
SURVIVAL RATE ANALYSIS - DATA VALORIZATION FOR FINTECH

MSC. DATA SCIENCE AND ECONOMICS

Gregorio Luigi Saporito

Department of Economics, Management and Quantitative Methods
Department of Computer Science
University of Milan
Milan, Italy

Tommaso Pessina

Department of Economics, Management and Quantitative Methods
Department of Computer Science
University of Milan
Milan, Italy

October 26, 2020

ABSTRACT

The aim of this analysis is the one of determine if a person will conclude the subscription process to a virtual piggy bank application and/or where he is stuck and try to learn something from this kind of data, like the time to conclude the process, the percentage of people that will really conclude the process (with respect to those who are stuck in the process) and try to build a model for predict if a people will conclude or not the subscription (basing on his/her behavior). The data were pre-anonymized by the company that provided them.

Keywords R · Statistics · Survival Analysis · Logistic Regression · More

1 Introduction

The original dataset provided by the company was organised as follow:

- ID: unique identifier for the event;
- COMPLETED_STEP: identify the completed step of the event in JSON format;
- DATE_EVENT: it contains the date and the time of the event;
- NETWORK_ID: represent the network by which the event occurred;
- USER_ID: unique identifier for the user;

Each step represent a different passage of the subscription process and they are:

- anagrafica-a, anagrafica-b, anagrafica-c, anagrafica-d: this are the step that ask to the client the personal data;
- codice fiscale: here the user should provide his/her fiscal code;
- antiriciclaggio: in this step is ask to the user if he/she is a politics or an "exposed" person and the source of his/her money;
- conclude: if a user arrives here he/she has concluded the subscription process;
- contract-activation, contract-subscription.

Starting from this dataset one can think about which kind of information it can deliver, but for see it we should run some operation over it. As matter of fact we have grouped users with their id in order to calculate all of this information. In more details, we create a new dataset which include:

- **CLIENT_id**: it is the unique identifier of the user;
- **surv_time**: we calculate the difference between the **DATE_EVENT** of the first step of a certain user and the **DATE_EVENT** of his last step in the process, and we see it as hour;
- **status**: is a flag that will assume value 1 if a user arrive through the step "conclude";
- **extra_attempts**: indicate if the user has done other attempt of subscription;
- **main_extra_attempt**: it will assume a value different from "none" if the previous item is different from zero and it represent the last step in which he/she is stuck;
- **stuck**: contain the step (of the subscription process) in which the user is stuck;
- **mean_time**: it is the mean time that the user spent in each of the steps that he/she has concluded.

All our analysis is based on this new dataset.

2 Survival analysis

In this chapter we will discuss about the survival analysis with the aim to see how many people go all the way through the subscription process (along a time scale) or where they are stuck. Said that, we can also see if a person will conclude the subscription after a while of begin stuck.

First of all, we see that there are many outliers for survival time and few clients go through the system for a long time. After that, we can visually analyze in Figure 1 where the client are stuck in the subscription process.

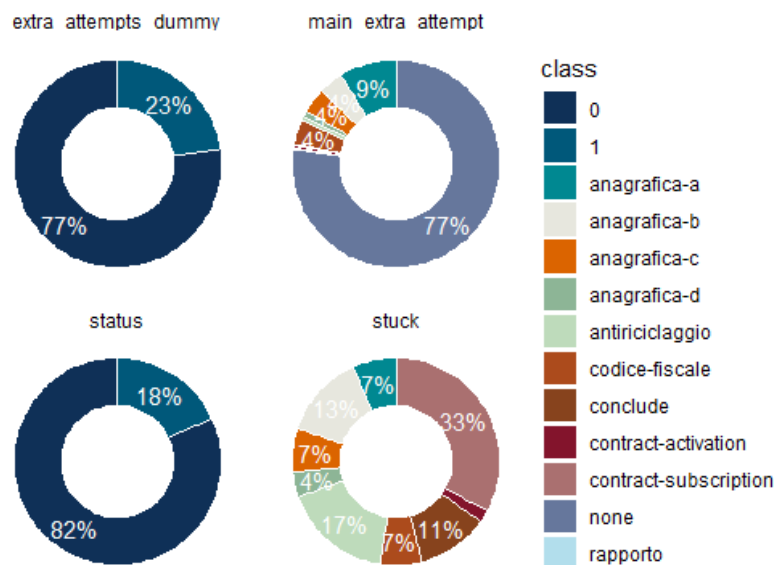


Figure 1: Donut graph of stuck users

Before discuss in detail our survival analysis, we should say few word about which kind of method we used. We used the "Kaplan–Meier estimator" that is a method used (especially in medical research) to estimate the probability that something will survive over time.

This is the basis of our analysis, that we will discuss in the next subchapter.

2.1 Kaplan-Meier estimates of the probability of survival over time

The first thing that we can say is that, as shown in Figure 2, as time passes the probability of concluding increases.

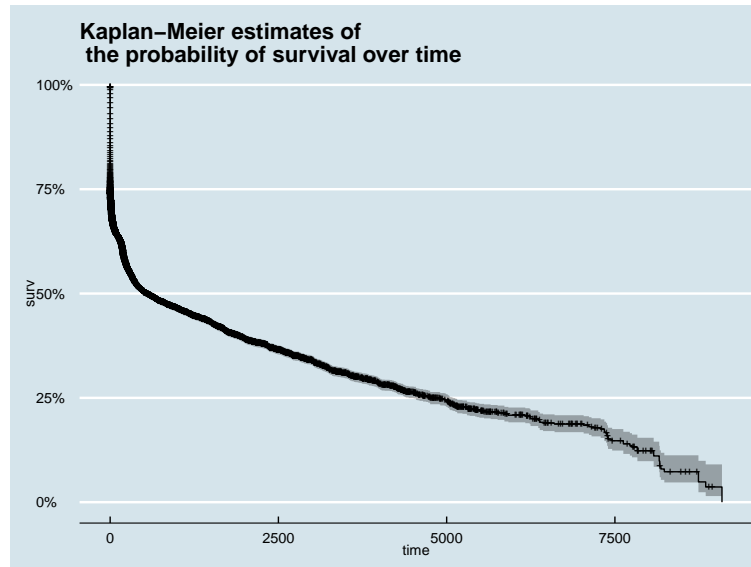


Figure 2: Probability of survival over time

Said that, we can now we look at the survival curve between who makes extra attempts and who does not. As shown in Figure 3, who makes extra attempts is more likely not to conclude and there are small overlaps between the curves.

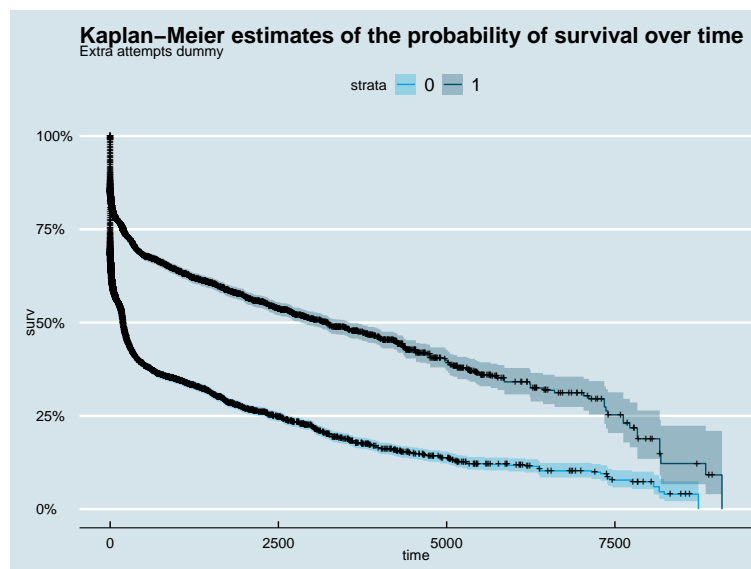


Figure 3: Probability of survival over time with extra attempt

Now, we can compare groups where the main stuck occurs and remove the ones with a too high confidence interval but there are still a lot of overlapps with confidence intervals though, as we can see in Figure 4.

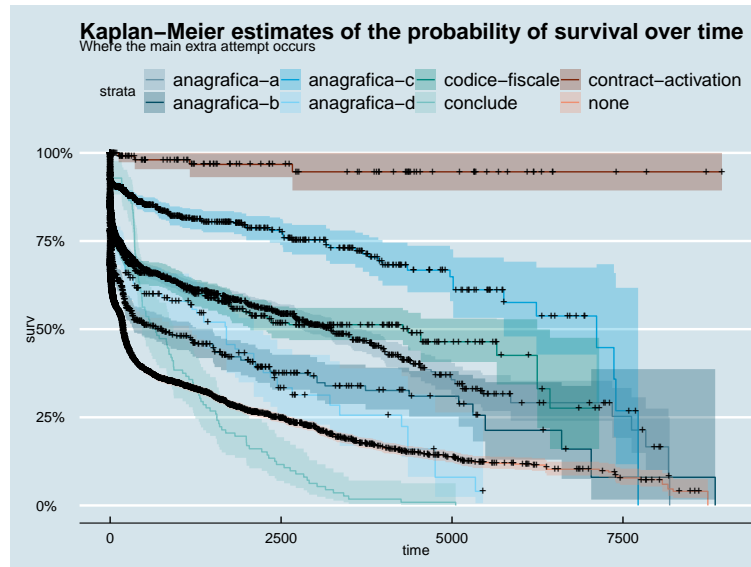


Figure 4: Total probability of survival over time with extra attempt

Finally, we look at where the client was last stuck in Figure 5.

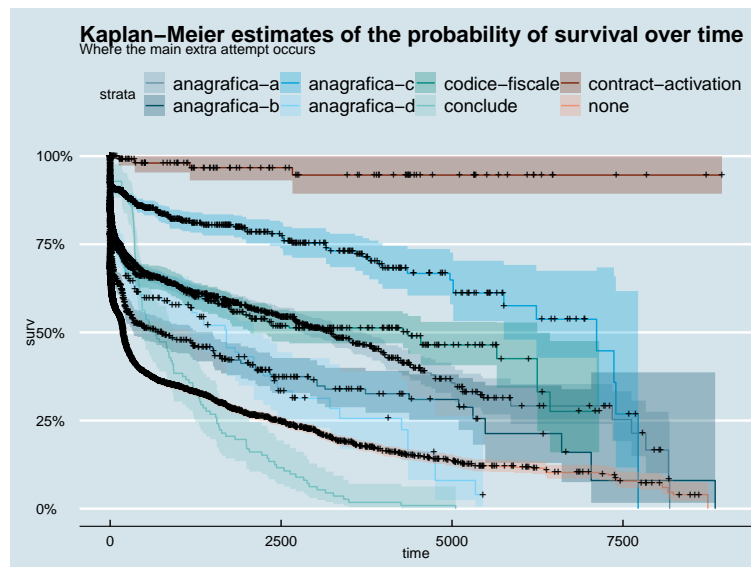


Figure 5: Probability of survival over time where the client is stuck

3 Predict the probability of conclude the subscription

logistic

4 Conclusion