# Modelling time-to-dropout via shared frailty Cox models. A trade-off between accurate and early predictions

Chiara Masci, Marta Cannistrà & Paola Mussida

Published online: 01 Sep 2023.

Submit your article to this journal ↗

View related articles ↗

View Crossmark data ↗

Routledge
Taylor & Francis Group

Check for updates

# Modelling time-to-dropout via shared frailty Cox models. A trade-off between accurate and early predictions

Chiara Masci [a], Marta Cannistrà [a,b] and Paola Mussida [c]

aDepartment of Mathematics, MOX – Modelling and Scientific Computing, Politecnico di Milano, Milan, Italy; bSchool of Management, Politecnico di Milano, Milan, Italy; cDepartment of Electronics, Information and Bioengineering, Politecnico di Milano, Milan, Italy

**ABSTRACT**

This paper investigates the student dropout phenomenon in a technical Italian university from a time-to-event perspective. Shared frailty Cox time-dependent models are applied to analyse the careers of students enrolled in different engineering programs with the aim of identifying the determinants of student dropout through time, predicting the time to dropout as soon as possible and to observe how the dropout phenomenon varies across time and degree programs. The innovative contributions of this work are methodological and managerial. First, the adoption of shared frailty Cox models with time-varying covariates is relatively new to the student dropout literature and it allows to consider the student career evolution and the heterogeneity across degree programs. Second, understanding the dropout pattern over time and identifying the earliest moment for obtaining its accurate prediction allow policy makers to set timely interventions for students at risk of dropout.

## 1. Introduction and motivation

The Italian Higher Education system is affected by a high level of dropout, with many students abandoning their Bachelor programs during the first or second year. According to ANVUR (Italian National Agency for the Evaluation of Universities and Research Institutes), the Italian dropout rate for the students from whom complete data is available is around 24%, with half (12%) of them dropping out in the first two years (ANVUR 2018).

This data is even more worrying considering that only 28% of graduates in Italy are from the 25–34 years old population, against a European average of 40 University dropout represents a worrisome phenomenon with both economic and social impacts.

From the economic standpoint, dropout represents a net waste of resources for universities, since education is a costly activity. From a social perspective, dropout affects students, who face a social stigma (e.g. fewer job opportunities and lower salaries), disconnecting them from their social environment (Alban and Mauricio 2019).

Hence, studying the dropout phenomenon and its determinants is paramount. Identifying students at risk and, in particular, the riskiest moment of their career is extremely important: only with timely interventions, universities would be able to retain their students, shepherding them towards graduation (Seidel and Kutieleh 2017).

Aligned with the motivations outlined, this paper has a dual objective. Firstly, it aims to investigate student dropout and the significant factors influencing it over time. Secondly, it seeks to identify the earliest point in time at which accurate dropout predictions can be obtained, considering both the occurrence and the timing of the event. These predictions can be used by universities to facilitate early interventions through appropriate preventive actions. The focus is on identifying not only who the students at risk are, but also when these students are at risk, discussing the effectiveness of an Early Warning System, seen as a tool for the early detection of at-risk students. The study is held at Politecnico di Milano (PoliMi) in Italy.

The paper is organized as follows: in Section 2, we present the Italian educational system and context of this study. We set up an overview of the academic literature about survival analysis and dropout in Section 3. In Section 4 we present the main features of the PoliMi dataset and the methodology adopted. Results and final considerations are detailed in Sections 5 and 6.

## 2. The context

### 2.1. University system

Since the implementation of the Bologna reform in the early 2000s, the university education system has been structured into three cycles, each corresponding to a degree: the bachelor's degree in the 1st cycle, the master's degree in the 2nd cycle, and the PhD in the 3rd cycle. As part of the reform, a university training credit (European Credit Transfer and Accumulation System, acronym ECTS) is introduced to measure the workload required of students to obtain a degree. This system simplifies the recognition of exams taken in other Italian or European universities, allowing the transfer of credits through the ECTS. Each university exam weighs a certain number of credits, with one credit equivalent to 25 h of student commitment. In an academic year, students typically earn 60 training credits, resulting in 180 credits for bachelor's degree and 120 for a master's degree.

It is important to note that there is no national centralized enrollment system for all these institutions, and each university is accountable for managing its own student data: universities collect and send data to the Ministry of Education. As a result, the combination of data across universities is a challenging task, so every university has a limited perspective on students' academic journeys. If a student leaves one university, it may be interpreted as either dropping out of the education system or enrolling in another institution. In this paper, the concept of dropout refers to leaving Politecnico di Milano rather than the educational system in general.

### 2.2. Our university: Politecnico di Milano

Politecnico di Milano provides study programs in Architecture, Design and Engineering and has an approximate enrollment of 47,000 students. The university has an internal IT department dedicated to the development of the necessary services for managing and reprocessing administrative data. The university's fees are determined by the income bracket of the student's family, and study grants are awarded considering a combination of family income and student merit.

Each academic year is divided into two semesters: the first spans from September to December, and the second from February to May. At the conclusion of each semester, there are the exam sessions. Two sessions are available shortly after the lessons' completion, followed by up to five additional ones. Students have the freedom to retake the exams, declining grades, without a restriction on the number of times.

Given the high dropout rate (28%) at Politecnico di Milano, the Rector is now interested in reducing the number of students leaving the university every year. To this end, the development of an Early Warning System to timely predict dropouts allows university managers to implement remedial interventions to retain those students.

## 3. Related literature: survival analysis for studying dropout

As part of the wide academic literature aiming at predicting dropout in education settings (Cannistrà et al. 2022; Hegde and Prageeth 2018; Kehm, Larsen, and Sommersel 2019), survival analysis is directed toward the deepening of when this event occurs, considering students' educational career complemented with its time dimension.

Singer and Willett (1993) were among the first to develop a model on Cox's seminal article (1972) for educational contexts to study discrete-time survival analysis. The idea behind this family of methods is to answer to research problems, such as dropout, concerning whether and when an event of interest occurs. If before the discrete time to event adapted well to educational problems, periodically observed, nowadays the information available about students, allows to consider it as continuous over time. Indeed, many studies applying survival analysis focus on digital learning (Xie 2020; Utami et al. 2020; Spitzer et al. 2021; Chen et al. 2020) for a simple reason: it is easy to track students over time. Researchers interested in studying whether and when dropout from the platform (i.e. last access) occurs may know the exact instant of time of this event. Another stream of literature is centered on the doctoral path: when and why PhD students drop out from their career (Van Der Haert et al. 2014; Booth and Satchell 1995; De Valero 2001; Grove, Dutkowsky, and Grodner 2007). Indeed, the investigation of this phenomenon is relevant since it gives the opportunity to understand the most effective type of support for retaining PhD students, given their value in our society (Van Der Haert et al. 2014).

Focusing on schools and universities, the academic contributions applying survival analysis to students' academic career progression aim at modelling the phenomenon by highlighting the most important underlying factors (Arulampalam, Naylor, and Smith 2004; Weybright et al. 2017; Thaithanan et al. 2021; Patacsil 2020; Min et al. 2011; Plank, DeLuca, and Estacion 2008; Barragaan, Gonza´lez, and Calder´on 2022; Vallejos and Steel 2017; No, Taniguchi, and Hirakawa 2016; Gury 2011; Lesik 2007).

Arulampalam, Naylor, and Smith (2004) and Barragaan, Gonza´lez, and Calder´on (2022) found academic performance to be an important dropout predictor, while according to Weybright et al. (2017), the student's background (e.g. being a male and not living with his mother) plays a significant role in predicting dropout (Barragaan, Gonza´lez, and Calder´on 2022). Soares et al. (2015) observed that the difficulties faced with particular subjects, the desire for a different school, the perception that those completing their studies will have better job opportunities, and the importance assigned to school choice influence dropout from secondary school. When looking at the university dropout phenomenon's time component, Min et al. (2011) found significant differences for early semesters across groups. White and/or female students tend to leave university earlier than other sub-populations. Engineering students mostly abandon their academic career during the third semester, but this can happen even during the second semester when the student has a low math grade.

In terms of adopted models, the majority of scholars (Weybright et al. 2017; Thaithanan et al. 2021; Min et al. 2011; Plank, DeLuca, and Estacion 2008; Barragaan, Gonza´lez, and Calder´on 2022; Vallejos and Steel 2017; Gury 2011; Lesik 2007; Arulampalam, Naylor, and Smith 2004) use Cox Proportional Hazards (PH) models to estimate the probability of dropping out, often comparing Kaplan-Mayer curves on different students' features. Interesting sources of innovation are related to the comparison between fixed and random effects, as in Arulampalam, Naylor, and Smith (2004), to model the effect of being enrolled in different degree programs; or to the combination between survival analysis and analytic hierarchy process methodologies, as in Barragaan, Gonza´lez, and Calder´on (2022), to model dropout as a decision subjected to multiple alternatives; or by handling covariates' selection within a Bayesian framework (Vallejos and Steel 2017). Generally, academic literature is moving toward modelling dropout and estimating its related factors with ever-increasing precision.

To contribute to this stream of research, this paper aims to study the dropout phenomenon with a time perspective, adding two sources of innovation. The first relates the methods adopted, where

the inclusion of frailties allows to account for the nested structure of students into degree programs, modelling the heterogeneity at the second level of the hierarchy, and the modelling of time-dependent covariates allows to update student information in time, building increasingly informed models. The second innovation regards the final collateral goal of the analysis: identifying the earliest moment in a student's career in which we can accurately predict his/her time-to-dropout. Indeed, early and accurate predictions are essential to effectively support at-risk students.

## 4. Data and methods

In the following two subsections we present the dataset and our methodological approach.

### 4.1. Polimi dataset

The PoliMi dataset contains administrative information about the careers of students enrolled between Academic Years 2010 and 2021 (12 years span period)[1] in Bachelor's degree programs of Engineering. The University collects information about students' demographics and previous studies and tracks their entire academic careers, making anonymized data available in real time (Mussida and Lanzi 2022). The demographics regard gender and age, residency and citizenship, and university's fee bracket paid by the student (as a proxy of socio-economic status). Then, high school track and final mark inform about student's previous career, while PoliMi admission test score is the first grade measured at the University. As regards career tracks, the number of credits obtained (ECTS, European Credit Transfer and Accumulation System) and the relative Grade Point Average (GPA) are collected for each student each semester. It is worth noting that while the observation of students' dropout may occur daily, their academic progression within the university is not continuously observed. Instead, exams are typically held at the end of each semester, resulting in periodic registration of grades and credits. Consequently, even though the outcome of interest is continuous, the covariates (such as grades and credits) are not.

The analysis excludes students who abandon their studies during the first semester of their first year (1700 students from 2010 to 2020, 16% of the total dropouts) since many students enroll at PoliMi while waiting to be admitted to other programs at other universities, or they immediately decide to abandon because they had different expectations. This heterogeneity behind these dropouts might bias the results and these are not the dropouts that we aim to identify and on which we want to act.[2]

The final dataset contains 49,501 students enrolled in 16 degree programs. Table 1 reports the selected student-level variables, collected at the time of enrolment, with their explanation and summary statistics. The target variable regards the status of the student's career at the end of the third year, which can be concluded with graduation, with a dropout or with the student still being active. Variables Status at 3y and Career duration at 3y, reported in Table 1, define the target variable. It is indeed important to consider the Career duration, defined by the time-to-dropout, as continuous, since students could formally leave the university at any point during their academic journey.[3] Table 1 reveals distinctive characteristics of the student sample enrolled at Politecnico di Milano. Specifically, 77% of the students are male, and 80% originate from a scientific high school background. Approximately one-third of the students pay the highest university fees, indicating that they come from affluent families. On average, the students achieved a grade of 75 out of 100 in their final high school exams. These findings suggest the possibility of self-selection among students before enrollment. The majority of Politecnico di Milano students exhibit strong academic performance and come from privileged socio-economic backgrounds. Regarding the career tracks, Table 2 reports the selected longitudinal information relative to each student's careers, semesterly updated, that are ECTS and GPA.

The distribution of the students within the 16 Engineering degree programs is reported in Table A2 in Appendix A2. For privacy reasons, we are not allowed to report the degree programs names but only their anonymized codes.

**Table 1.** Student-level variables adopted in the analysis, their description, type, and summary statistics.

| Name | Description | Type | Summary info |
|---|---|---|---|
| Gender | Student gender (F/M) | Categorical | Male = 77,5%, Female = 22,5% |
| Admission age | Age as of the day of enrolment | Numerical | mean = 18.72, median = 19, sd = 1.22, range = [16-61] |
| Income | University fee bracket: *High, Medium, Low* or *Study Grant (SG)* | Categorical | High = 32.8%, Medium = 23.5%, Low = 30.6%, SG = 13.1% |
| Origins | *Milanese* living in Milan, *Commuter* living outside Milan, *Offsite* have moved to Milan | Categorical | Commuter = 67.5%, Milanese = 25.7%, Offsite = 6.8% |
| Highschool type | Field of study at high school: *Scientific, Classic, Technical. Foreigner* if he/she got his/her diploma abroad and *Other* if none of the above | Categorical | Classic = 5.4%, Other = 1.6%, Scientific = 80.5%, Foreigner = 0.7%, Technical = 11.8% |
| Highschool grade | Grade obtained in high school | Numerical | Mean = 84.87, median = 85, sd = 11.61, range = [60-100] |
| Admission score | Score obtained on the PoliMi admission test | Numerical | Mean = 73.22, median = 71.55, sd = 9.36, range = [60-100] |
| Department | Study program of the student | Categorical | 16 faculties |
| Status at 3y | Student career status considering a follow up time of 3 years, grouped by G (graduated), A (active), and D (dropout) | Categorical | Graduated = 10.9%, Active = 77.4%, Dropout = 12.7% |
| Career duration at 3y | Length of the student career considering a follow up time of 3 years, expressed in semester | Numerical | Mean = 5.08, median = 6, sd = 1.47, range = [1,6] |

Note: The Table shows the descriptive statistics of the time-invariant covariates used in subsequent analysis. In detail, for categorical variables it shows the distribution in each category, for Numerical variables their mean, median, standard deviation (sd) and range. Variables Status at 3y and Career duration at 3y are used to build the outcome of interest.

## 4.2. Models and methods

In this subsection, we briefly recall the basics of survival analysis and we describe the statistical models adopted in the study.

### 4.2.1. Basics of survival analysis

Survival analysis regards the group of statistical procedures for the modelling of the time until an event of interest occurs (Kleinbaum and Klein 1996; David and Mitchel 2012). For each unit of analysis, the event (i.e. *dropout*) might occur during the follow-up (i.e. the period of observation – in our case, three years) or not. In the second case, we refer to the observation as censored. For each unit $i = 1, \ldots, N$, the target variable is defined as the couple of the survival time $T_i = \min(T_i^*, C_i)$ and the censoring indicator $\delta_i = (T_i^* \leq C_i)$, where $C_i$ is the censoring time and $T_i^*$ is the observed event time, if any. $\delta_i$ is the indicator function that indicates whether the event occurred ($\delta_i = 1$) or not ($\delta_i = 0$) for the individual $i$. Censoring is assumed independent of survival time. Being $T$ a non-negative random variable, the survival function.

$$S(t) = P(T > t) = 1 - P(T \leq t) = 1 - F(t) \tag{1}$$

represents the probability of survival until time $t$, while the hazard function describes the instantaneous risk of failure and is defined as

**Table 2.** Student-level variables related to the student career, measured each semester until the end of the third year.

| | Description | Type | Summary info |
|---|---|---|---|
| Exa_Ay | Academic year corresponding to the observation | Categorical | Range = 2010-2021, |
| Exa_Semester | Semester corresponding to the observation. 1 if first semester, 2 if second semester | Categorical | 1 = 58.6%, 2 = 41.4% |
| ECTS | ECTS obtained by the student during each semester | Discrete | mean = 18.48, median = 20, range = 0-40, sd = 12.6 |
| GPA | Weighted average grade measured for each semester | Numeric | mean = 18.97, median = 22.8, range = 0-30, sd = 10.43 |

Note: the Table shows the structure of the time-dependent dataset, in which the GPA and ECTS are measured within each semester and Academic Year.

$$h(t) = \lim_{\Delta t \to 0} \frac{P(t \le T \le t + \Delta t | T \ge t)}{\Delta t} \qquad (2)$$

The survival function $S(t)$ can be estimated through the Kaplan-Meier estimator (KM) (Kaplan and Meier 1958), which represents the probability of surviving in a given length of time while considering time in many small intervals. In the case of two or more groups, the Log- Rank Ratio test (Mantel 1966) can be used to test statistical differences across the estimated KM curves.

### 4.2.2. Shared frailty Cox PH models with time-invariant and time-varying covariates

Cox regression models are the most popular mathematical modelling approach to estimate the survival curves when considering several explanatory variables simultaneously. When the units are not *i.i.d.* but they are nested within groups, Shared Frailty Cox models introduce a frailty term, shared among units within the same group (in our case, students within degree programs), to take the structure into account (Kleinbaum and Klein 1996; David and Mitchel 2012).

The *Shared Frailty Cox Proportional Hazards (PH) model* assumes the hazard function for the $i$−th individual, for $i = 1, \dots, N$ within the $j$−th group, for $j = 1, \dots, J$, to be modelled as follows:

$$h_{ij}(t, \, \boldsymbol{x}_{ij}) = h_0(t) \times \omega_j \times exp\left\{ \sum_{p=1}^{P} \beta_p x_{p,ij} \right\} \qquad (3)$$

where $h_0(t)$ is the baseline hazard function, $\boldsymbol{\beta x}_{ij}$ is the linear predictor, where $\boldsymbol{x}_{ij}$ is the vector containing the P covariates relative to the $i$−th individual within the $j$-th group and $\boldsymbol{\beta}$ is the P-dimensional vector of corresponding coefficients, $\omega_j$ is the frailty term for the $j$−th group. To better quantify the effect of the covariates, Hazard Ratios (HRs) can be derived from the vector of coefficients $\beta$. The modelling is based on the following assumptions: the effect of each covariate is constant across time (PH assumption), all failure times are independent given the frailties, and the values of the random effects $\omega_j$ are constant over time and common to all the individuals belonging to the same group. The frailties $\omega_j$ have a positive unobserved multiplicative effect on the hazard function. They are *i.i.d.* following a Gamma distribution with $E(\omega) = 1$ and $Var(\omega) = \theta$, where $\theta$ is the unknown parameter. Larger values of $\theta$ mean greater heterogeneity among the groups. Individuals belonging to a group with $\omega_j > 1$ have an increased hazard and decreased probability of survival compared to those with average frailty ($\omega_j = 1$). Similarly, individuals belonging to a group with $\omega_j < 1$ have a decreased hazard and increased probability of survival compared to those with average frailty.

This modelling can be extended to handle time-varying covariates. The shared frailty Cox model with both time-invariant and time-varying covariates, with respect to the $i$-th individual within the $j$-th group, assumes the following form:

$$h_{ij}(t, \, \boldsymbol{x}_{ij}(t)) = h_0(t) \times \omega_j \times exp\left\{ \sum_{p=1}^{P} \beta_p x_{p,ij} \sum_{q=1}^{Q} \gamma_q x_{q,ij}(t) \right\} \qquad (4)$$

where P and Q are the number of time-invariant and time-varying covariates, respectively, and $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are the P-dimensional and Q-dimensional vectors of coefficients associated to these covariates, respectively. This modelling assumes that the effect of time-varying covariates $\boldsymbol{x}_q(t)$ on the survival probability at time $t$ depends on the value of this feature at time $t$ and not on its value at previous times. The PH assumption is no longer satisfied and the Hazard Ratio between two individuals $i$ and $j$ varies across time, depending on the covariates' values.

### 4.2.3. Goodness-of-fit indices

To evaluate the goodness of fit of our models, we rely on the most common metrics, the Concordance index (C-index) (Steck et al. 2007), which is defined as the proportion of concordant pairs, i.e. pairs of individuals for which the expected event times are predicted in the correct ordering,

divided by the total number of possible evaluation pairs. The closer to one, the more accurate the Cox model. We support the C-index with a further evaluation, obtained by treating our survival models as classification models: by looking at the estimated survival probability at a fixed time t*, we compute classification performance indices, e.g. precision, recall and ROC curve.

## 5. Results

The event of our interest is the *failure event of student dropout* from university. A follow up period of five semesters is considered: a student dropping out between the end of the first and sixth semester is labelled as *dropout*, while all other students, i.e. students who drop out after 3 years from the enrolment, who graduate or who have an active career at the end of the 3$^{rd}$ year, are marked as *censored*.

This section is divided into three main parts. In Section 4.1, we report the results of a preliminary analysis to describe the cohort of students and the dropout distribution across time. In Sections 4.2.1 and 4.2.2, we show the results of shared frailty Cox models, first with only the time-invariant covariates and, then, with the addition of the time-varying ones. Results focus on the interpretation of the effect of student-level characteristics on the dropout risk, on the quantification of the heterogeneity across degree programs and on the models' predictive power. Lastly, Section 4.3 reports a comparison of Cox models fitted by sequentially adding students' information in time in order to identify the best trade-off between accurate and early predictions.

### 5.1. Preliminary analysis

As reported in Table 1, 12.7% of the students in our sample dropped out during the five semesters after the first one. Figure 1 reports the estimated survival function and the distribution of the time to dropout, measured in semesters. As expected, most of the dropouts occur in correspondence with the end of academic years, mainly during the first two ones.

The hazard function presents three major peaks (represented by jumps in the survival function), which correspond to moments with a high frequency of dropout, at the end of each of the three academic years. The highest peak is in correspondence of the second semester, which marks the end of the first year of university; therefore, preventive interventions before this time are needed.

Appendix A3 reports a detailed descriptive univariate analysis to investigate the association between student characteristics and dropout risk.



**Figure 1.** Estimated survival function and time-to-dropout distribution. Note: The figure in the right panel reports the distribution of the times, expressed in semesters, in which students definitely abandon PoliMi during a follow-up of 5 semesters (3 years, except for the first semester). Mean = 3.02, median = 2.20). 0 corresponds to the enrolment.

## 5.2. Shared frailty Cox PH models

In this section, we fit two Shared Frailty Cox models, considering students (level 1) nested within degree programs (level 2), in order to estimate the student time to dropout between the end of first and sixth semester, by exploring the effects of student characteristics and of the degree programs. The first is a Shared Frailty Cox model with time-invariant covariates, while in the second time-varying covariates about students' academic results are added.

The inclusion of a random effect allows us to model the heterogeneity across degree programs and to quantify the effect of each degree program on the dropout risk of their students. With respect to the more common inclusion of a fixed effect, this modelling has three main advantages. First, the frailty takes account of the dependence between the students enrolled within the same degree course. We expect those students to share the same curriculum, the same teachers and the same environment. Fixed effects assume all individuals to be independent while taking account of this dependence is important and avoids biased estimates. Second, by using a Gamma frailty, we are not forced to take one of the degree programs as a reference and estimate the statistical differences between this category and the others, but we are able to measure how the effect of each degree program statistically differs from the average. From an interpretative point of view, this provides a valuable improvement. Third, in terms of prediction, from a shared frailty model we can predict the expected dropout time of a student, given his/her characteristics and net to the effect of the degree program attended and, then, evaluate how this estimate does change across the degree programs.

For both models, we randomly divide the dataset into training and test sets, containing 70% and 30% of the observations, respectively.
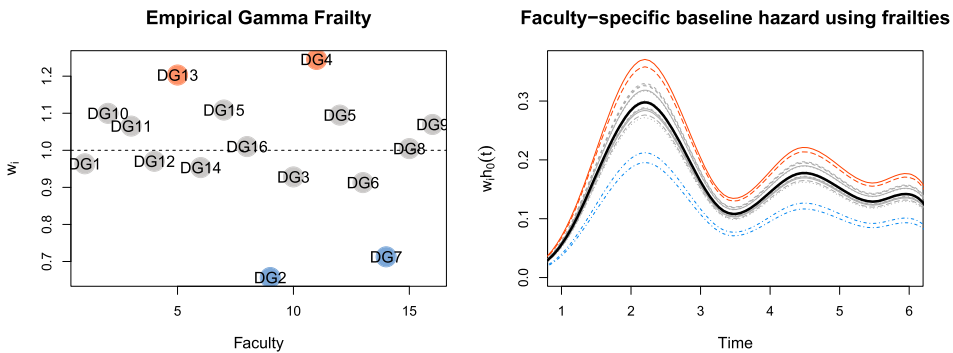
### 5.2.1. Shared Frailty Cox PH model with time-invariant covariates

The Shared Frailty Cox model includes as time-invariant covariates Gender, Income, Origins, HighschoolType, HighschoolGrade, AdmissionScore, Age19, and ECTS[4] of the first semester. Table 3 shows the summary of the model estimated on the training set. Results show that females have an average lower risk of dropout than males (HR = 0.84), and students with SG income category are less likely to drop out with respect to students in the Medium category (HR = 0.772). Commuters are more likely to drop out than Milanese students (HR = 1.144), being that a student who attended a Technical school or other types of high schools is associated with a

**Table 3.** Shared Frailty Cox model with time-invariant covariates, output of the summary.

| | Coefficient | Standard error | Hazard Ratio | 95%CI for HR |
|---|---|---|---|---|
| Gender:F | −0.177** | 0.042 | 0.840 | (0.77 - 0.91) |
| Income:High | 0.031 | 0.038 | 1.031 | (0.96 - 1.11) |
| Income:Low | −0.038 | 0.039 | 0.962 | (0.89 - 1.04) |
| Income:SG | −0.258** | 0.056 | 0.772 | (0.69 - 0.86) |
| Origins:Commuter | 0.135** | 0.034 | 1.144 | (1.07 - 1.22) |
| Origins:Offsite | −0.074 | 0.069 | 0.929 | (0.81 - 1.06) |
| HighschoolType:Classical | −0.013 | 0.062 | 0.987 | (0.87 - 1.12) |
| HighschoolType:Foreigner | −0.144 | 0.157 | 0.866 | (0.64 - 1.18) |
| HighschoolType:Others | 0.279** | 0.096 | 1.322 | (1.10 - 1.59) |
| HighschoolType:Technical | 0.085** | 0.045 | 1.088 | (1.00 - 1.19) |
| HighschoolGrade | −0.003* | 0.002 | 0.997 | (0.99 - 1.00) |
| AdmissionScore | −0.002 | 0.002 | 0.998 | (0.99 - 1.00) |
| Age19: > 19 | −0.050 | 0.046 | 0.951 | (0.87 - 1.04) |
| ECTSP | −0.123** | 0.002 | 0.884 | (0.88 - 0.89) |
| Number of events | 4,549 | | | |
| Observations | 34,651 | | | |
| Frailty | $\hat{\theta} = 0.029$ | $se(\hat{\theta}) = 0.0104$ | pval = 0.010 | |
| Concordance | 0.816 | | | |
| Log Likelihood | −17118.22 | | | |

Note: *$p < 0.1$; **$p < 0.05$. Estimated baseline survival and hazard functions are reported in Figure A4a in Appendix A4.

**Empirical Gamma Frailty**  **Faculty−specific baseline hazard using frailties**



**Figure 2.** Estimated frailty terms and degree programs-specific baseline hazard functions in the time-invariant case. Note: Left panel shows the empirical Gamma Frailty terms for the 16 degree courses estimated by the shared frailty Cox model with time-invariant covariates. Red and blue points identify the faculties that have a frailty term significantly higher and lower than 1, respectively. Right panel reports the faculty-specific baseline hazard functions for the 16 specific degree courses.

higher dropout risk with respect to students who attended Scientific schools (HR = 1.088 and 1.322, respectively), and the higher the high school final grade, the lower the risk of drop out, on average (HR = 0.997). Lastly, the number of credits obtained in the first semester is confirmed to be an important protective factor. This output confirms again how the early academic results obtained by the student have an important role in a student's choice to withdraw from studies. The admission score at PoliMi and the age as of enrolment are not result to be significant.

Regarding the degree program effect, $\hat{\theta} = 0.029$ is the estimated variance of the frailty parameter. The variance of the frailty term $\hat{\theta}$ is significantly different from 0 (*p*-value of the Wald test 0.01), confirming the presence of heterogeneity between degree programs. The estimated frailty terms $\omega_j$, $j = 1, \ldots, 16$, which denotes the effect of each particular study program on the baseline hazard function, are shown in the left panel of Figure 2. Among the 16 degree programs, two results are associated with a higher dropout risk with respect to the average, net to the effect of student characteristics ($\omega_{DG4} = 1.245$ and $\omega_{DG13} = 1.203$). On the opposite, two programs results are associated with lower dropout risks ($\omega_{DG2} = 0.656$ and $\omega_{DG7} = 0.712$). In the plot, the groups are colored depending on the asymptotic 95% confidence interval[5] [$\omega \pm 1.96 \times \sigma(\omega_j)$].

The impact of these estimated values on the survival probability can be easily visualized in the department-specific baseline hazard functions, as shown in the right panel of Figure 2.

In terms of model predictive performance, the C-index computed both on the training and test set are 0.816 and 0.814, respectively.

### 5.2.2. Shared Frailty Cox PH model with time-varying covariates

We now extend the previous model by including time-varying covariates. In particular, we consider GPA and ECTS measured at the end of each semester as time-progressive information. Model results are reported in Table 4.

By including the career tracks over time, some of the personal student characteristics change their significance with respect to the first model. Here, gender is no more significant; with respect to a Medium income, having a Low income and having a scholarship (SG) are protective factors; with respect to *Milanese* students, *Commuters* and *Offsite* students have on average a higher dropout risk; with respect to scientific high school, having attended a foreigner or a technical school is a protective factor; having obtained a good high school grade is a risk factor, and being a student older than the average is a protective factor. As regards the career track, both ECTS and GPA are very significant and are protective factors. It is worth noting that, net to the effect of progressive ECTS and GPA, we still observe many significant student characteristics.

Regarding the frailty term, its estimated variance $\hat{\theta} = 0.012$ again results to be significantly different from 0. The distribution of the 16 estimated frailties $\hat{\omega}_j$ and the program-specific baseline

**Table 4.** Shared Frailty Cox model with time-dependent covariates, output of the summary.

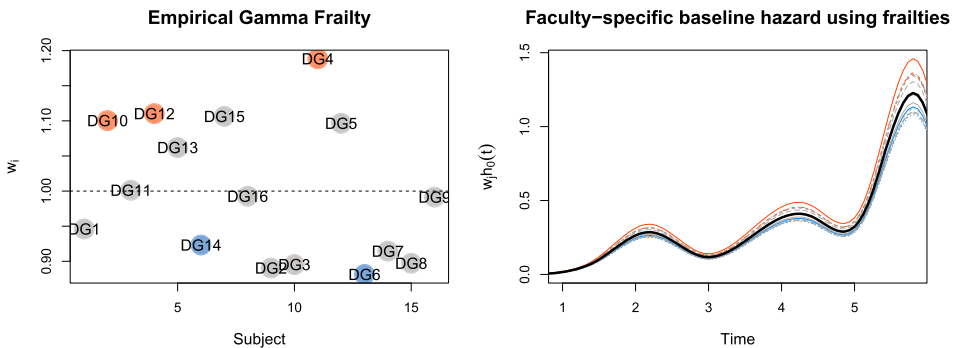|  | Coefficient | Standard Error | Hazard ratio | 95%*CI* for HR |
|---|---|---|---|---|
| Gender:F | −0.06 | 0.043 | 0.945 | (0.87-1.03) |
| Income:High | −0.053 | 0.038 | 0.948 | (0.88-1.02) |
| Income:Low | −0.105** | 0.039 | 0.899 | (0.83-0.97) |
| Income:SG | −0.293** | 0.056 | 0.746 | (0.67-0.83) |
| Origins:Commuter | 0.156** | 0.034 | 1.169 | (1.09-1.25) |
| Origins:Offsite | −0.128* | 0.069 | 0.880 | (0.77-1.01) |
| HighschoolType:Classical | 0.028 | 0.062 | 1.029 | (0.91-1.16) |
| HighschoolType:Foreigner | −0.337** | 0.157 | 0.713 | (0.52-0.97) |
| HighschoolType:Others | 0.144 | 0.096 | 1.155 | (0.96-1.39) |
| HighschoolType:Technical | −0.144** | 0.045 | 0.866 | (0.79-0.95) |
| HighschoolGrade | 0.007** | 0.001 | 1.007 | (1.00-1.01) |
| AdmissionScore | −0.003 | 0.002 | 0.997 | (0.99-1.00) |
| Age19: > 19 | −0.384** | 0.046 | 0.681 | (0.62-0.75) |
| ECTSPprog | −0.061** | 0.001 | 0.941 | (0.94-0.94) |
| GPAprog | −0.028** | 0.002 | 0.972 | (0.97-0.98) |
| Number of events | 4549 | | | |
| Observations | 197591 | | | |
| Frailty | $\hat{\theta} = 0.012$ | $se((\hat{\theta})) = 0.006$ | pval = 0.020 | |
| Concordance | 0.857 | | | |
| Log Likelihood | −14528.61 | | | |

Note: ∗$p < 0.1$; ∗∗$p < 0.05$. Estimated baseline survival and hazard functions are reported in Figure A4b in Appendix A4.

hazard functions are reported in Figure 3. Except for *DG*4, that confirmed to be associated with a higher dropout risk both in the time-invariant and time-dependent frameworks, the other departments with an effect significantly different from 1 differ from the ones identified in the time-invariant framework. Here, *DG*10 and *DG*12 are associated with higher dropout risks, while *DG*6 and *DG*14 are associated with lower ones, suggesting that, net to the effect of the entire student career in the first three years, there are heterogeneous dropout dynamics across these departments.

The C-index measured on the training and on the test sets are both equal to 0.857. As expected, the inclusion of the career tracks over time improves the model accuracy and the predictive power, leading to a powerful model. Nonetheless, in order to promptly help at-risk students, early predictions are needed. In this perspective, in the next subsection, we conduct a comparative analysis in order to estimate the best trade-off between accurate and early predictions.

## 5.3. Definition of an efficient early warning system

In order to evaluate the trade-off between early and accurate predictions, we perform a comparative analysis in which we build several shared frailty Cox models by including student information
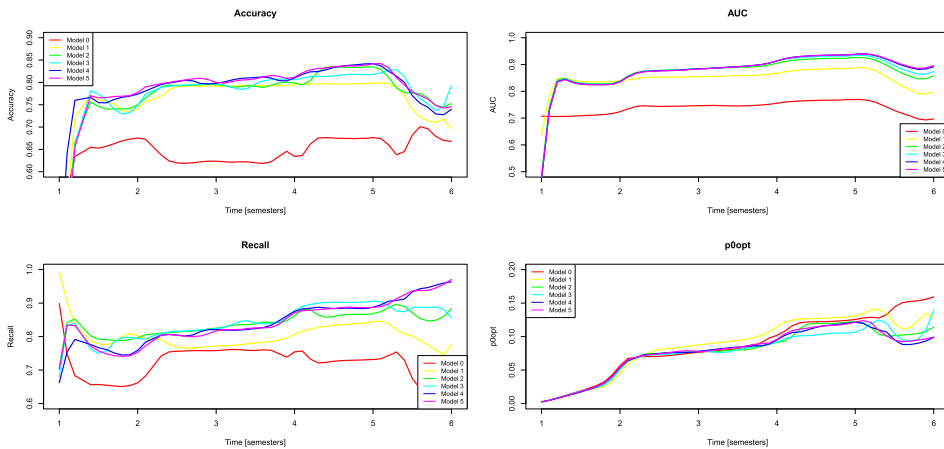


**Figure 3.** Estimated frailty terms and faculty-specific baseline hazard functions in the time-varying case. Note: Left panel shows the empirical Gamma Frailty terms for the 16 degree programs estimated by the shared frailty Cox model with time-varying covariates. Red and blue points identify the departments that have a frailty term significantly higher and lower than 1, respectively. Right panel reports the department-specific baseline hazard functions for the 16 specific degree programs.

measured until different time points and we evaluate their predictive performance, in terms of C-index and classification indices. In particular, we build six subsequent models, numbered from 0 to 5. *Model 0* only includes the student information measured at the time of enrolment (all time-invariant, no student career progress information is considered). *Model 1* includes the student information measured at the time of enrolment plus the number of ECTS obtained during the first semester (all time-invariant). Lastly, for *s = 2, … ,5, Model s* includes the student information measured at the time of enrolment plus the progress of the number of ECTS and GPA obtained during the first *s* semesters.

As we did in the previous section, we randomly divided the sample into training (70%) and test sets (30%). The predictive performance is measured in terms of C-index, accuracy, recall, and Area Under the ROC Curve (AUC), measured on the test set. For each of the six models, the classification indices are built at different time instants $t* = \{1.0, 1.1, 1.2, … , 5.9, 6.0\}$ by classifying a student as dropout or not standing on his/her predicted dropout probability at time $t*$. At each time $t*$, the optimal threshold $p_0(t*)$ for the classification is found (on the training set) and students in the test set are classified accordingly. Figure 4 shows the six trends of accuracy, recall, AUC, and optimal classification threshold in time, while Table 5 reports the C-index of the six models computed on the test set.

From Figure 4 and Table 5, we observe a first significant improvement in the models predictive performance when we move from Model 0 to Model 1 and a second less pronounced one when we move from Model 1 to Model 2. The difference in the performances between the last four models is almost negligible. This result suggests that student information at the time of enrolment is not



**Figure 4.** Accuracy, recall, AUC, and optimal classification threshold in time for the six Shared frailty Cox models, estimated on the test set. Note: Figures show Accuracy, Recall, AUC and values of optimal p for the predictions of dropout in different moments.

**Table 5.** Concordance Index computed on the test set, comparison between the 6 different time-dependent shared frailty Cox models.

| Model | C-index |
|---|---|
| *Model 0* | 0.682 |
| *Model 1* | 0.813 |
| *Model 2* | 0.849 |
| *Model 3* | 0.851 |
| *Model 4* | 0.855 |
| *Model 5* | 0.857 |

Note: The Table presents the Concordance Indexes at different moments of the students' academic careers (from first semester to fifth semester) to detect the optimal moment for predicting time-to-dropout.

sufficient to provide a good prediction for the dropout risk (C-index = 0.682, accuracy between 0.65 and 0.7, AUC between 0.7 and 0.75). With the inclusion of first-semester information, we become much more confident in identifying students at risk (C-index = 0.813, accuracy between 0.7 and 0.75, AUC approximately 0.8), and with the entire first-year information we reach a level that is comparable to the one that we obtain by observing the complete student career of the first six semesters.

This evidence, together with the high frequency of dropout during the first year, suggests that first-year career is already extremely informative and is enough to outline targeted interventions. The end of the first and second semesters represent two pivotal moments in implementing preventive actions.

## 6. Concluding remarks and policy implications

The need to deal with the dropout issue is particularly relevant for scholars and policy makers due to its important consequences at the personal, social, and economic levels (Castro-Lopez et al. 2022). Early Warning System is a promising approach aiming at reducing educational withdrawal, predicting the phenomenon as soon as possible. However, academic literature focuses much on identifying the 'who', while less is done about the 'when'. Indeed, the key research goals of this paper are identifying the time when dropout occurs and the optimal time to predict it.

To pursue these goals, we developed a set of shared frailty Cox models with time-invariant and time-varying covariates for predicting student dropout at different engineering faculties of PoliMi. The main innovation of this work relies on the methodological approach adopted and on its advantages: the time-to-event approach allows to predict of the time to dropout, while the frailty and the time-varying covariates allow to fit the data and their complexity. The first aspect is relevant since it represents clear insights for universities and program managers, who can effectively use these predictions to intervene on time. In our case, dropout mainly occurs at the end of every year, but particularly after the first one. This means that students face difficulties especially at the beginning of their career. Potential reasons could be found in the low pre-academic preparation or in a misalignment in students' expectations about university career. In this perspective, empowering the selection procedure and enriching the set of student information collected at the time of enrolment would help in providing more accurate and timely predictions. The second key takeaway relates to the characteristics of the most resilient (and, on the contrary, the most at risk) students. Girls, study grant recipients, and offsite students are those who retain more than their counterparts. The interpretation could be found in their (expected) higher motivation, that represents the main latent factor related to students' retention (Tinto 2017). STEM disciplines often suffer from a lack of female representation. Students who receive study grants may feel a sense of responsibility and duty due to the opportunity they have been given, which is supported by Modena, Rettore, and Tanzi (2020). Additionally, offsite students have typically relocated to another city, Milan, with the highest rental prices in Italy, likely due to the sacrifices made by their families.

In the case of Politecnico di Milano, there is evidence of a significant phenomenon of self-selection among students during the enrollment process. Despite the admission test being less challenging compared to other universities, Politecnico di Milano holds the top position in the QS ranking of Italian universities for 2023. However, the university's high dropout rate tends to discourage potential students from enrolling. This observation is further supported by the descriptive statistics of the sample, which indicate that the majority of students come from a scientific high school background, achieve high grades in their final high school exams, and come from non-disadvantaged socio-economic backgrounds.

The last consideration relates to the adoption of an Early Warning System for the detection of students at risk of dropout. This paper aims to the stage for a discussion about the timing of predictions as the result of an optimization problem to balance their accuracy and their timing. Findings indicate that as the student's career progresses, predictions' precision improves (as expected), but waiting for too long may lead the university to not have enough time to retain students. Evidence suggests that

a possible optimal moment for prediction is the end of the first year since the improvement in accuracy for the following semesters is nearly negligible. This approach enables to perform proactive interventions in a prioritized manner when limited academic resources are available. As in the case of Wayne State University (Ameri et al. 2016), for the student retention problem it is critical to not only correctly classify whether a student is going to dropout but also when this is going to happen. This approach is crucial for a focused intervention. Another interesting aspect to consider, as discussed by Gury (2011) for a national sample of French university students, is the distinction between early and late leavers. This dropout characterization allows policy managers to develop different strategies based on the category. Early dropouts would have benefited from receiving additional preexisting information about their likelihood of success and a greater emphasis on instilling the social and academic values essential for pursuing a higher education degree. The situation for late dropouts is different. Individuals who struggle in university are not provided opportunities for remediation, and the educational system fails to facilitate their pursuit of a more suitable academic program aligned with their abilities and interests.

Possible and interesting further development directions regard two main aspects. The former concerns the investigation of the heterogeneity at the degree program level. Indeed, the dropout dynamics across degree courses might differ across time (e.g. the baseline hazard function of a degree program might be higher during the first year but lower during the second, with respect to the average), and time-invariant frailties are not able to catch this source of variability. Developing Cox models with time-varying frailties and degree program-specific parameters of covariates would significantly help the research in this direction. The latter regards the possibility of enriching the student-level dataset by including information about student motivation, psychological and personal aspects that would help the prediction allowing for even earlier accurate estimates. To this end, the research group is collecting survey's information from dropout students about motivations and future perspectives about their career.

## Notes

1. The choice of this time span is motivated by the fact that the university administrative database remained unchanged over time. In this way, we have the possibility the larger amount possible of homogenized data. We are aware that in this time window, different events occurred, among all Covid-19. However, we are not interested in considered them in our analysis, since we are interested in predictions of time-to-dropout.
2. In the dataset, students who dropped out during their first semester are 1,700 (21.33% of the total dropout students). The University does not have time to take targeted preventive actions on these dropouts; therefore, their prediction is neither attractive nor valuable. Also, these dropouts have different characteristics compared to students leaving Politecnico di Milano afterwards, as it is represented in Appendix A1. For this reason, this population requires a different analysis, since the dynamics behind this phenomenon are different.
3. Students who want to leave the university, need to fill an online form for administrative reasons (e.g., to not pay the next fee), anytime during their academic career.
4. We do not include *GPA* at this stage because *ECTS* and *GPA* are highly correlated, due to all students that have 0 *GPA* and 0 *ECTS* at first semester.
5. The groups whose lower bound of the confidence interval is greater than 1 are red, while the groups whose higher bound of the confidence interval is lower than 1 blue. In grey we find the departments whose confidence interval contains 1, suggesting that they are not significantly different from the average.

## Disclosure statement

No potential conflict of interest was reported by the authors.

## ORCID

*Chiara Masci* 🆔 http://orcid.org/0000-0002-9208-3194
*Marta Cannistrà* 🆔 http://orcid.org/0000-0001-7631-7790
*Paola Mussida* 🆔 http://orcid.org/0000-0001-6364-5381

## References

Alban, M., and D. Mauricio. 2019. "Predicting University Dropout Trough Data Mining: A Systematic Literature." *Indian Journal of Science and Technology* 12 (4): 1–12. https://doi.org/10.17485/ijst/2019/v12i4/139729.

Ameri, S., M. J. Fard, R. B. Chinnam, and C. K. Reddy. 2016, October. "Survival Analysis Based Framework for Early Prediction of Student Dropouts." Proceedings of the 25th ACM international on conference on information and knowledge management, 903–912.

ANVUR. 2018. *Rapporto biennale sullo stato del sistema universitario e della ricerca*. https://www.anvur.it/wp-content/uploads/2018/11/ANVUR-Completo-con-Link.pdf.

Arulampalam, W., R. A. Naylor, and J. P. Smith. 2004. "A Hazard Model of the Probability of Medical School Drop-Out in the UK." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 167 (1): 157–178.

Barragaan, S., L. Gonza´lez, and G. Calder´on. 2022. "Modelling Student Dropout Risk Using Survival Analysis and Analytic Hierarchy Process for an Undergraduate Accounting Program." *Interchange* 53 (3-4): 407–427.

Booth, L. L., and S. E. Satchell. 1995. "The Hazards of Doing a PhD: An Analysis of Completion and Withdrawal Rates of British PhD Students in the 1980s." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 158 (2): 297–318.

Cannistrà, M., C. Masci, F. Ieva, T. Agasisti, and A. M. Paganoni. 2022. "Early-Predicting Dropout of University Students: An Application of Innovative Multilevel Machine Learning and Statistical Techniques." *Studies in Higher Education* 47 (9): 1935–1956. https://doi.org/10.1080/03075079.2021.2018415.

Castro-Lopez, A., A. Cervero, C. Galve-Gonza´lez, J. Puente, and A. B. Bernardo. 2022. "Evaluating Critical Success Factors in the Permanence in Higher Education Using Multi-Criteria Decision-Making." *Higher Education Research & Development* 41 (3): 628–646. https://doi.org/10.1080/07294360.2021.1877631.

Chen, C., G. Sonnert, P. M. Sadler, D. Sasselov, and C. Fredericks. 2020. "The Impact of Student Misconceptions on Student Persistence in a MOOC." *Journal of Research in Science Teaching* 57 (6): 879–910. https://doi.org/10.1002/tea.21616.

David, G. K., and K. Mitchel. 2012. *Survival Analysis: A Self-Learning Text*. New York: Spinger.

De Valero, Y. F. 2001. "Departmental Factors Affecting Time-to-Degree and Completion Rates of Doctoral Students at one Land-Grant Research Institution." *The Journal of Higher Education* 72 (3): 341–367. https://doi.org/10.2307/2649335.

Grove, W. A., D. H. Dutkowsky, and A. Grodner. 2007. "Survive Then Thrive: Determinants of Success in the Economics PH.D. Program." *Economic Inquiry* 45 (4): 864–871. https://doi.org/10.1111/j.1465-7295.2007.00041.x.

Gury, N. 2011. "Dropping Out of Higher Education in France: A Micro-Economic Approach Using Survival Analysis." *Education Economics* 19 (1): 51–64. https://doi.org/10.1080/09645290902796357.

Hegde, V., and P. P. Prageeth. 2018. "Higher Education Student Dropout Prediction and Analysis Through Educational Data Mining." 2018 2nd International conference on inventive systems and control (ICISC). (p. 694-699). https://doi.org/10.1109/ICISC.2018.8398887

Kaplan, E. L., and P. Meier. 1958. "Nonparametric Estimation from Incomplete Observations." *Journal of the American Statistical Association* 53 (282): 457–481. https://doi.org/10.1080/01621459.1958.10501452.

Kehm, B. M., M. R. Larsen, and H. B. Sommersel. 2019. "Student Dropout from Universities in Europe: A Review of Empirical Literature." *Hungarian Educational Research Journal* 9 (2): 147–164. https://doi.org/10.1556/063.9.2019.1.18.

Kleinbaum, D. G., and M. Klein. 1996. *Survival Analysis a Self-Learning Text*. Springer.

Lesik, S. A. 2007. "Do Developmental Mathematics Programs Have a Causal Impact on Student Retention? An Application of Discrete-Time Survival and Regression-Discontinuity Analysis." *Research in Higher Education* 48 (5): 583–608. https://doi.org/10.1007/s11162-006-9036-1.

Mantel, N. 1966. "Evaluation of Survival Data and Two New Rank Order Statistics Arising in its Consideration." *Cancer Chemotherapy Reports* 50: 163–170.

Min, Y., G. Zhang, R. A. Long, T. J. Anderson, and M. W. Ohland. 2011. "Nonparametric Survival Analysis of the Loss Rate of Undergraduate Engineering Students." *Journal of Engineering Education* 100 (2): 349–373. https://doi.org/10.1002/j.2168-9830.2011.tb00017.x.

Modena, F., E. Rettore, and G. M. Tanzi. 2020. "The Effect of Grants on University Dropout Rates: Evidence from the Italian Case." *Journal of Human Capital* 14 (3): 343–370. https://doi.org/10.1086/710220.

Mussida, P., and P. L. Lanzi. 2022. "A Computational Tool for Engineer Dropout Prediction. *IEEE Global Engineering Education Conference*.

No, F., K. Taniguchi, and Y. Hirakawa. 2016. "School Dropout at the Basic Education Level in Rural Cambodia: Identifying Its Causes Through Longitudinal Survival Analysis." *International Journal of Educational Development* 49: 215–224. https://doi.org/10.1016/j.ijedudev.2016.03.001.

Patacsil, F. F. 2020. "Survival Analysis Approach for Early Prediction of Student Dropout Using Enrollment Student Data and Ensemble Models." *Universal Journal of Educational Research* 8 (9): 4036–4047. https://doi.org/10.13189/ujer.2020.080929.

Plank, S. B., S. DeLuca, and A. Estacion. 2008. "High School Dropout and the Role of Career and Technical Education: A Survival Analysis of Surviving High School." *Sociology of Education* 81 (4): 345–370. https://doi.org/10.1177/003804070808100402.

Rondeau, V., J. Gonzalez, Y. Mazroui, A. Mauguen, A. Diakite, A. Laurent, C. Sofeu. 2019. Frailty pack: General frailty models: Shared, Joint and Nested Frailty Models With Prediction; Evaluation of Failure-Time Surrogate Endpoints. Rpackage version 3.0.3. https://cran.r-project.org/package = frailtypack.

Seidel, E., and S. Kutieleh. 2017. "Using Predictive Analytics to Target and Improve First Year Student Attrition." *Australian Journal of Education* 61 (2): 200–218. https://doi.org/10.1177/0004944117712310.

Singer, J. D., and J. B. Willett. 1993. "It's About Time: Using Discrete-Time Survival Analysis to Study Duration and the Timing of Events." *Journal of Educational Statistics* 18 (2): 155–195.

Soares, T. M., N. d. S. Fernandes, M. C. N´obrega, and A. C. Nicolella. 2015. "Fatores Associados ao Abandono Escolar no Ensino Médio Público de Minas Gerais." *Educação e Pesquisa* 41: 757–772. https://doi.org/10.1590/S1517-9702201507138589.

Spitzer, M. W. H., R. Gutsfeld, M. Wirzberger, and K. Moeller. 2021. "Evaluating Students' Engagement with an Online Learning Environment During and After COVID-19 Related School Closures: A Survival Analysis Approach." *Trends in Neuroscience and Education* 25: 100168. https://doi.org/10.1016/j.tine.2021.100168.

Steck, H., B. Krishnapuram, C. Dehing-oberije, P. Lambin, and V. C. Raykar. 2007. "On Ranking in Survival Analysis: Bounds on the Concordance Index." In *Advances in Neural Information Processing Systems (Vol. 20)*, edited by J. Platt, D. Koller, Y. Singer, and S. Roweis. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2007/file/

Thaithanan, J., O. Thinnukool, M. Chaichana, and W. Wanishsakpong. 2021. "Using Survival Analysis to Investigate Undergraduate Student Dropout Rates in the College of Arts, Media and Technology, Chiang Mai University." *Multicultural Education* 7 (10).

Tinto, V. 2017. "Through the Eyes of Students." *Journal of College Student Retention: Research, Theory & Practice* 19 (3): 254–269. https://doi.org/10.1177/1521025115621917.

Utami, S., I. Winarni, S. K. Handayani, and F. R. Zuhairi. 2020. "When and Who Dropouts from Distance Education?" *Turkish Online Journal of Distance Education* 21 (2): 141–152. https://doi.org/10.17718/tojde.728142.

Vallejos, C. A., and M. F. Steel. 2017. "Bayesian Survival Modelling of University Outcomes." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 180 (2): 613–631.

Van Der Haert, M., E. Arias Ortiz, P. Emplit, V. Halloin, and C. Dehon. 2014. "Are Dropout and Degree Completion in Doctoral Study Significantly Dependent on Type of Financial Support and Field of Research?" *Studies in Higher Education* 39 (10): 1885–1909. https://doi.org/10.1080/03075079.2013.806458.

Weybright, E. H., L. L. Caldwell, H. Xie, L. Wegner, and E. A. Smith. 2017. "Predicting Secondary School Dropout among South African Adolescents: A Survival Analysis Approach." *South African Journal of Education* 37 (2): 1–11. https://doi.org/10.15700/saje.v37n2a1353.

Xie, Z. 2020. "Modelling the Dropout Patterns of MOOC Learners." *Tsinghua Science and Technology* 25 (3): 313–324. https://doi.org/10.26599/TST.2019.9010011.

# Appendices

## *Appendix A1. Comparison between 1st semester dropout and other dropouts*

**Table A1.** Dropout descriptive statistics and their comparison between 1st semester dropouts and later dropouts.

|  | Later dropout (N = 6537) | 1st semester dropout (N = 1700) | Total dropout (N = 8237) | *p* value |
|---|---|---|---|---|
| Gender |  |  |  | < 0.001 |
| F | 1098 (16.8%) | 589 (34.6%) | 1687 (20.5%) |  |
| M | 5439 (83.2%) | 1111 (65.4%) | 6550 (79.5%) |  |
| Admission age |  |  |  | < 0.001 |
| Mean (SD) | 0.129 (0.335) | 0.069 (0.253) | 0.117 (0.321) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| Income |  |  |  | < 0.001 |
| High | 2186 (33.4%) | 1457 (85.7%) | 3643 (44.2%) |  |
| Low | 1883 (28.8%) | 120 (7.1%) | 2003 (24.3%) |  |
| Medium | 1823 (27.9%) | 92 (5.4%) | 1915 (23.2%) |  |
| Study Grant | 645 (9.9%) | 31 (1.8%) | 676 (8.2%) |  |
| Origins |  |  |  | 0.001 |
| Commuter | 4267 (65.3%) | 1178 (69.3%) | 5445 (66.1%) |  |
| Milanese | 1847 (28.3%) | 444 (26.1%) | 2291 (27.8%) |  |
| Offsite | 423 (6.5%) | 78 (4.6%) | 501 (6.1%) |  |
| Highschool type |  |  |  | < 0.001 |
| Classica | 417 (6.4%) | 140 (8.2%) | 557 (6.8%) |  |
| Others | 160 (2.4%) | 44 (2.6%) | 204 (2.5%) |  |
| Scientifica | 4909 (75.1%) | 1301 (76.5%) | 6210 (75.4%) |  |
| Straniera | 74 (1.1%) | 9 (0.5%) | 83 (1.0%) |  |
| Tecnica | 977 (14.9%) | 206 (12.1%) | 1183 (14.4%) |  |
| Highschool grade |  |  |  | < 0.001 |
| Mean (SD) | 78.362 (11.427) | 86.150 (11.706) | 79.969 (11.909) |  |
| Range | 60.000 - 100.000 | 60.000 - 100.000 | 60.000 - 100.000 |  |
| Admission score |  |  |  | < 0.001 |
| Mean (SD) | 70.441 (8.103) | 71.969 (8.894) | 70.757 (8.295) |  |
| Range | 60.000 - 100.000 | 60.030 - 100.000 | 60.000 - 100.000 |  |

Note: The Table presents proportions across levels for categorical variables and mean, standard deviation and range for numerical variables. The *p*-value refers to the chi-sq test for categorical variables and to the anova tests for numerical variables.

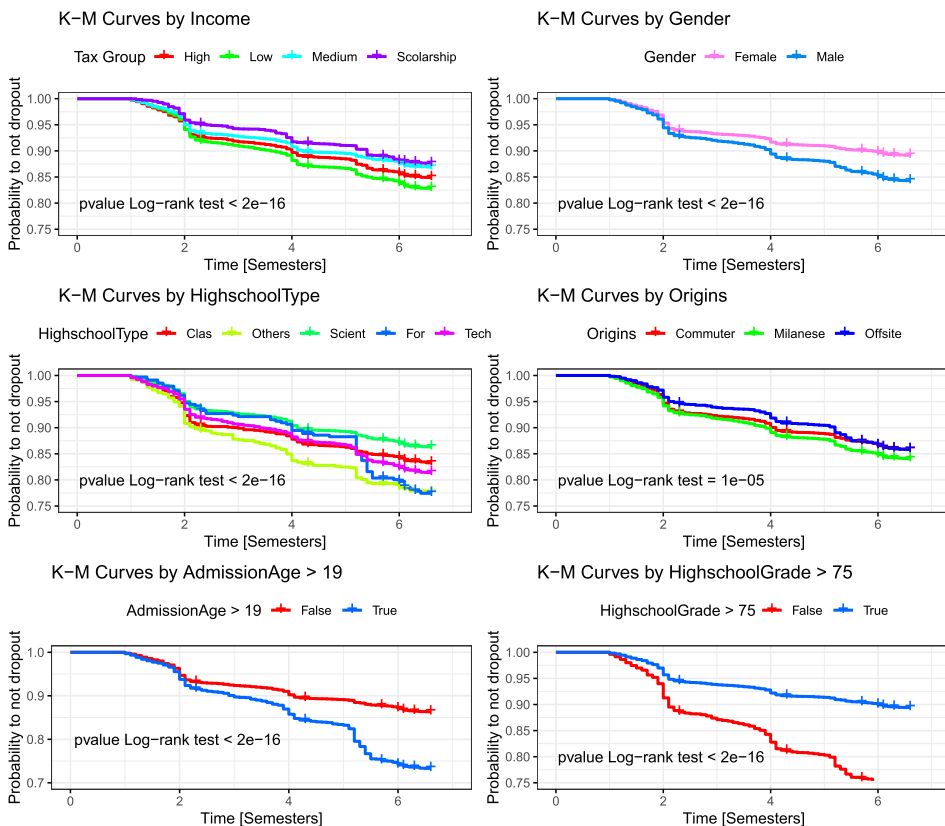## *Appendix A2. Students' distribution within programs*

**Table A2.** Distribution of students within the 16 programs and relative dropout percentage.

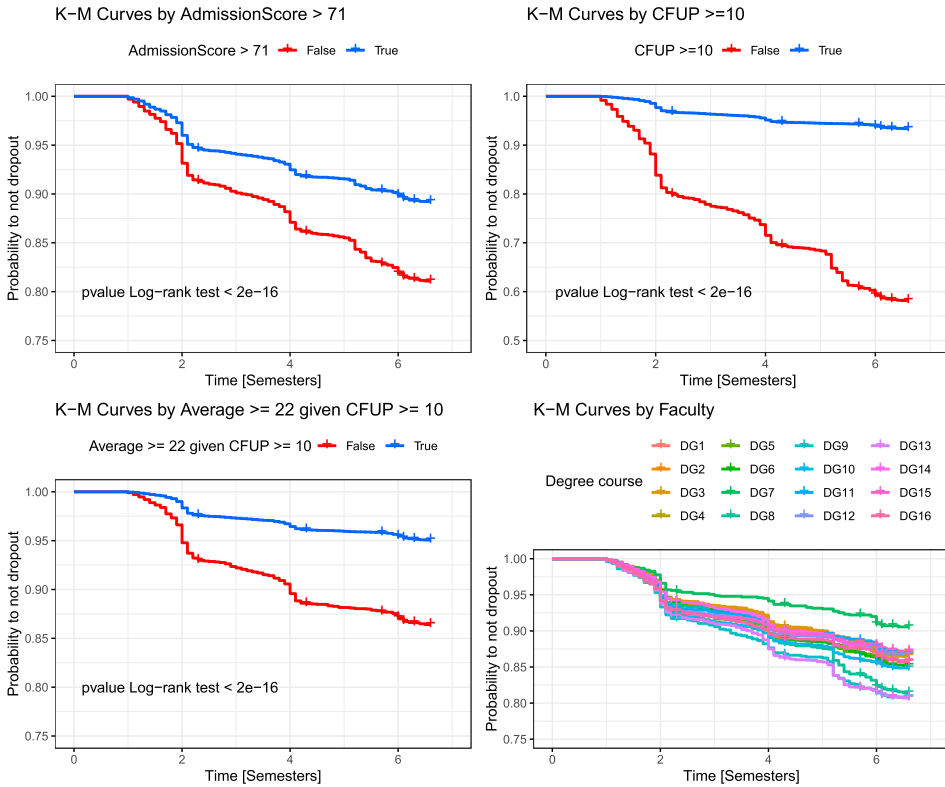| Degree Program Code | Number of students | % dropout |
|---|---|---|
| DG1 | 4,392 | 11.04 |
| DG2 | 1,208 | 11.04 |
| DG3 | 2,265 | 11.03 |
| DG4 | 4,374 | 12.05 |
| DG5 | 2,112 | 12.64 |
| DG6 | 1,767 | 13.02 |
| DG7 | 1,207 | 7.71 |
| DG8 | 1,586 | 15.38 |
| DG9 | 1,002 | 16.67 |
| DG10 | 3,948 | 13.52 |
| DG11 | 1,635 | 10.95 |
| DG12 | 6,659 | 11.61 |
| DG13 | 6,719 | 16.55 |
| DG14 | 6,020 | 12.09 |
| DG15 | 2,486 | 11.06 |
| DG16 | 2,121 | 12.21 |

## Appendix A3. Dropout risk across time by students' characteristics

From KM curves in Figure A3a, we observe that, despite the number of males is widely larger than the number of females (77.5% vs 22.5%), males are more likely to drop out. Regarding the family income, students with administrative support (*SG* category) are less likely to drop out across time. This could depend from the fact that students with SG are more motivated and feel the responsibility for having obtained a grant. The highest risk category, especially right after the end of first semester, is that of high income students. Nonetheless, on the long term, also low income students show a higher dropout probability with respect to the other categories, which suggests that students with a more disadvantaged background, who do not receive administrative support, are more exposed to dropout, especially on the long term. The dropout probability of the Medium category reaches results very close to the SG group at the end of the follow up time. For what concerns student origins, Offsite students (i.e. students coming from other regions who moved to Milan to study at PoliMi) have on average a lower dropout risk with respect to Milanese and Commuter students. Regarding the type of high school attended before the enrolment at PoliMi, most of the students come from Scientific school (80.5%) and result to be the ones with lowest dropout risk. Students coming from Classical schools present a significant higher risk of dropout during the first year, while, on the opposite, students who attended a high school abroad are less likely to dropout at the beginning but more likely to dropout during their third year. At the end of the follow-up, technical schools and all other types of high schools show a relatively high dropout probability. Furthermore, also the high school grade results to be a determinant of the dropout risk. We identify 75 as the threshold that differentiates the most the two populations, highlighting that students with a high school final mark lower than 75 have on average higher risk of dropout with respect to the others. Students enrolling at PoliMi later than the standard age (19), tent to drop more than younger students, especially after the first year.

Focusing on the early performance at PoliMi (KM curves in Figure A3b), we observe a lower dropout risk for those students obtaining an admission score higher than 71, that resulted to be the most significant threshold. In terms of ECTS and GPA measured at the end of first semester, we observe that obtaining less than 10 ECTS is an extremely predictive risk factor. The sharp difference between the two KM curves highlights the importance of this information and its predictive power. Among the students who obtained at least 10 ECTS, having a GPA lower than 22, i.e. the most
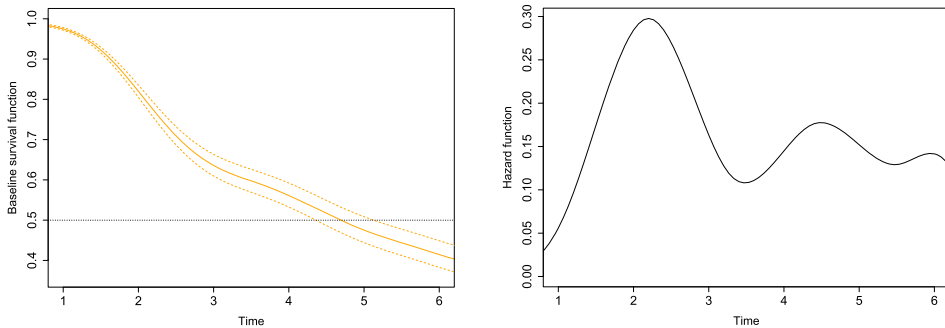


**Figure A3a.** Kaplan-Meier Curves for Gender Income, Origins, HighschoolType, High schoolGrade, and AdmissionAge.
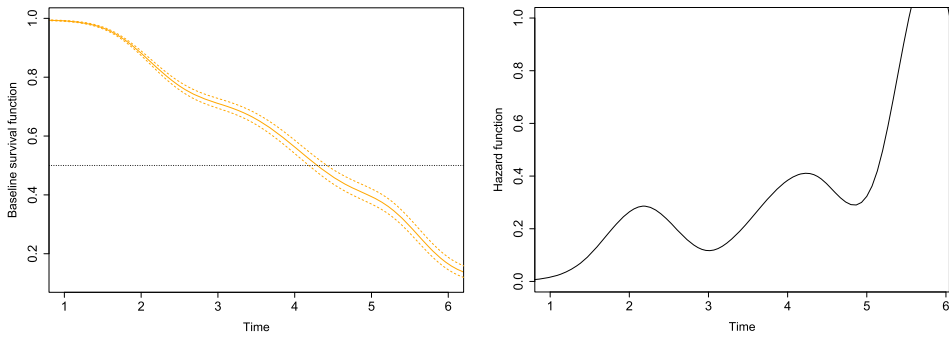
### K–M Curves by AdmissionScore > 71

### K–M Curves by CFUP >=10

### K–M Curves by Average >= 22 given CFUP >= 10

### K–M Curves by Faculty

**Figure A3b.** Kaplan-Meier Curves for *AdmissionScore*, *ECTS*, *GPA,* and *Degree Program*. Note: For each numerical covariate, the threshold represents the value for which the difference between the two Kaplan-Meyer curves is maximized.

discriminant value, constitutes a further risk factor. Lastly, given our interest in investigating the difference across degree programs, we observe that the 16 KM curves show heterogeneous dropout dynamics across faculties, detecting up to a 13% difference in the dropout percentage at the end of follow-up across faculties.

## Appendix A4. Baseline survival and hazard functions of the shared frailty Cox models

**Figure A4a.** Estimated baseline survival and hazard functions of the shared frailty Cox model with time-invariant covariates.

**Figure A4b.** Estimated baseline survival and hazard functions of the shared frailty Cox model with time-varying covariates. Note: At the end of the follow-up, the baseline survival function reaches very low value (and in parallel, the hazard reaches very high ones). The number of progressive ECTS mainly drives this trend since surviving until the end of the third year with 0 ECTS is very unlikely.