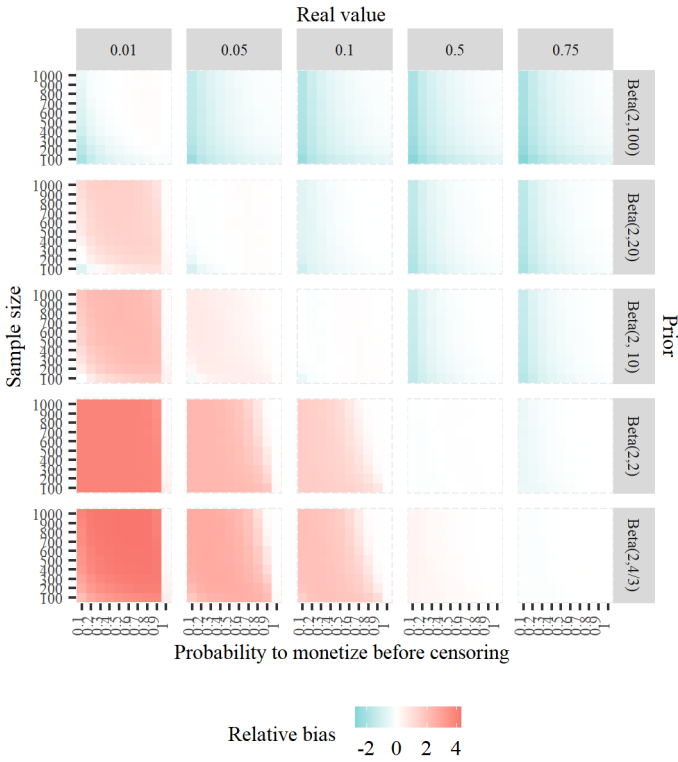


Fig. 7. Relative bias of monetization percentage π on a logarithmic scale as a function of censoring time and sample size.

bias the estimate a lot. If the mode of the prior distribution is a lot smaller than the real monetization percentage, it is anticipated that there is sufficient information in the data that decreases the effect of the prior and the estimate might only slightly underestimate the real value. Then again, if the mode of the prior distribution is a lot bigger than the real value, the sample information will not affect the estimate until the complete data case and the real value is greatly overestimated. Caution is needed when choosing the prior.

B. Game data

The results for the game data are shown in Fig. 8. A prior distribution $\text{Beta}(2, 50)$, that corresponds to adding 1 monetizing and 49 unmonetizing players, is used. We show both the computed estimate $\hat{\pi}$ and the percentage of monetized players so far at each censoring time for every game version. Obviously the percentage of monetized players increases as there are new monetized players. The value decreases if the number of monetized players does not increase at the same rate as the total number of players increases. At first, the predicted monetization percentages are greater than the observed monetized percentages, but as the follow-up time is enlarged, both estimates become equal to the supposed true value.

In Fig. 7 we can see that a prior with a greater value for parameter β is better for predicting the small monetization percentages. Now the prior is chosen so that its mode is 0.02, which is only a little higher than most of the real monetization

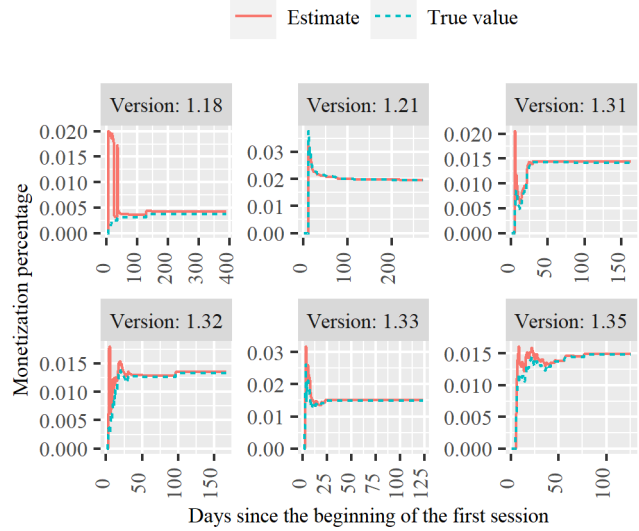


Fig. 8. Estimates for the game data with $\text{Beta}(2, 50)$ prior.

percentages shown in table I. The monetization percentage of version 1.18 is still greatly overestimated at the beginning but for the other versions the results seem to be similar to the ones shown before in [4]. Based on our simulation study, the peak in the beginning again seems to be caused by prior having a stronger effect than data on the estimate. Later when the estimate is greater than the true value, it is because the model predicts that there are players that will monetize although they have not yet done so.

Despite that we assumed that the number of event times that are not along the lines in figure 4 is sufficiently small for satisfying the assumption of the event times following exponential distribution, those purchase times may affect the reliability of the method. Some other Weibull distribution could be considered instead of the exponential distribution. Another potentially beneficial modification could be that the observed time was the total time spent playing the game instead of the calendar time from the beginning of the first session.

C. CDNOW data

The results for CDNOW data are presented in Fig. 9. We can see that there are peaks in the percentage estimate in the beginning when using informative priors. This is again caused by that there is more information in the prior than in the data. We can see that the stronger the prior is, the longer it affects the estimate. When the data starts to affect the estimate, the estimated percentage is only little higher than the true value.

It is a surprising result that the method does not predict better for this data set although there are a lot more data and the percentage of returning customers is a lot bigger than the monetization percentage in free-to-play games. We thought that small sample size and very small monetization percentage are the main reasons, why the model is not able to make better

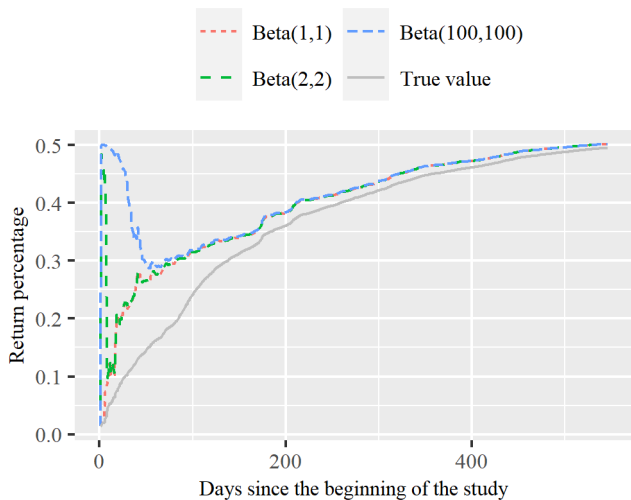


Fig. 9. Estimates for the CDNOW data with several priors compared to the real percentage of monetized individuals.

predictions, but based on the experiment with this data set, those qualities of data do not seem to affect the results that much. Then again, assuming the event times to follow another distribution could result in better estimates.

VI. CONCLUSION

The goal of this paper was to investigate a method that can be used to predict the percentage of monetizing players in limited follow-up data sets. The method uses a survival analysis based mixture cure model together with the Expectation Maximization algorithm, which results in an iterative algorithm that returns the monetization percentage and the conversion rate in an unlabelled data set. Most real-world data sets are probably like this, because game developers would like to use game analytics to improve the game during development or shortly after launch. Our approach suggested a new way to develop predictive models for this task.

The updated version of the method also takes prior information into account. We improved the idea of using pseudo-observations by replacing it with Bayesian inference, which formalizes prior knowledge about a data set. This approach decreased both the bias and the variance that were observed in the original method. The results of the simulation studies with generated data seemed promising: when a more informative prior was used, the variance decreased a lot and the estimates seemed to be less biased if the prior was well chosen.

However, with real data the prior information has an effect only in the beginning when very short follow-ups imply little information in the data. With sufficient information, the prior distribution scarcely affects the values of the estimates. The method was now also tested in the CDNOW data set. We formulated a similar prediction task by asking how many first time buyers become recurrent customers. Both the sample size and the return percentage are a lot higher in the CDNOW data set. We assumed that the method might work even better for

data with these qualities, but the results were very similar to those in the game data.

Our findings suggest that it is very difficult to predict the monetization percentage with unlabelled data sets when there is no extensive historical follow-up. The required parametric assumptions are close but do not exactly match real world data sets and this has a surprisingly large effect on predictive performance. Additionally, statistical theory guarantees that the parameters will converge to the true values, but in practise this happens extremely slowly as a function of sample size. The experiments suggest that parametric methods and asymptotic guarantees should be viewed with skepticism. We suggested a method of how unlabelled data can be handled, but additional research is needed to develop predictive models to this setting, which was found to be more difficult.

ACKNOWLEDGMENT

This work was supported by the Academy of Finland (grant 311273).

APPENDIX

We derive the iteration formulas (8) and (9) using the EM-algorithm in this section. These formulas find the maximum likelihood estimate in the incomplete data case. EM-algorithm finds the maximum a posteriori estimate $\hat{\Psi}_{\text{MAP}}$ for incomplete data by taking advantage of conditional expectation value of complete data likelihood function L_c and prior density function $p(\Psi)$ [24]. The algorithm consists of two parts:

- 1) In E-step the conditional expectation of the log complete data posterior density is calculated:

$$\mathbb{E}_{\Psi^{(k)}} \{\log p(\Psi|\mathbf{x})|\mathbf{y}\} = Q(\Psi; \Psi^{(k)}) + \log p(\Psi), \quad (11)$$

where

$$Q(\Psi; \Psi^{(k)}) = \mathbb{E}_{\Psi^{(k)}} \{\log L_c(\Psi) | \mathbf{y}\} \quad (12)$$

estimates the missing data by taking the conditional expectation of the complete data likelihood function.

- 2) In M-step a vector $\Psi^{(k+1)} \in \Omega$ which maximizes the (11) is found.

These two phases are repeated until convergence is achieved. It is shown in [25] and [31] that the algorithm converges to a local maximum of the likelihood function. In our case the log-likelihood function is concave, which implies that there is only one maximum and it is the global maximum. The additional term, the prior density function, in (11) almost always makes it more concave [24]. Thus these estimates converge to the global maximum of the log complete data posterior density function when $k \rightarrow \infty$. We first derive a formula for (12), then a formula for (11) and finally find the maximum by finding the roots of the partial derivatives of (11) with respect to π and λ separately. At first the conditional probabilities are shown in table III.

TABLE III
PROBABILITIES OF EVENT TIME T CONDITIONED ON MONETIZING INDICATOR ζ

$P(T \zeta)$	$T = t$	$T > t$
$\zeta = 1$	0	1
$\zeta = 0$	$\lambda e^{-\lambda t}$	$e^{-\lambda t}$

Then the formulas of marginals of survival and density functions in this mixture cure model case are

$$\begin{aligned}
S(t) &= P(T > t) \\
&= P(T > t | \zeta = 0) P(\zeta = 0) + P(T > t | \zeta = 1) P(\zeta = 1) \\
&= e^{-\lambda t} \cdot \pi + 1 \cdot (1 - \pi) \\
&= 1 - \pi + \pi e^{-\lambda t}
\end{aligned} \tag{13}$$

and

$$\begin{aligned}
f(t) &= P(T = t) \\
&= P(T = t | \zeta = 0) P(\zeta = 0) + P(T = t | \zeta = 1) P(\zeta = 1) \\
&= \lambda e^{-\lambda t} \cdot \pi + 0 \cdot (1 - \pi) \\
&= \pi \lambda e^{-\lambda t}.
\end{aligned} \tag{14}$$

Now the formulas of incomplete data likelihood and log-likelihood functions are functions (4) and (5) having functions (13) and (14) substituted.

The EM-algorithm requires the complete data likelihood function which is

$$\begin{aligned}
L_c(\Psi | \mathbf{t}, \delta, \zeta) &= \prod_{i=1}^n (\pi \lambda e^{-\lambda t_i})^{1-\delta_i} \\
&\quad \cdot \left[(1 - \pi)^{\mathbb{I}(\zeta_i=1)} (\pi e^{-\lambda t_i})^{\mathbb{I}(\zeta_i=0)} \right]^{\delta_i}
\end{aligned} \tag{15}$$

and the logarithm of it is

$$\begin{aligned}
l_c(\Psi | \mathbf{t}, \delta, \zeta) &= \sum_{i=1}^n \{ (1 - \delta_i) [\log \pi + \log \lambda - \lambda t_i] \\
&\quad + \delta_i [\mathbb{I}(\zeta_i = 1) \log(1 - \pi) \\
&\quad + \mathbb{I}(\zeta_i = 0) (\log \pi - \lambda t_i)] \}.
\end{aligned} \tag{16}$$

The last thing needed to define (12) is the probability to be an unmonetizing player conditioned on purchase time. These probabilities are calculated with the Bayes' theorem:

$$\begin{aligned}
P(\zeta_i = j | T > t_i, \Psi^{(k-1)}) \\
&= \frac{P(T > t_i | \zeta_i = j, \Psi^{(k-1)}) P(\zeta_i = j)}{P(T > t_i | \zeta_i = 0, \Psi^{(k-1)}) P(\zeta_i = 0) + P(T > t_i | \zeta_i = 1, \Psi^{(k-1)}) P(\zeta_i = 1)}
\end{aligned} \tag{17}$$

and

$$\begin{aligned}
P(\zeta_i = j | T = t_i, \Psi^{(k-1)}) \\
&= \frac{P(T = t_i | \zeta_i = j, \Psi^{(k-1)}) P(\zeta_i = j)}{P(T = t_i | \zeta_i = 0, \Psi^{(k-1)}) P(\zeta_i = 0) + P(T = t_i | \zeta_i = 1, \Psi^{(k-1)}) P(\zeta_i = 1)},
\end{aligned} \tag{18}$$

and the results are shown in Table IV.

TABLE IV
PROBABILITIES OF MONETIZING INDICATOR ζ CONDITIONED ON EVENT TIME T

$P(\zeta T)$	$\zeta = 1$	$\zeta = 0$
$T > t$	$\frac{1-\pi}{1-\pi+\pi e^{-\lambda t}}$	$\frac{\pi e^{-\lambda t}}{1-\pi+\pi e^{-\lambda t}}$
$T = t$	0	1

The resulting function (12) in a reduced form is

$$\begin{aligned}
Q(\Psi | \Psi^{(k-1)}) &= \mathbb{E}_{\zeta|T, \Psi^{(k-1)}} [l_c(\Psi | \mathbf{t}, \delta, \zeta)] \\
&= \sum_{i=1}^n \sum_{j=0}^1 P(\zeta_i = j | T > t_i, \Psi^{(k-1)})^{\delta_i} \\
&\quad \cdot P(\zeta_i = j | T = t_i, \Psi^{(k-1)})^{1-\delta_i} l_c(\Psi | t_i, \delta_i, \zeta_i) \\
&= \sum_{i:\delta_i=0} \left[\frac{1-\pi^{(k-1)}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} \log(1-\pi) \right. \\
&\quad \left. + \frac{\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} (\log \pi - \lambda t_i) \right] \\
&\quad + \sum_{i:\delta_i=1} [\log \pi + \log \lambda - \lambda t_i].
\end{aligned} \tag{19}$$

Now that $\pi \sim \text{Beta}(\alpha, \beta)$, the prior density function is

$$p(\pi) = \frac{1}{B(\alpha, \beta)} \pi^{\alpha-1} (1 - \pi)^{\beta-1}, \tag{20}$$

where $B(\alpha, \beta)$ is a beta function. By combining (19) and logarithm of (20) we get

$$\begin{aligned}
\mathbb{E}_{\Psi^{(k)}} \{ \log p(\Psi | \mathbf{x}) | \mathbf{y} \} &= \\
&\sum_{i:\delta_i=0} \left[\frac{1-\pi^{(k-1)}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} \log(1-\pi) \right. \\
&\quad \left. + \frac{\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} (\log \pi - \lambda t_i) \right] \\
&\quad + \sum_{i:\delta_i=1} [\log \pi + \log \lambda - \lambda t_i] \\
&\quad + \log \frac{1}{B(\alpha, \beta)} + \log \pi^{\alpha-1} + \log (1 - \pi)^{\beta-1}
\end{aligned} \tag{21}$$

Finally the formulas (8) and (9) are obtained by solving the roots of the partial derivatives of (21) with respect to λ and π separately.

REFERENCES

- [1] "2018 year in review," https://adindex.ru/files2/access/2019_01/230617_SuperData%202018%20Year%20in%20Review.pdf, SuperData, A Nielsen Company, 2019, accessed: May 2019.
- [2] AppsFlyer, "The state of in-app spending," http://cdn2.hubspot.net/hubfs/597489/IAP_Guide/The_State_of_In-App_Spending_AppsFlyer.pdf, 2016.
- [3] E. Lee, Y. Jang, D.-M. Yoon, J. Jeon, S.-i. Yang, S. Lee, D.-W. Kim, P. P. Chen, A. Guitart, P. Bertens *et al.*, "Game data mining competition on churn prediction and survival analysis using commercial game log data," *IEEE Transactions on Games*, 2018.

- [4] R. Numminen, M. Viljanen, and T. Pahikkala, "Predicting the monetization percentage with survival analysis in free-to-play games," in *2019 IEEE Conference on Games (CoG)*. IEEE, 2019.
- [5] J. Metcalf and K. Crawford, "Where are human subjects in big data research? the emerging ethics divide," *Big Data & Society*, vol. 3, no. 1, p. 2053951716650211, 2016.
- [6] B. Mikkelsen, C. Holmgård, and J. Togelius, "Ethical considerations for player modeling," in *Workshops at the Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [7] E. L. Neely, "Come for the game, stay for the cash grab: The ethics of loot boxes, microtransactions, and freemium games," *Games and Culture*, 2019.
- [8] U. Endriss and J. Leite, "Predicting players behavior in games with microtransactions," in *STAIRS 2014: Proceedings of the 7th European Starting AI Researcher Symposium*, vol. 264. IOS Press, 2014, p. 230.
- [9] H. Xie, S. Devlin, D. Kudenko, and P. Cowling, "Predicting player disengagement and first purchase with event-frequency based data representation," in *2015 IEEE Conference on Computational Intelligence and Games (CIG)*. IEEE, 2015, pp. 230–237.
- [10] R. Sifa, F. Hadiji, J. Runge, A. Drachen, K. Kersting, and C. Bauckhage, "Predicting purchase decisions in mobile free-to-play games," in *Eleventh Artificial Intelligence and Interactive Digital Entertainment Conference*, 2015.
- [11] T. Debeauvais and C. V. Lopes, "Gate me if you can: The impact of gating mechanics on retention and revenues in jelly splash." in *International Conference on the Foundations of Digital Games*, 2015.
- [12] M. Viljanen, A. Airola, J. Heikkonen, and T. Pahikkala, "Playtime measurement with survival analysis," *IEEE Transactions on Games*, vol. 10, no. 2, pp. 128–138, 2018.
- [13] P. S. Fader and B. G. Hardie, "Probability models for customer-base analysis," *Journal of interactive marketing*, vol. 23, no. 1, pp. 61–69, 2009.
- [14] N. Hanner, K. Heppner, and R. Zarnekow, "Counting customers in mobile business - the case of free to play." in *PACIS*, 2015, p. 174.
- [15] M. Viljanen, A. Airola, A.-M. Majanoja, J. Heikkonen, and T. Pahikkala, "Measuring player retention and monetization using the mean cumulative function," *arXiv preprint arXiv:1709.06737*, 2017.
- [16] S. Voigt and O. Hinz, "Making digital freemium business models a success: Predicting customers' lifetime value via initial purchase information," *Business & Information Systems Engineering*, vol. 58, no. 2, pp. 107–118, 2016.
- [17] D. R. Cox and D. Oakes, *Analysis of survival data*. Chapman and Hall Ltd, 1984.
- [18] R. Sifa, C. Bauckhage, and A. Drachen, "The playtime principle: Large-scale cross-games interest modeling," in *2014 IEEE Conference on Computational Intelligence and Games*. IEEE, 2014, pp. 1–8.
- [19] A. Isaksen, D. Gopstein, and A. Nealen, "Exploring game space using survival analysis." in *International Conference on the Foundations of Digital Games*, 2015.
- [20] R. G. Miller Jr, *Survival analysis*. John Wiley & Sons, 1981, vol. 66.
- [21] D. F. Moore, *Applied Survival Analysis Using R*. Springer International Publisher Switzerland, 2016.
- [22] J. Klein, H. van Houwelingen, J. Ibrahim, and T. Scheike, *Handbook of Survival Analysis*, ser. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press, 2016. [Online]. Available: <https://books.google.fi/books?id=t1vOBQAAQBAJ>
- [23] X. Liu, *Survival Analysis: Models and Applications*. Higher Education Press, 2012.
- [24] G. McLachlan and T. Krishnan, *The EM algorithm and extensions*. New York: John Wiley & Sons, 1997.
- [25] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977. [Online]. Available: <http://www.jstor.org/stable/2984875>
- [26] D. Hiemstra, *Probability Smoothing*. Boston, MA: Springer US, 2009, pp. 2169–2170. [Online]. Available: https://doi.org/10.1007/978-0-387-39940-9_936
- [27] J. O. Berger, *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag New York Inc., 1985.
- [28] R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2019. [Online]. Available: <https://www.R-project.org/>
- [29] P. S. Fader and B. G. Hardie, "Forecasting repeat sales at cdnow: A case study," *INFORMS Journal on Applied Analytics*, 2001. [Online]. Available: <https://doi.org/10.1287/inte.31.3s.94.9683>
- [30] H. Akaike, "A new look at the statistical model identification," *IEEE Transactions on Automatic Control*, vol. 19, no. 6, pp. 716–723, 1974.
- [31] C. F. J. Wu, "On the convergence properties of the EM algorithm," *The Annals of statistics*, pp. 95–103, 1983.