

Supervised Machine Learning Predictive Analytics for Alumni Income

Daniela A. Gomez-Cravioto

Tecnologico de Monterrey: Instituto Tecnologico y de Estudios Superiores de Monterrey
<https://orcid.org/0000-0001-9286-9480>

Ramon E. Diaz-Ramos (✉ a01133921@itesm.mx)

Tecnologico de Monterrey: Instituto Tecnologico y de Estudios Superiores de Monterrey
<https://orcid.org/0000-0001-7324-7205>

Neil Hernandez Gress

Tecnologico de Monterrey: Instituto Tecnologico y de Estudios Superiores de Monterrey

Jose Luis Preciado

Tecnologico de Monterrey: Instituto Tecnologico y de Estudios Superiores de Monterrey

Hector G. Ceballos

Tecnologico de Monterrey: Instituto Tecnologico y de Estudios Superiores de Monterrey

Research

Keywords: Machine Learning, Income Prediction, Alumni Survey Analysis, Knowledge Discovery, Explainable Artificial Intelligence.

Posted Date: October 5th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-880103/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Supervised Machine Learning Predictive Analytics for Alumni Income

Daniela A. Gomez-Cravioto · Ramon E. Diaz-Ramos* · Neil Hernandez Gress · Jose Luis Preciado · Hector G. Ceballos ·

Received: date / Accepted: date

Abstract Background: This paper explores different machine learning algorithms and approaches for predicting alum income to obtain insights on the strongest predictors for income and a ‘high’ earners’ class. **Methods:** The study examines the alum sample data obtained from a survey from Tecnológico de Monterrey, a multicampus Mexican private university, and analyses it within the cross-industry standard process for data mining. Survey results include 17,898 and 12,275 observations before and after cleaning and pre-processing, respectively. The dataset includes values for income and a large set of independent variables, including demographic and occupational attributes of the former stu-

dents and academic attributes from the institution’s history. We conduct an in-depth analysis to determine whether the accuracy of traditional algorithms in econometric research to predict income can be improved with a data science approach. Furthermore, we present insights on patterns obtained using explainable artificial intelligence techniques. **Results:** Results show that the gradient boosting model outperformed the parametric models, linear and logistic regression, in predicting alum’s current income with statistically significant results ($p < 0.05$) in three tasks: ordinary least-squares regression, multi-class classification and binary classification. Moreover, the linear and logistic regression models were found to be the most accurate methods for predicting the alum’s first income. The non-parametric models showed no significant improvements. **Conclusion:** We identified that age, gender, working hours per week, first income after graduation and variables related to the alum’s job position and firm contributed to explaining their income. Findings indicated a gender wage gap, suggesting that further work is needed to enable equality.

Keywords Machine Learning, Income Prediction, Alumni Survey Analysis, Knowledge Discovery, Explainable Artificial Intelligence.

1 Introduction

Higher education institutions seek to boost their alumni outcomes after graduation. To validate whether this goal is being accomplished, there is value in collecting data from their alumni and identifying patterns between those who achieved their expected outcomes and those who did not. The results from this analysis can help guide stakeholder’s decision to support future

D. Gomez-Cravioto
School of Engineering and Sciences, Tecnológico de Monterrey, 64849 Monterrey, N.L.
Tel.: +52 81 8358 2000
E-mail: a01181520@itesm.mx

R. Diaz-Ramos
School of Engineering and Sciences, Tecnológico de Monterrey, 64849 Monterrey, N.L.
Tel.: +52 81 8358 2000
E-mail: a01133921@itesm.mx

N. Hernandez
School of Engineering and Sciences, Tecnológico de Monterrey, 64849 Monterrey, N.L.
Tel.: +52 81 8358 2000
E-mail: ngress@tec.mx

J. Preciado
School of Engineering and Sciences, Tecnológico de Monterrey, 64849 Monterrey, N.L.
Tel.: +52 81 8358 2000
E-mail: jlpreciadoarreola@tec.mx

H. Ceballos
School of Engineering and Sciences, Tecnológico de Monterrey, 64849 Monterrey, N.L.
Tel.: +52 81 8358 2000
E-mail: ceballos@tec.mx

1 alumni. In this study, we exploit the data from Tec-
2 nologico de Monterrey’s alumni survey to obtain two
3 main insights: assessing students’ economic outcomes
4 and validating if there are gaps with relation to diverse
5 backgrounds (e.g. gender, education and socioeconomic
6 diversity). Understanding the factors that may favour
7 some alumnus will help give them equal opportunities
8 to achieve their economic objectives.

9
10 Many institutions survey their graduates to collect
11 information on their post-graduation outcomes, such
12 as their income and socioeconomic status [1–3]. These
13 studies are beneficial not only to evaluate an institu-
14 tion’s effectiveness but also to support institutional plan-
15 ning and future students’ achievements. Unfortunately,
16 the actions taken to analyse the results rarely include
17 data mining to obtain insights regarding features that
18 can have a higher relationship with the outcome. This
19 is especially true for actionable features that can be
20 boosted with activities performed during students’ lives
21 on campus. In this work, a data-based model is built
22 for understanding the main factors that can influence
23 alumni income prediction. The study uses data science,
24 advanced analytics and machine learning techniques.
25 While career success can be evaluated as intrinsic or
26 extrinsic [4], this study will focus solely on extrinsic suc-
27 cess; based explicitly on the objective rating of salary.

28
29 The data on which this work is focused comes from
30 a survey carried out in 2018 by Tecnologico de Monter-
31 rey in the university’s approach to measure their grad-
32 uates’ social and economic impact. The survey was sent
33 through email to the total alumni population who grad-
34 uated from 1953 and 2017, and advertisement for this
35 survey was promoted on social media. The overall re-
36 sponse rate was 7% of the total population, which ac-
37 counts for 17,896 former students. The obtained data
38 set provides an excellent opportunity to supplement
39 the university with knowledge about previously hidden
40 trends and patterns regarding the factors that affect
41 alumni salary attainment. This study’s primary pur-
42 pose is to identify if factors such as age, gender, major,
43 graduate studies, the overall grade achieved, and par-
44 ent’s education and occupation can influence the alum-
45 nus’s first income after graduation and their current
46 monthly salary. Furthermore, we aim to identify the
47 variables that also impact the first income after gradu-
48 ation, which have resulted in a significant predictor for
49 the former.

50
51 The contribution of this study is threefold. Firstly,
52 it contributes to the modern field of machine learning
53 research applied to econometric studies by exploring in-
54 come distribution and comparing traditional economet-
55 ric techniques, such as Quantile Regression, Linear Re-
56 gression and Logistic Regression with machine learning

non-parametric tree-based algorithms, Random Forest
and Gradient Boosting to find the best method for ap-
proaching the problem of income prediction. Secondly,
the study adds to the existing literature in Educational
Data Analytics with a data-driven approach and ma-
chine learning algorithms applied to an alumni impact
survey dataset. Finally, the study adds to the applica-
tion of Knowledge Discovery in Data and Explainable
Artificial Intelligence by identifying rule-based patterns
in the dataset, identifying feature importance with SHAP
values, and performing a sensitivity analysis on the vari-
ables detected as having the most important relation-
ships with income.

1.1 Related Work

Over the past several decades, many studies have esti-
mated how the final grades, college major, demograph-
ics and occupation characteristics affect individuals’ in-
come. However, very few studies have combined all these
characteristics in a single model. This research builds
on previous works that examined college students’ fu-
ture income to determine the most important features
and use machine learning as a tool to assess these fea-
tures. A table showing the most recent studies on indi-
vidual income prediction with a multivariate model can
be seen in Table 1.

Alina Lazar [5] proposed the use of Support Vec-
tor Machines to predict income. She used the Current
Population Survey (CSP) from the U.S. Census Bureau
as a database for her study. This dataset contained so-
cial, demographic and economic characteristics of U.S.
citizens 16 years and older. The author used Principal
Component Analysis (PCA) to reduce the number of
features in the dataset and then fed this to a Support
Vector Machine classifier. With this, she was able to
achieve an accuracy score as high as 84%.

The study from Hartog and Webbink [6] analysed
both expectation and realisation of incomes from for-
mer students who graduated from high schools or uni-
versities in the Netherlands. The variables analysed in-
cluded background variables (Gender, age, parent’s ed-
ucation, parent’s income), higher education variables
(year of education, student’s status), and secondary ed-
ucation variables (school marks), potential work expe-
rience (time since graduation). One of the experiments
conducted in this study included a prediction of realised
earnings. This model was performed with OLS regres-
sion and achieved a 16% R².

The study from Lee and Lee [7] investigated the
wage determinants in the Korean labour market. The
researchers used quantile regression methods. They in-
dicated that the advantage of quantile regressions is

Table 1: Recent studies on income prediction summary recollection

Source	Task	Methods	Results
Lazar [5]	classification	SVM	Acc = 0.84
Hartog, Webbink [6]	regression	OLS	R2 = 0.14
Lee, Lee [7]	quantile regression	5th	Pseudo-R2 = 0.29
		25th	Pseudo-R2 = 0.33
		50th	Pseudo-R2 = 0.34
		75th	Pseudo-R2 = 0.34
		95th	Pseudo-R2 = 0.32
Oehlein [8]	regression	OLS	R2 = 0.37
Stran, Truong [9]	regression	Lasso OLS	USD \$6,394.64 (RMSE)
Figueiredo, Fontainha [10]	quantile regression	10th	Pseudo-R2 = 0.27
		50th	Pseudo-R2 = 0.45
		90th	Pseudo-R2 = 0.50
Sharath et al. [11]	classification	NB	Acc = 0.48
		C4.5	Acc = 0.51
		Boosted	Acc = 0.53
		C4.5	
Khongchai, Songmuang [12]	multi-class classification	DT	Acc = 0.73
		SVM	Acc = 0.43
		MLP	Acc = 0.38
		KNN	Acc = 0.84
		NB	Acc = 0.43
Chen, Sun, Thakuriah [13]	multi-class classification	SVM	Acc = 0.74
		DT	Acc = 0.74
		LR	Acc = 0.72
		RF	Acc = 0.71
		GBM	Acc = 0.70
		NN	Acc = 0.68
		LSTM	Acc = 0.65
		DNN	Acc = 0.65

that it allows examining a more comprehensive picture for different quantile wage groups. The results obtained from their study showed that age is the most important factor for wage determination. The authors also found that female workers are significantly underpaid compared to their male counterparts.

Oehlein [8] attempted to determine the aspects of college that impacted students' future income. He focused on deciding whether or not their GPA was an influencer. In this study, OLS regression was used with the R-squared obtained for the prediction being 0.374. The author's findings include that grades, natural ability and major significantly affect income. He found that the highest paying major was engineering and that the attribute female was negatively correlated with income.

The research study from Stran and Truong [9] evaluated different demographic features to predict earnings by comparing the results of students graduating from several colleges. The most important features identified in this study were the percentage of students who received a Pell grant, the number of female students, the rate of first-generation students, and the percentage of students who had sent a FAFSA application to multiple schools before entering. The best performance, considering the MSE, was the one from Lasso Regression and Random Forest.

The research performed by Figueiredo and Fontainha [10] studied the distinct wages for men and women in Portugal with an OLS and a quantile regression approach. The results from this study showed that quantile regression obtained better results than OLS. The findings indicated that the levels of education have a higher impact on wage determination. Also, the variables that contributed the most in the model were related to the firm, while those related to family only contributed to explaining men's wage. Finally, the study indicated a significant difference between men's and women's wages, indicating that further studies are required to explain the gender wage gap.

Sharath et al. [11] performed a machine learning study with the US Census Bureau dataset. The focus of the study was to obtain insights into the financial status of the people in US. The results obtained showed inequality in the society due to a gender wage gap and showcased one of the root causes for these inequalities by determining the relationship level between income and education level. Furthermore, a classification model was obtained to predict economic class categories with an accuracy of 53%.

Khongchai and Songmuang [12] presented their work using classification to predict future students' income. Their initial dataset contained 108 attributes obtained from graduate student history data collected for ten years from a university in Thailand. The features included gender, faculty-student ratio, programme, workplace type, work experience, certifications, total Grade Point Average (GPA), and salary. The best model obtained by the authors was the KNN with an 84% accuracy.

During this research, the most recent work was the one from Chen, Sun, and Thakuriah in 2019 [13]. The authors used many metadata in the web and relational attributes such as job descriptions, locations, job content and job-related features to predict individuals' salaries. They compared Support Vector Machines, Decision Trees, Logistic regression, Random Forest, Gradient Boosting, Graph Convolutional Networks and Deep Learning to classify six predefined categories (0-20,000; 20,000-30,000;

30,000–40,000; 40,000–50,000; and >50,000). The best overall accuracy obtained was the one from SVM with a 0.74 accuracy. The metadata features had a significant contribution to reaching this accuracy.

1.1.1 Dataset

This study analyses data from an alumni impact study survey conducted by Tecnológico de Monterrey university. The survey invitation was electronically sent to all former students since the inception of the university in 1943 (269,482 individuals). From the total population of alumni who graduated between 1953 and 2017, 7% responded to the survey; this accounts for 17,896 graduates across different generations. Tecnológico de Monterrey provided the original dataset collected from the survey for this study. The dataset contains no personally identifiable information, and the dataset contains all the salary figures normalised and reported in Mexican Pesos.

The records include 72 columns with demographic information from the alumni such as major, gender, graduation date, campus, age, occupation, level of education attained, parents' education, parents' occupation, as well as information related to their accomplishments such as businesses created, type of business, salary and score reported based on their satisfaction in their professional lives, as well as other variables.

In regards to the unintended bias inherent in this dataset, the first bias we identified was the one towards younger adults, specifically for those between 27 and 60 years old. Since the survey was sent by email, it accounted for fewer elder respondents (older than 60 years old). Hence, it is important to note that our results will not account for alumni older than 60.

2 Methodology

The methodology followed in this study is an adaptation of the Cross-Industry Standard Process for Data Mining (CRISP-DM). The CRISP-DM consists of a general model of a data mining project. This was developed in 1996 [14, 15] and has been widely used since then. The steps followed in this project to transform raw data into insights are shown in Fig.1. This diagram shows specific actions performed during the application of the CRISP-DM process.

2.1 Methods

Quantile Regression (QR): can be used when asymmetries and heavy tails exist in data distributions. The

advantage over linear regression is that this method is more robust to outliers and more flexible to the linear assumptions. The main difference between these two is that while least-squares regression is focused on minimising the sums of squared residuals to estimate models for conditional mean functions, QR models the conditional quantile of the response variable for some quantity of

$$\tau \in (0, 1)$$

, where $\tau = 0.5$ is the median [16]. For example, when trying to predict income in countries where the income is highly skewed, we can predict the median or the quantile instead of the mean. For this reason, the QR method is highly used in econometrics studies for wage determinants, discrimination effects and income inequality trends.

Ensemble Methods: to counteract the decision tree issues of stability [17] and accuracy [18], successful approaches include the ensemble of decision trees. The ensemble approach integrates multiple predictors and is built by two specific methods: bagging and boosting [19]. One of the best performing applications of the bagging method is Random Forest (RF), and a practical algorithm based on the boosting notion is the widely used ensemble method Gradient Boosting (GB).

Random Forest: algorithm consists of building B random samples, and for each of these samples, building a decision tree model f_b [20]. The final prediction is obtained by taking into account the vote of each of the models for a classification task 1 and the average prediction for a regression task 2.

$$\hat{C}(x) = \text{majority vote}\{\hat{C}_b\} \quad (1)$$

$$\hat{f}(x) = \frac{1}{B} \sum_{b=1}^B f_{b(x)} \quad (2)$$

The advantages of RF are mainly inherited from the decision trees, previously explained. For instance, they can be used for both classification and regression tasks; their nature enables them to handle categorical predictors; they are non-parametric models, so they do not need a formal distribution assumption. Additionally, they can manage non-linear relationships between the covariates and target variable, and they can perform feature selection automatically. However, unlike decision trees, random forests are harder to interpret, as the model is built with multiple decision trees, making it hard to visualise in a plot. Another limitation of this model is that it can become highly computationally complex when having a large number of trees.

Gradient Boosting combines multiple simple decision trees. The trees are joined sequentially, each tree

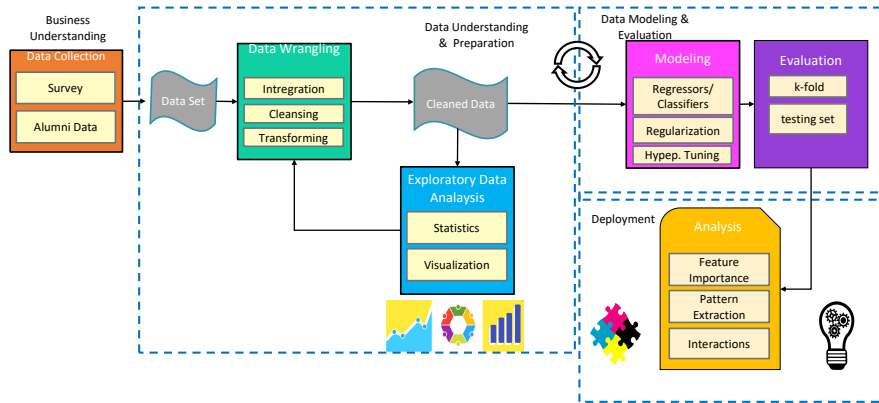


Fig. 1: The Knowledge Discovery Process flow diagram employed in this work and linked to the CRISP-DM methodology

trying to amend the errors of the previous one, $f_i^{(j)}$ (3). Frequently, this method has a better performance than Random Forest while having similar properties; however, careful tuning is required to avoid over-fitting the data [21].

$$f_i = f_i + \alpha f_i^{(j)} \quad (3)$$

2.2 Data Preparation

The data present in the survey results contain several missing values and an excessive number of attributes. To give the data the proper format for data mining, we performed a series of steps to clean the data.

- Data Integration:** First, the original dataset was compiled with a dataset with data from the university's planning department, which contained students' information upon graduation. The data included final GPA, number of semesters in which the student was involved in co-curricular activities (sports, leadership, and cultural activity), their English score, and whether they had previous work experience before graduating (internships). In this step, we noted that most campuses track recent students' participation in co-curricular activities and store their scores in a database; however, not all campuses held this information for alumni from older cohorts; this was among the fields with the largest percentage of missing values.
- Correcting Inconsistencies:** Since inconsistencies were present in the survey, subject matter expertise was needed to correct errors/inconsistencies. We

first cleaned the data by translating all the questions to variable names and translating all the data to English. There were many different words in the responses which referred to the same term, so we grouped them in a single word. Typos and misspelling were corrected. Finally, punctuation such as commas, apostrophes, quotes, question marks and others was removed.

- Handling Missing Values:** The next step performed was the handling of missing values. We eliminated all the records which had no information regarding the target variable. Then, we eliminated follow-up questions with more than 80% of missing values, as they were not of central importance to our analysis. Variables for extra-curricular activities and work experience previous to graduation had more than 80% of missing values. However, as these variables were of interest for our analysis, we decided to split the dataset into two. The first split was the information from all graduates; we used this later to predict their current income. The second one was a subset of the original dataset; we preserved all records with information regarding their school activities (co-curricular activities, internships, etc.); we used this later to predict the First Income after graduation. The age of the respondents from the subset is exclusively between 21 and 28 years old at the time of the survey. Hence, this analysis was considered exclusive for recent graduates (alumni graduated between 2012 and 2017). After this step, we had less than 40% of missing values in both datasets and no values missing for the tar-

get variables. We then performed a missing values imputation. We completed the imputation by using a Nearest Neighbours (NN) imputation, considering three distinct neighbours. We selected a K-NN model for this process as it has proven to be a useful technique for predicting missing values; it has surpassed the efficacy of average or median imputation in previous research [22]. The value imputed considered the three distinct neighbors' weighted average closer in the distance by using the Euclidean distance metric.

We identified extreme values in both target variables, current income and first income after graduation. In this regard, a winsorization method was used to mitigate the effect of the extreme values. The difference between just trimming the data and winsorizing it is that the latter will retain the observations but changes the numeric outliers to fall on the edge of the distribution [23].

With the winsorisation, we bounded the data to the 0.05 and 0.95 percentiles. The resulting distribution of the target variables, Current Income and First Income after graduation, after the winsorisation process, are observed in Fig. 2 and 3. We can see that there are still some outliers in the distribution. However, as these observations are ascertained as genuine, they are not removed. The data should be transformed for its use in data analysis, specifically in parametric models, as these require that the distribution is symmetric. This change is performed in the transformation step in this section.

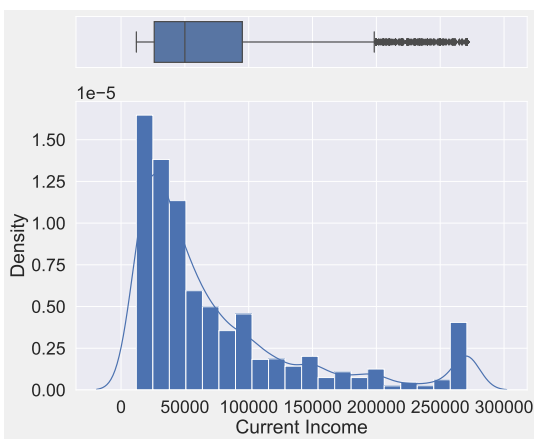


Fig. 2: Density bar plot and boxplot of the winsorized result of the Current Income variable

From the figures above, we can observe that both distributions are highly skewed to the right. To use these variables in the linear regression, we must per-

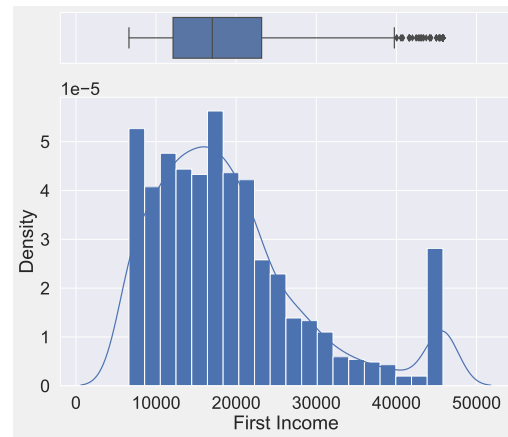


Fig. 3: Density bar plot and boxplot of the winsorized result of the First Income variable

form an additional transformation to the data. The transformation that we approached in this step was a box-cox transformation that could approximate normality assumptions. Even though the resulting transformations exhibits symmetry, they do not resemble a normal distribution. Both of the histograms show ‘heavy tails’, a common topic in income distributions [24–27]. Since having fat tails makes it of interest to understand the distribution, we decided to start modelling with a QR model (experiment A) instead of linear regression. The QR and the non-parametric machine learning models explored in this study make no assumptions about the distribution of the residuals; hence, they can be used when asymmetries and heavy tails exist in data distributions. The transformations that yielded the most symmetric distributions were then used for the linear regression model.

- Data Binning:** Since the unequal representation of the different groups could lead to unfair outcomes towards individuals or demographics, in this step, we seek to drop this difference by binning categories and reduce this imbalance as much as possible. In this sense, we grouped the predictor variables that contained more than five category labels. We performed this according to business acumen. An example of this was the variable ‘Campus’. Initially, the variable contained 33 categories, we reduced these to only six based on the economic regions in which the ‘Campus’ are located across Mexico and an additional one for the ‘Virtual Campus’, which represent those students that did not have a physical Campus but took all their courses online. The categories for the variables ‘Current Location’ and ‘Pre-Study Location’ were also binned based on these economic regions as well as an additional label for those living

overseas. Subsequently, these variables were binned based on their frequency.

5. **Dealing with Multicollinearity:** To determine whether the independent variables were highly correlated, we used the Variance Inflation Factor (VIF) and chose the score five as a threshold. VIF measures the correlation inflation between the independent variables. When the score was above the threshold, we dropped the variable from the dataset. The final results indicated that the remaining variables do not have a VIF above 5, which implies no multicollinearity issues in the dataset.
6. **Categorical Encoding:** Concerning variable encoding, for ordinal categorical variables, the assignment was done with incremental ordering, starting with the lowest category (i.e. 0 years was assigned 0, 1 to 3 years was set one, and more than three years was given 2). The categorical variables without a natural ordering, we transformed into dummy variables with one-hot-encoding. To deal with the ‘dummy variable trap’ [28] we dropped one of the dummy variables from each categorical feature.
7. **Data Standardization:** A standardisation of the data was performed in numerical variables using the standard normal or z-score normalisation method. This process is necessary so that the machine learning models treat all variables equally, and a variable is not considered more important because it has a higher range of values [29].

2.3 Exploratory Data Analysis

After the data wrangling step, the ‘Current Income’ dataset contains 12,275 observations and 65 variables, both continuous and categorical. Six of these variables are numeric, and 59 are categorical; however, all of them are now numeric values. On the other hand, the ‘First Income’ dataset contains 2264 observations and 39 variables; 2 of them numeric and 37 categorical. The first step for exploring these variables was to analyse correlation to identify those variables having a higher linear relation with the target variables. A heatmap showing the resulting Spearman coefficients for the continuous and ordinal variables is presented in Figure 4. The results showed a moderate relation ($0.3 \leq |r| \leq 0.5$) between the target variable and ‘Age’, ‘First Income’, ‘People in Charge’ and ‘Years Worked Foreign’. Then, statistics and visualisations were obtained to measure the marginal effect of variables of interest as per previous studies and their relation with the target variable. The money currency for the current and first income variables is in Mexican Pesos (MXN).

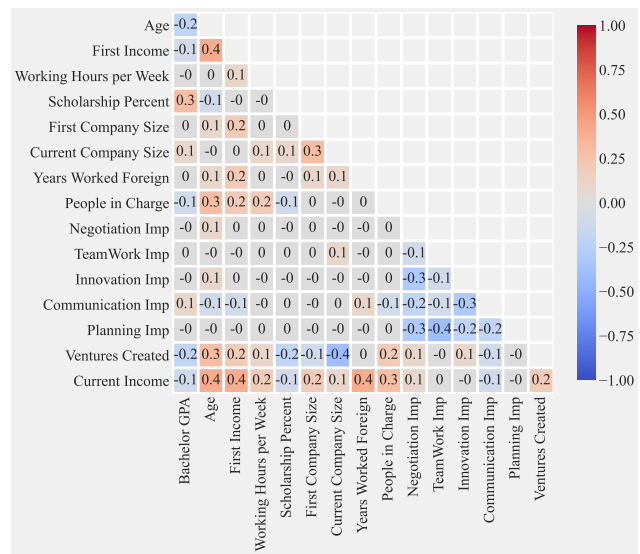


Fig. 4: Current Income Spearman correlation coefficient matrix of the continuous variables of the dataset

Salary Based on Gender The aggregation table in Table 2 depicts a comparison between Gender and ‘First Income’ and Gender and ‘Current Income.’ Overall, there is a gap between the results obtained for each gender; however, the gap is more prominent in the latter target variable. The gap for ‘First Income’ and ‘Current Income’ is \$3,151 and \$26,845 respectively. When looking closer into the results with the plot in Fig. 5 and after performing a Mann-Shitney-Wilcoxon test two-sided with Bonferroni correction hypothesis testing, we determined that the medians are significantly different, with a significance level of 0.01% in both analyses.

Table 2: First Income and Current Income median, mean, and standard deviation statistics by Gender

First Income			
Category	Mean	Std	Median
F	\$17,335.73	\$8,699.79	\$15,615.00
M	\$20,940.73	\$10,494.72	\$18,766.00
Current Income			
Category	Mean	Std	Median
F	\$55,933.13	\$54,427.00	\$37,155.00
M	\$89,726.13	\$73,106.86	\$64,000.00

Salary Based on School The salary variable was examined concerning the school groups. Table 3 shows that the alumni who graduated from ‘Engineering’ have a higher median than the other categories in both cases. The plot in Fig. 6 exhibits a significant difference between the School Variable medians. The difference between ‘Engineering’ and ‘Business’ is not very signifi-

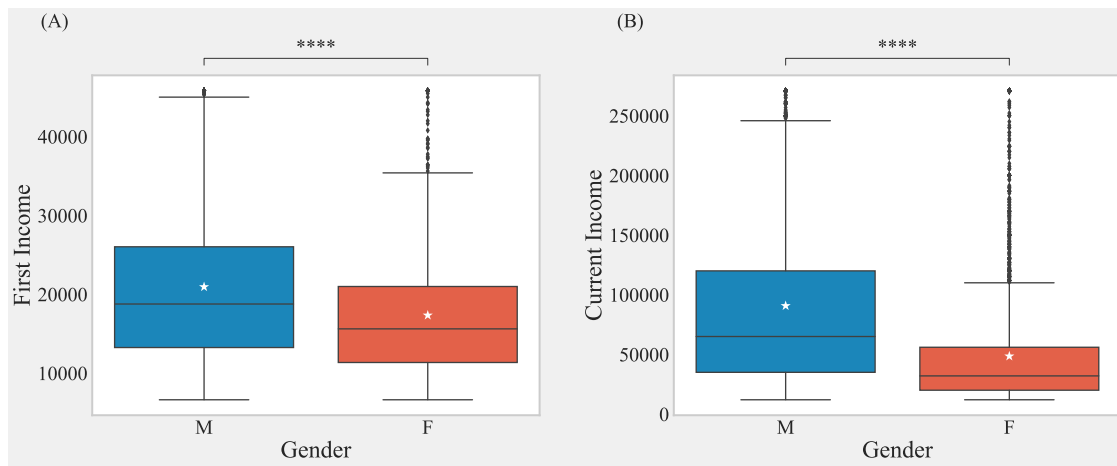


Fig. 5: First Income and Current Income boxplot based on Gender (male M, and female F) distribution count

cant in the ‘First Income’ result; however it becomes more significant in the latter score, where all of the comparisons obtained a significance level of 0.01%.

Table 3: First Income and Current Income median, mean, and standard deviation statistics by School variable

First Income			
Category	Mean	Std	Median
Business	\$ 18,398.68	\$ 9,164.97	\$ 16,558.00
Engineering	\$ 21,555.90	\$ 10,068.39	\$ 19,870.00
Other	\$ 15,922.15	\$ 8,990.18	\$ 13,533.00
Current Income			
Category	Mean	Std	Median
Business	\$ 78,838.27	\$ 68,905.98	\$ 53,707.50
Engineering	\$ 81,871.54	\$ 69,803.65	\$ 56,000.00
Other	\$ 55,154.67	\$ 57,171.84	\$ 34,657.00

Salary based Current Employment Characteristics

The box-plots in Fig. 7 and Fig. 8 presents a comparison between the most critical variables identified in the correlation analysis, which describe ‘Current Employment’ characteristics. These variables were not evaluated in the ‘First Income’ analysis as these are characteristics of the alumnus’s current status. In this analysis, we can see a significant difference between the number of years that the alumni have lived in a foreign country (outside of Mexico), showing higher values for those that have lived (and presumably worked) outside the longest (Table 4). The study did not present information about where these alumni have lived outside of Mexico. However, based on previous academic descriptive analysis, it was determined that 70% of the former students have migrated to North America. Fi-

nally, whether the former student has obtained a graduate degree or not has a significant difference (Table 5), showing a positive outcome for those that have achieved higher educational attainment (a Masters or Ph.D. degree).

Table 4: First Income and Current Income median, mean, and standard deviation statistics by Years Worked Foreign variable

Current Income			
Category	Mean	Std	Median
none	\$ 62,883.47	\$ 59,282.52	\$ 40,100.00
1-3	\$ 80,345.03	\$ 65,678.08	\$ 60,000.00
>3	\$ 137,887.29	\$ 81,580.06	\$ 117,362.00

Table 5: First Income and Current Income median, mean, and standard deviation statistics by Graduate Degree variable

Current Income			
Category	Mean	Std	Median
No	\$ 66,801.01	\$ 63,662.08	\$ 42,000.00
Yes	\$ 83,083.64	\$ 70,611.37	\$ 58,036.00

2.4 Modeling

In the previous section, the marginal analysis provided us a basic picture of the interrelation between selected variables in the survey with income. However, the results are limited as they only provide a descriptive statistic of bivariate association; they do not reflect relationships between covariates and their impacts. Therefore, a

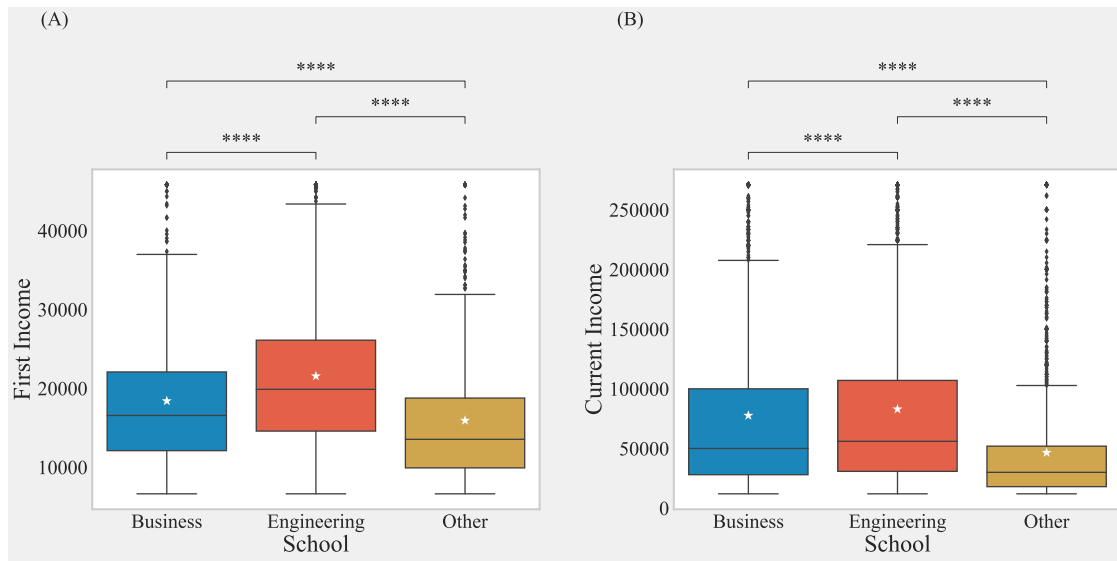


Fig. 6: First Income and Current Income boxplot based on School distribution count

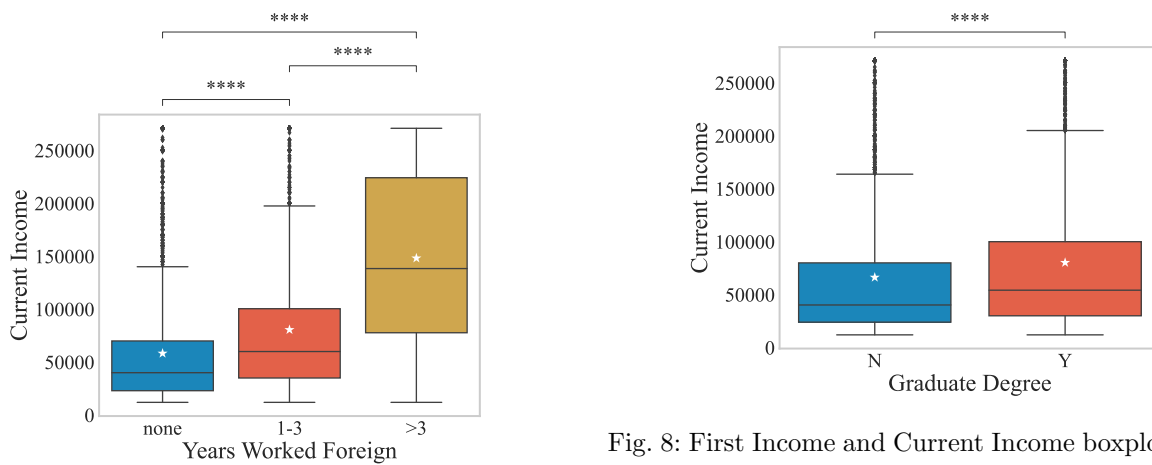


Fig. 7: First Income and Current Income boxplot based on Years Worked Foreign variable distribution count

multivariate analysis is explored to give us more precise assertions and the greatest predictive power. For this end, this section presents the different experiment configurations that were approached to analyse the alumni monthly income.

Since the most popular techniques used in econometric studies for income prediction are QR and linear regression, this research starts by evaluating these techniques and then comparing them with modern non-parametric machine learning algorithms, RF QR and GB QR. Then, to explain the most important factors related to alumni income, the study proposes exploring the data through a classification setting. This is done by discretisation of the dependent variables with multiple quartile categorisations and a median split.

Fig. 8: First Income and Current Income boxplot based on Graduate Degree variable distribution count

The modeling experimentation for this analysis was performed in four different configurations, labeled with letters: *Experiment A* involves the Quantile Regression; *Experiment B* involves the Traditional OLS Regression, *Experiment C* includes a Multi-Class Classification, and *Experiment D* is the Binary Classification.

The learning algorithms employed in each experiment are: QR, Multiple-Linear Regression, Logistic Regression, RF and GB. In the last experiment, the binary classification, we also integrated other machine learning models to compare their performance with the best achieving models. The evaluation contemplated a total of eight classifiers: Logistic Regression (LR), Linear Discriminant Analysis (LDA), K-Nearest Neighbours (KNN), a simple C4.5 Decision Tree (DT), Support Vector Machines (SVM), Naive Bayes Classifier (NB)

1 Random Forest Classifier (RFC) and Gradient Boost-
 2 ing Classifier (GBC).
 3

6 2.5 Regularization and Feature Selection

8 For the linear, QR, and LR models, we used a lasso
 9 regularisation. For the selection of the lambdas used in
 10 the penalty term, we performed a 10-fold-cross valida-
 11 tion in each of the models and used the loss function
 12 as the absolute error metric. While regularisation was
 13 obtained with this process, we found that the number
 14 of variables selected was still too many to design con-
 15 crete policies based on them. Therefore, to further re-
 16 duce the identified variables, we conducted a Sequen-
 17 tial Forward Floating Selection (SFFS) algorithm and
 18 evaluated the performance of the model with the top
 19 twenty most important variables. The SFFS is a float-
 20 ing variant of the traditional stepwise variable selection
 21 method [30]. It involves searching for the best subset
 22 of variables by adding and removing features at each
 23 step and evaluating the loss function, which can be the
 24 RMSE for regression or accuracy for classification.
 25

26 On the other hand, for the tree-based methods, we
 27 performed the Recursive feature elimination (RFE) [31]
 28 method with a 10-fold-cross-validation in the training
 29 set. This technique implements a backward selection; it
 30 starts with a model with all predictors and continuously
 31 evaluates the model’s score when removing each one of
 32 them. Those features with less importance are then re-
 33 moved from the final model. This method is frequently
 34 used with tree-based ensemble models since they can
 35 leverage the RF and GB internal methods for measur-
 36 ing feature importance [32].
 37

38 With the SFFS (Fig. 9) and RFE (Fig. 10) methods,
 39 we were able to select the most important variables for
 40 the model; twenty variables were selected for the ‘Cur-
 41 rent Income’ model and 16 for the ‘First Income’ model.
 42 The subset of features was selected based on the opti-
 43 mizing the loss function.
 44

48 2.6 Cross-Validation and Evaluation Metrics

50 For model development evaluation purposes, the datasets
 51 were split into a training set and a testing set, with
 52 80% and 20% of the data. The latter set was left out
 53 for the last evaluation process, and the training set was
 54 split into several splits with the use of stratified 10-
 55 fold cross-validation. This split was done to estimate
 56 the regressors and classifiers’ performance and to carry
 57 out the hyper-parameter tuning. The complete dataset
 58 was stratified uniformly so that there were all different
 59

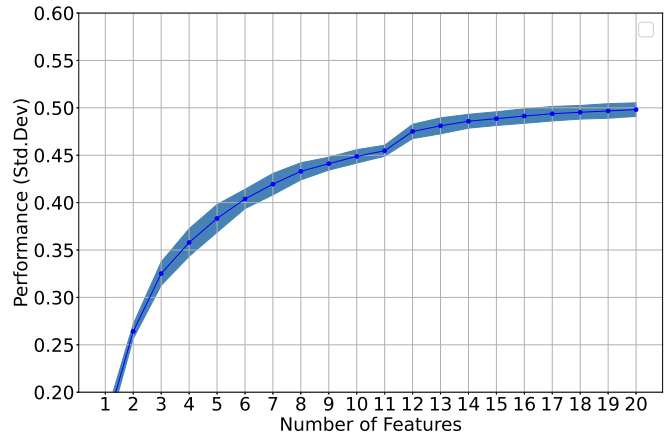


Fig. 9: Sequential forward floating selection linear regression with top 20 most important variables

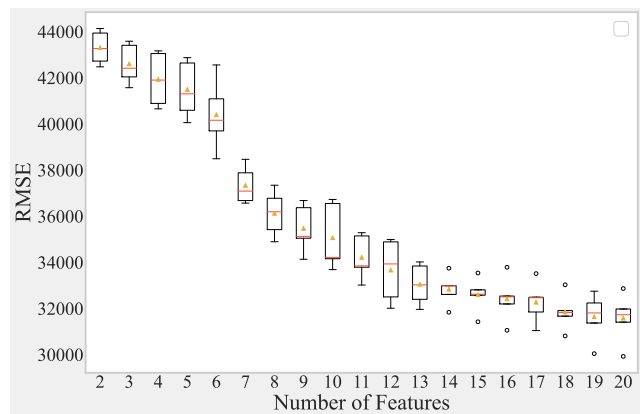


Fig. 10: Recursive feature elimination of gradient boosting model for the 20 most important variables

types of attributes’ values in both the Training set and the Test set.

The metrics evaluated with the cross-validation method were accuracy and Area Under the Curve (AUC) for the classification task and Root Mean Squared Error (RMSE) and adjusted-R2 for regression.

When working with a sample with high dimensionality, it is preferable to use the adjusted-R-squared-statistic as it penalizes the use of predictors that do not help explain the variation of the dependent variable [33, 34]. Equation 4 describes the adjusted-R-squared statistic, where n represents the sample size and k the number of features for the given observations in the analysis.

$$R_{adj}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - k - 1} \quad (4)$$

For the use of R^2 in quantile regression models, we use the pseudo-R-squared (Equation 5) defined by Koenker & Manchado in 1999 [35]. This metric allows the measurement of variability for a particular quantile

defined by τ . $\hat{V}(\tau)$ represents the pseudo R-squared for an unrestricted quantile regression model, while $\tilde{V}(\tau)$ is an intercept-only model. The pseudo-R-square metric value, such as in the traditional R-square, ranges between $[0,1]$. Still, it is a local measure of how well a particular quantile fits the model, not a global measure of goodness of fit for the total distribution [36].

$$R^1(\tau) = 1 - \frac{\hat{V}(\tau)}{\tilde{V}(\tau)} \quad (5)$$

To evaluate the loss function for linear regression, we measured the RMSE (Equation 6). To measure this in quantile regression, the quantile-loss error is used [37, 38]. This is also called the pinball loss and is similar to the Mean absolute Error (MAE) loss; however, it is not based on the mean but in the conditional quantile. The formula to obtain this value is shown in Equation 7.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (6)$$

$$L(y, t) := \begin{cases} (1 - \tau)(t - y) & \text{if } y < t \\ \tau(y - t) & \text{if } y \geq t \end{cases} \quad (7)$$

To measure classification models' performance, the confusion matrix and the following metrics are computed: overall accuracy (Equation 8), and the AUC. The latter measures the two-dimensional area that is underneath the receiver operating characteristic curve (ROC). The ROC curve is the graph that counts the number of correct positive classification gains in each of its thresholds; the curve plots the True Positive Rate (TPR) and the False Positive Rate (FPR) as defined below in Equation 9 and Equation 10.

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (8)$$

$$TPR = \frac{TP}{TP + FN} \quad (9)$$

$$FPR = \frac{FP}{FP + TN} \quad (10)$$

2.7 Model Interpretation Methods

With complex methods, such as ensemble methods or deep learning, more complex interactions are found by the algorithm; hence a higher accuracy can be obtained. The problem is that complex machine learning methods explaining how the prediction is assessed are not straightforward. The collection of post hoc interpretability methods that seek to convert 'black-box models' to 'glass-box models' are referred to as Explainable

Artificial Intelligence or XAI techniques [39]. We obtained the most important variables from the best performing model for all of the different approaches presented above. If the best performing model was a tree-based model, we used Shapley Additive exPlanations (SHAP). If it was the LR model, we obtained the most important features based on each variable's weighted coefficients.

SHAP helps explain how complex machine learning models make predictions and provides global interpretability using game theory and providing each feature with a SHAP value. SHAP values were introduced by Shapley [40]. They provide a way to distribute contributors' total gain (attribute's marginal contribution), assuming that all features contribute. The greater the Shapley value, the more positive effect it has on the objective function. SHAP values give feature attribution to each feature with the classical Shapley values from game theory.

For the last approach presented, we performed two XAI methods to explain the alumni income results. For this, we employed the binary model since categorisation can ease the presentation of variable effects. The first strategy consisted of visualising the interactions between two variables and their relation with the target variable with the use of SHAP dependence plots. SHAP dependence plots are a popular visualisation technique to summarise model predictions. This method is similar to the Partial Dependence Plot (PDP) introduced by Friedman [41]. They show how a feature relates to the model's target value. In the SHAP dependence plot, each observation is plotted as a scatter-plot point; the y-axis corresponds to the SHAP value and the x-axis to the attribute's value. By defining a different colour for each feature and showing them in a 2-D graph, we can visualise two variables' interaction effects.

For this study, the SHAP values were calculated and plotted in log-odds. Log-odds create a logistic transformation to the function, which provides visual attractiveness. When plotting the prediction's log-odds, we can see the effect between the feature inputs and the output value. With this unit, we can observe the change in the value of the target value when the predictor analysed is changed by one log-odd, and all the other variables are fixed. When the ratio is greater than 1, it indicates that the event is more likely to happen as the independent variable increases. In contrast, when the odds ratio is less than 1, the event is less likely to occur as the independent variable increases. The second strategy we used to interpret the results was mining rule-based patterns with the PBC4cip algorithm [42]. We used this method to identify insight regarding the

1 decision rules identified for better discrimination of the
 2 classes.

3 The PBC4cip algorithm is a model-specific method
 4 that uses an ensemble of decision trees and converts
 5 them into multivariate decision rules. As an example,
 6 a multivariate contrast pattern for income prediction
 7 could be the following: [IF Marital Status = Married
 8 AND Gender = Female AND Education = High School
 9 THEN Class = Low].

10 During the training phase, the PBC4cip algorithm
 11 weights the sum of the supports in each of the classes as
 12 stated in equation 11, where C represents the number
 13 of instances belonging to the class c , T the number of
 14 instances in the training dataset, P the set of patterns
 15 found for the class c , and $Sup(p, c)$ the support of the
 16 pattern p into the class c .

$$w_c = \frac{1 - \frac{C}{T}}{\sum_{p \in P} Sup(p, c)} \quad (11)$$

17 Then in the classification stage, the sum of each
 18 class's supports is multiplied by the weight w_c of its
 19 corresponding class. This is done to punish the high
 20 sum of supports computed by the majority class. Then,
 21 the instance evaluated is classified based on the class
 22 with the highest value according to Equation 12.

$$w(p, c) = w_c \sum_{p \in P} Sup(p, c) \quad (12)$$

23 Finally, in order to select the most relevant patterns,
 24 filtering is performed based on two constraints: sup-
 25 port difference and confidence above a minimal defined
 26 threshold. If the support difference or the confidence is
 27 not large enough, it is assumed that the pattern is not
 28 worthy of consideration. The rules obtained from this
 29 model were filtered by considering only those which had
 30 a support difference between both classes 40% or higher
 31 and confidence above 65% to ensure the rules were rel-
 32 evant for the prediction task.

3 Results

33 Table 6 and Table 7 presents the results for the four
 34 different experiments conducted in this paper for the
 35 prediction of current income. All of these methods were
 36 tuned with respect to their specific parameters and hyper-
 37 parameters by using Grid Search, and they all consid-
 38 ered a subset of the topmost important features selected
 39 in a pipeline process, using the SFFS method for the lin-
 40 ear models and the RFE for tree-based methods. When
 41 comparing the results, we show that the performance
 42 is better for the GB Model; this is true for both the
 43 regression and classification tasks. However, in the QR

model's statistical analysis, we identified that the im-
 44 provement was not significant.

45 To determine the significance, we used the data of
 46 the quantile loss obtained from the 10-fold CV for the
 47 three different models and performed a post hoc test to
 48 identify the pair of algorithms that do not have equal
 49 performance. For this statistical analysis, we used the
 50 post hoc Tukey HSD test. For Experiment A, the results
 51 of the Tukey HSD test showed there are no significant
 52 differences between the performance of the following
 53 models: QR50 and QGB50, we could not reject the hy-
 54 pothesis (p-value < 0.05) in any of the quantiles; thus,
 55 there was no sufficient statistical evidence to confirm
 56 that the results have a different distribution. Hence, by
 57 the parsimony theorem [43], we recommend using the
 58 traditional QR model for the QR approach. In contrast,
 59 the significance of the GB model in the traditional re-
 60 gression, the multi-class classification and the binary
 61 classification was significantly better than the rest of
 62 the models.

63 We integrated additional machine learning techniques
 64 in the binary approach and compared their accuracy
 65 scores in the 10-fold cross-validation. The current in-
 66 come model results for this approach are shown in Fig. 11.
 67 Based on the post hoc Tukey HSD test, we infer no sig-
 68 nificant differences within the following groups: GB and
 69 RF; SVC and LDA; KNN, DT and LR. All other dif-
 70 ferences were significant.

71 It can be observed in the results of the first income
 72 model shown in Table 8 and Table 9, and in the results
 73 of the post hoc Tukey HSD test for the binary models
 74 (Fig. 12), that for this model our hypothesis that non-
 75 parametric methods can perform better in this data
 76 does not hold. The results indicated that the linear and
 77 logistic regression were the most adequate to describe
 78 the variables' relationship with first income after grad-
 79 uation.

3.1 Feature Importance

80 For the 'Current Income' model, the ranking of the most
 81 important features and their overall contribution was
 82 plotted in a SHAP-values graph. This graph shows in
 83 red the variables that negatively impact the model and
 84 in green the ones that impact positively. The graphs in
 85 Fig. 13 show the features that impact the class 'High'.
 86 This technique is similar to obtaining the coefficients
 87 in a linear model and can bring transparency to our
 88 machine learning model. In this graph, we can see how
 89 17 of the subsets of variables impact the model posi-
 90 tively. For instance, 'Age' is the most important vari-
 91 able for 'Current Income' and impacts in a positive way;
 92 the 'Gender' variable follows this, the number of 'Years

Table 6: Regression models results of pseudo R2, quantile loss, adjustes R2, and root mean squared error for Current Income variable

	A			B	
	<i>Pseudo R2</i>	<i>Q-Loss</i>		<i>R2-adj</i>	<i>RMSE</i>
QR50	0.23	18,659.52	OLS	0.44	50,431.45
QRF50	0.37	14,719.26	RFR	0.51	47,325.67
QLGB50	0.38	14,301.58	GBC	0.54	45,892.69

Table 7: Classification models results of accuracy and area under the curve for Current Income variable

	C			D	
	<i>Accuracy</i>	<i>AUC</i>		<i>Accuracy</i>	<i>AUC</i>
LR	0.48	0.749	LR	0.79	0.870
RFC	0.50	0.762	RFC	0.82	0.890
GBC	0.53	0.796	GBC	0.83	0.910

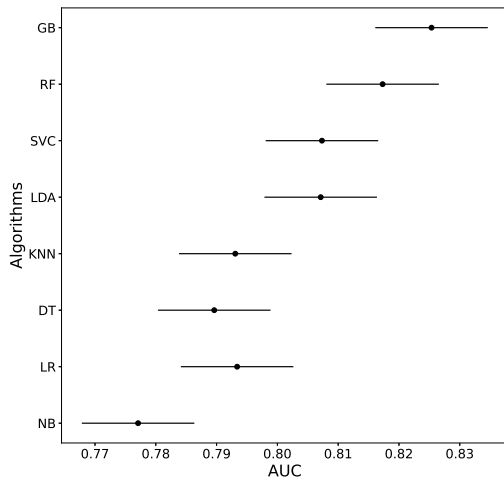


Fig. 11: Tukey HSD Post Hoc Test for the Binary Current Income classification models area under the curve results

worked Foreign’ and the ‘First Income’ variable. On the other hand, working in the ‘Tertiary Industry’ sector is impacting negatively. An interesting insight that can be noted is that having a high ‘Scholarship’ percentage during the alumni studies impacts negatively in their ‘Current Income’. However, this is affected by the proportion of the people with a scholarship vs. alumni without a scholarship; therefore, there might be an unfair bias for this variable.

On the other hand, something positive that can be observed is that ‘Bachelor GPA’ affects positively in many scenarios. Finally, a proxy variable is showing up in this graph, the *Negative Importance* variable, which shows that overall, giving importance to negotiation skills can boost the income of the alumni.

Regarding the ‘First Income’ model, a feature importance plot was obtained based on the estimated coef-

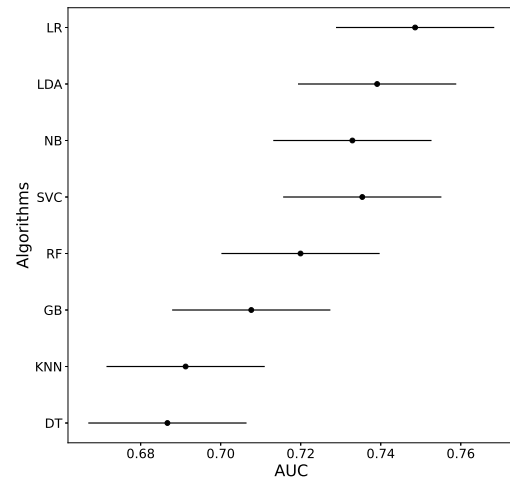


Fig. 12: Tukey HSD Post Hoc Test for the Binary First Income classification models area under the curve results

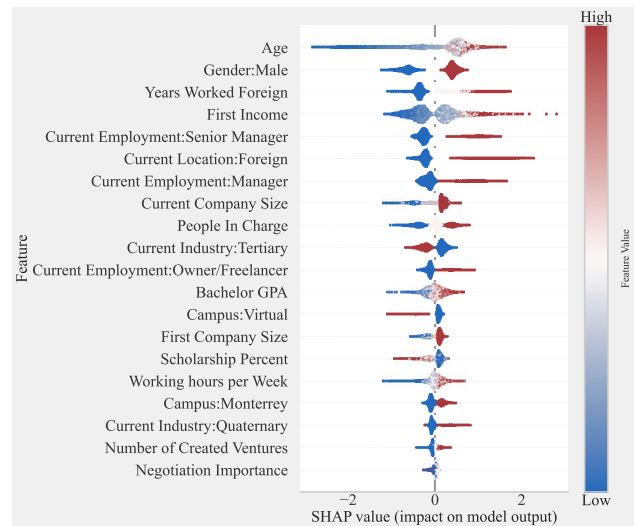


Fig. 13: GB Feature Importance with SHAP values

ficients of the 13 selected features for the LR model. The plot is shown in Fig. 14. When taking a close look at the coefficients, we can see that the features which impacted the model the most were the size of the company and having attended Engineering school, and having worked Foreign in their first job, all of these variables

Table 8: Regression models results of pseudo R2, quantile loss, adjustes R2, and root mean squared error for First Income variable.

	A			B	
	<i>Pseudo R2</i>	<i>Q-Loss</i>		<i>R2-adj</i>	<i>RMSE</i>
<i>QR50</i>	0.13	2,266.01	<i>OLS</i>	0.20	8,554.05
<i>QRF50</i>	0.10	3,167.97	<i>RFR</i>	0.19	8,609.08
<i>QLGB50</i>	0.13	3,090.62	<i>LGBR</i>	0.17	8,710.96

Table 9: Classification models results of accuracy and area under the curve for First Income variable

	C			D	
	<i>Accuracy</i>	<i>AUC</i>		<i>Accuracy</i>	<i>AUC</i>
<i>LR</i>	0.42	0.681	<i>LR</i>	0.69	0.75
<i>RFC</i>	0.37	0.6377	<i>RFC</i>	0.66	0.72
<i>LGBC</i>	0.37	0.6335	<i>LGBC</i>	0.65	0.71

impacted positively, whereas working in the quaternary sector, having lived in Mexico in a region different from the North or Central area, and being a First-Generation student impacted negatively in the income-class prediction.

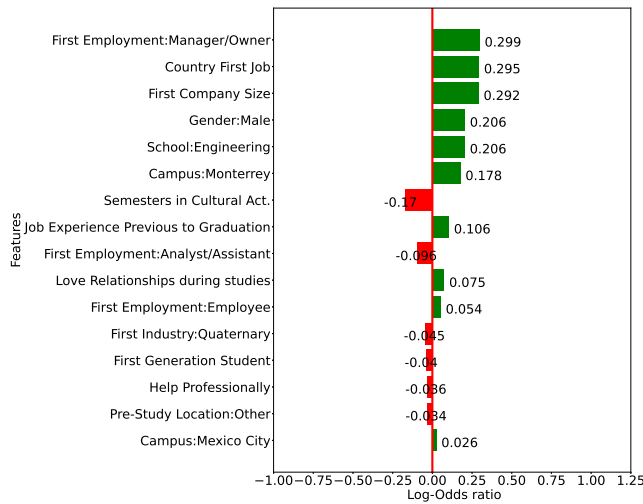


Fig. 14: First Income LR Feature Importance with coefficient weights

3.2 Exploring Feature Interactions

The SHAP partial dependence plots exhibit the marginal effect between two features on predicting the target variable. This visually shows covariates' relationships with the target variable besides being linear, monotonic or more complex. This section used SHAP partial dependence plots to show the stronger covariate relationships with income for the GB model for 'Current Income'.

Fig. 15 shows a strong interaction between the 'First Income' feature and Age. When looking at the observations, we can also notice this interaction, which is stronger for those former students between 35 and 40, seconded by those between 30 and 35. The interactions seem to be very weak for those older than 50 years old, which seems logical as these are people that graduated a long time ago, and thus the impact of their 'First Income' is not so relevant. Thus there is much likely a mix of other features not shown in this dataset that explain their variation.

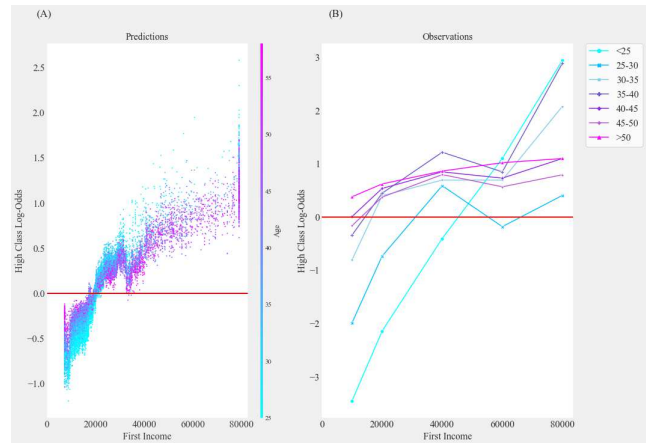


Fig. 15: First Income and Age SHAP dependency plot of the gradient boosting model, where (A) are the predictions and (B) are the real observations

The second interaction implies the relation between Working Hours and the Current Location of the alumnus. We can observe from Fig. 16 hat there is a linear relationship for working hours per week, which reaches a plateau. However, there is an interaction of the alumni's location. The predictions show that those working between 35 and 40 hours reach the peak if they live out-

side Mexico; yet, those living in Mexico do not reach this peak until working between 55 and 60.

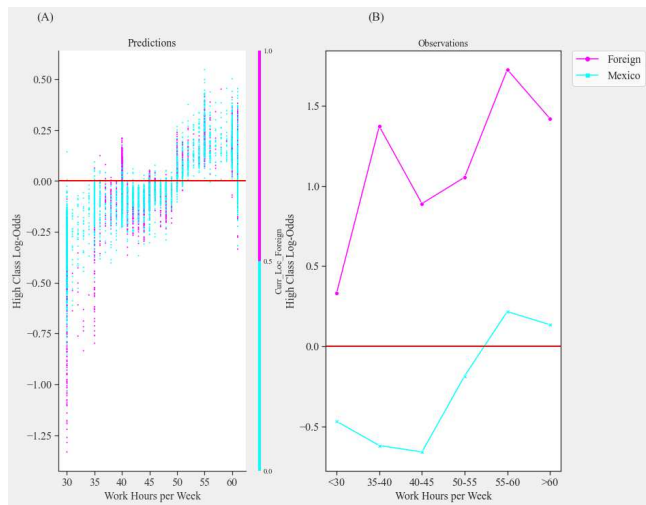


Fig. 16: Working hours per week and Current Location SHAP dependency plot, where (A) shows the predictions of the gradient boosting model and (B) the current observations

Another interesting insight that can be explained with this model is the interaction between gender and age. Graph B of Fig. 17 depicts how female alumni have a positive linear interaction when they are between 20 and 30 years old. Then this plateaus for the following years. On the other hand, graph A obtained from the model predictions shows how after 30, the variable of Gender impacts negatively in female alumni. Furthermore, the graph can be interpreted as gender having a negative effect on female alumni. We can see that the gap between Male and Female odds increases as the alumnus is older.

With the graphs presented in this section, we have identified that the most important variables for Income prediction identified by our model do not affect solely but interact with other covariates. The graph presented show the primary interaction relationships for the predictions along with a comparison with the observed data points versus log-odds.

3.3 Mining Contrast Patterns

The last technique used to analyse the 'Current Income' prediction model results was the contrast-pattern extraction with PBC4cip. In this section, an experiment was conducted to obtain contrast patterns that could give additional insight regarding the data analysed.

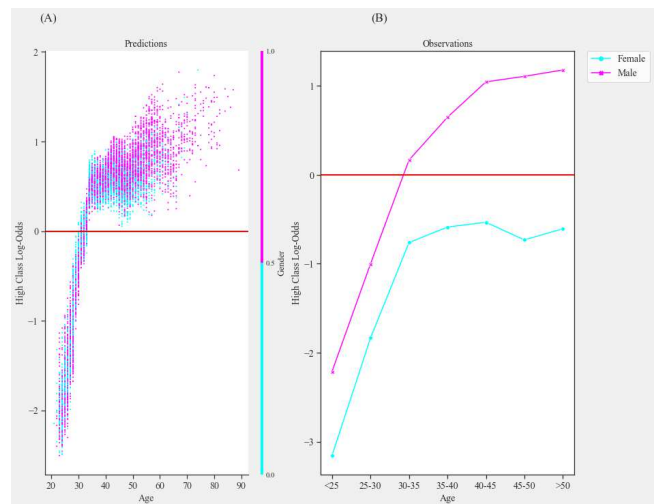


Fig. 17: Gender and Age SHAP dependency plot of the gradient boosting model, where (A) are the predictions and (B) are the observations

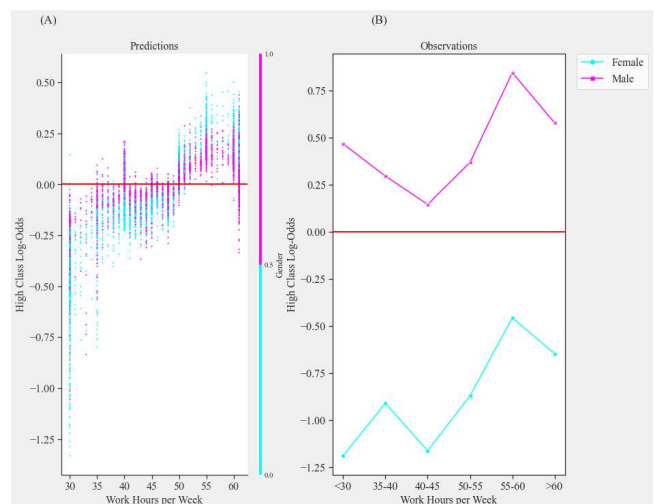


Fig. 18: Gender and Working Hours SHAP dependency plot of the gradient boosting model, where (A) are the predictions and (B) are the observations

One of the advantages of machine learning techniques over linear models is accounting for interactions between features. PBC4cip constructs rules by decomposing decision trees in a RF model, and any path that leads to a node can be transformed into a decision rule. The advantage of this is that the rules created are easy to interpret because, in our problem, they are binary decision rules. A limitation of data mining is that although it can identify patterns that are not obvious from data, not all of the patterns extracted may be useful. This is the most reason why data mining requires human intervention.

The rules obtained from this model were filtered by considering only those which had a support difference between both classes 40% or higher and confidence above 65%. This ensures that the rules are relevant for the prediction task. In addition, the redundant atoms obtained in the extracted patterns were removed with the automation filter in PBC4cip.

Three patterns were obtained that complied with these constraints and are shown in Table 10. This table shows that the set of patterns extracted each contain three features. To better comprehend the mathematical representation of the contrast patterns obtained, we used bar plots to visualise these three variables' impact on the 'High' class. The bar plots are shown in Fig. 19, Fig. 20 and Fig. 21.

In the first visualisation, we can observe how studying an Engineering or Business bachelor degree, having a job title different from Employee and being older than 28 years old gives the alumnus a higher probability for the 'High' class. As noted in Table 10, the support of this pattern for the 'High' class is 74%. This means that the pattern describes 74% of the observations with class 'High' (from the total dataset of 12,275 observations, where 5,877 belong to class 'High', 4,396 objects comply with this pattern). This pattern has a great coverage since the observation that it describes represents 35.6% of the overall observations in the dataset. Furthermore, the pattern confidence indicates that the probability that an object fulfils the property class 'High' given that the object fulfils the pattern is 64%.

Next, the second visualisation shows how being a Male alumnus, having a job title different from Employee, and being older than 28 years old, gives an alumnus a higher probability for the 'High' class. As noted in Table 10, the support of this pattern for the 'High' class is 61%, which indicates that the pattern describes 61% of the observations in class 'High'. There is also an extensive coverage since these observations represent 29.3% of the overall observations. The confidence given to this pattern is slightly higher than the previous one; it is 70%.

Finally, the visualisation in Fig. 21 depicts that having People in Charge, being older than 28 years old having a job title different from Employee also gives alumnus higher probabilities for the 'High' class, with support of 70%. Moreover, this pattern covers 33.7% of the observations in the dataset. The confidence for this pattern is 66%.

In this analysis, we mined three important patterns to contrast the two classes in our target variable, 'Current Income'. The variables that became evident in the obtained patterns were: School of Bachelor Degree, Job Title, Age, having people in charge and Gender. While

the first three are understandable variables to explain income, the latter variable has made evident the gender bias for the 'High' Class in the alumni population.

We note that the factors Age and the Job Title appear in all the patterns, and each one has a distinct variable. We can also see that the pattern which receives the most confidence is pattern 2, is where the distinct variable is gender; this shows us the importance that gender has been for the tree-based miner to determine the class.

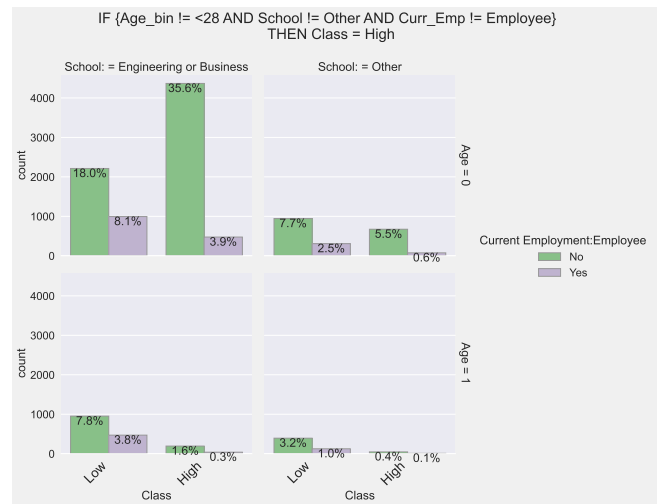


Fig. 19: Visualization of the first contrast pattern of Age, School, and Current Employment variables in bar plot count

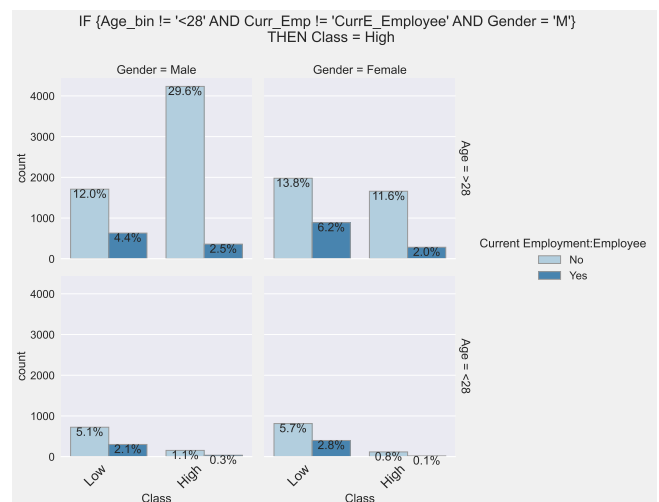


Fig. 20: Visualization of the second contrast pattern of binarized Age, Current Employment, and Gender variables in bar plot count

Table 10: PBC4cip three contrast patterns for current Income high class

ID	Pattern	Support by Class		Confidence
		Low	High	
CP1	IF Age_bin != '<28' AND School != 'School.Other' AND Curr_Emp != 'CurrE.Employee' THEN Class = 'High'	0.32	0.74	0.64
CP2	IF Age_bin != '<28' AND Curr_Emp != 'CurrE.Employee' AND Gender = 'M' THEN Class = 'High'	0.21	0.61	0.70
CP3	IF People_in_Charge != '0' AND Curr_Emp != 'CurrE.Employee' AND Age_bin != '<28' THEN Class = 'High'	0.3	0.7	0.66



Fig. 21: Visualization of the third contrast pattern of People in Charge, Current Employment, and binarization of Age variables in bar plot count

4 Discussions

When comparing the related work's results with the results from this thesis, we can see that in the QR approach, our results achieved better pseudo-R2 results than those obtained by Lee and Lee [7]. In contrast, the results from Figueiredo and Fontainha [10] had considerably better results. The main variables that were detected by the researchers and that were not available in the data set used in this thesis were: marital status, children status, the observation's firm's foreign capital, and the years of tenure at the current employer. Including these variables in the analysis as future work could serve positively to our model's performance.

For the traditional regression model approach, the results obtained by this study for the current income model were significantly better than the literature anal-

ysed. Both of the analysed related work used OLS to predict income, and with this study, we have shown that the decision-tree ensembles can yield significantly better results.

Unfortunately, in relation to the multi-class classification model, our results were worse than the literature analysed. Khongchai and Songmuang [12] obtained the best results using K-Nearest-Neighbours and Chen, Sun and Thakuriah [13] using Decision Trees. The former research used these additional dependent variables, which were not available in the dataset analysed in this thesis: specific degree programme and type of work performed in the company. While similar variables are included in this study, the analysed study variables consider more characteristics about the type of work that the students performed; other patterns could have been identified with the specific degree and work type. Therefore, this thesis hypothesis that follow-up research with more data can build models based on the specific degree and the type of work of the alumni. The latter research considered features provided by job descriptions from job posting sites. This included more work-related features such as location, contract versus permanent type, job content and job relationship features. While this work's objective was different from this thesis, it provided insights into how the detail of the job that the individual performs can be effectively used to predict their income.

Finally, for the binary model, this thesis obtained very similar results to related work. Lazar [5] showed that SVM could achieve high performance when predicting income. The author used the following predictors that were different from those used in this thesis: work class, marital status, race, capital gain, and capital loss. Further work can be done by including these variables in our GBC model to improve the performance. On the other hand, Sharath [11] achieved good results with boosted trees and various demographics

1 as predictors; however, our study achieved significantly
 2 better performance with the GBC and the variables
 3 identified as the most important for income prediction.
 4

5 Conclusions and Future Work

6
 7
 8 With the appearance of the digital transformation and
 9 the big data era, advanced analytics and data science
 10 has been increasingly used in many industries. In ed-
 11 ucation, it has been used to improve the learning pro-
 12 cess and evaluate academic institutions' efficiency. In
 13 econometric sciences, these techniques have been used
 14 to explain the links between economic, financial and so-
 15 cial effects. The differences between data analytics and
 16 data science are mainly that the latter makes use of
 17 machine learning techniques. These methods can pro-
 18 vide more accurate predictions than the traditional sta-
 19 tistical models used in data analytics. However, these
 20 methods' disadvantage is that they do not provide a
 21 clear interpretation of individual factors compared to
 22 conventional statistical methods. Hence, data analytics
 23 continues to dominate in education and econometric
 24 studies because of the ease of interpretation and the
 25 ability to distinguish variable effects.

26
 27 In this study, we show an application of the data sci-
 28 ence project life cycle to predict and identify the vari-
 29 ables with a strong relationship with alumni income.
 30 For this, we use 'the CRISP-DM methodology'. We fol-
 31 lowed the standard steps in the strategy and imple-
 32 mented additional steps to explain the results. Given
 33 this, we illustrate the flexibility that CRISP-DM can
 34 provide to data science projects based on the business's
 35 needs or research.

36
 37 We showed the importance of cleansing and trans-
 38 forming the data during this project's data understand-
 39 ing and preparation phase. Before modelling, we showed
 40 the importance of an exploratory analysis to under-
 41 stand the data, detect bias and identify specific pre-
 42 processing needs through the cleansing and transforma-
 43 tion process. The data exploration included descriptive
 44 statistics, visualisations through box-plots, correlation
 45 analysis, and the application of hypothesis testing for
 46 comparing two-factor levels and determining marginal
 47 effects of the independent variables with income.

48
 49 We compared different modelling techniques based
 50 on a distinction between parametric and non-parametric
 51 models during the modelling phase and utilised XAI
 52 techniques to interpret the results. The purpose of the
 53 study was to investigate the relationship between the
 54 target variables 'Current Income' and 'First Income'
 55 with demographical attributes obtained from an alumni
 56 survey. For this purpose, this research created and anal-
 57 ysed several machine learning methods to predict the
 58
 59
 60
 61
 62
 63
 64
 65

first income after graduation and former students' cur-
 rent income.

This study identified that for the best performing
 classification task, which discerns between low and high
 earners were, the top most important variables were:
 years worked foreign, first income, age, employment ti-
 tle, gender, employer's characteristics (company size,
 industry), the number of people in charge, the bache-
 lors GPA, and the working hours per week. While most
 of these variables are control variables, we identified the
 following actionable variables: bachelor's GPA, years
 worked foreign, working hours per week and first income
 after graduation. Hence, these variables can be paid
 more attention by those students seeking to achieve a
 high expected salary. Furthermore, this study's insights
 can be used to influence changes in the work sector
 and academic institutions, mainly to drive salary trans-
 parency and reduce the gender wage gap.

There are some interesting directions in which this
 work could be extending:

1. In this study, we only focused on comparing tradi-
 tional econometric algorithms with ensemble tree-
 based algorithms; it will be interesting to learn the
 performance of neural networks and explore the power
 of XAI techniques in deep learning.
2. Other educational institutions can use the method-
 ology followed in this study to perform a similar
 analysis to evaluate their alumni outcomes, iden-
 tify bias and provide them additional opportunities
 for obtaining their expected earnings.
3. Future work can consider the variables identified as
 more important in this study and augment the vari-
 ables provided in the related studies analyzed, such
 as marital status, children status, as well as more
 job-related characteristics.

Declarations

Ethical approval: This article does not contain any
 studies with human participants or animals performed
 by any authors.

Consent for publication: Not applicable.

Availability of data and materials: The data that
 support the findings of this study are available from Tec
 de Monterrey but restrictions apply to the availability
 of these data, which were used under license for the
 current study, and so are not publicly available.

Competing interests: The authors declare that they
 have no competing interests. All authors certify that
 they have no affiliations with or involvement in any or-
 ganization or entity with any financial interest or non-

1 financial interest in the subject matter or materials dis-
 2 cussed in this manuscript.

3 **Funding:** Not applicable.

4 **Authors' contributions:** DGC analyzed and inter-
 5 preted the alumni data regarding their income, per-
 6 formed the experiments reported, and was a major con-
 7 tributor in writing the manuscript. RDR performed a
 8 review and examination of the state of the art, and con-
 9 tributed in writing the manuscript. JP provided valu-
 10 able guidance for the statistical analysis and machine
 11 learning modeling. HC and NH collected the data used
 12 in the manuscript and provide feedback regarding the
 13 writing of the manuscript. All authors read and ap-
 14 proved the final manuscript.

15 **Acknowledgments:** The authors express their thanks
 16 to the School of Engineering and Science, Tecnológico
 17 de Monterrey, Monterrey, México, for providing the data
 18 for the analysis and to the PNPC CONACYT program
 19 with its financial support for DGC and RDR to study
 20 their master's degree.

21 References

- 22 1. C.R. Pace, (1979)
- 23 2. A.M. Delaney, Tertiary Education and Management
 24 **10**(2), 89 (2004)
- 25 3. J.F. Volkwein, New Directions for Institutional Research
 26 **2010**(S1), 125 (2010)
- 27 4. J.C. Rode, M.L. Arthaud-Day, C.H. Mooney, J.P. Near,
 28 T.T. Baldwin, International journal of selection and as-
 29 sessment **16**(3), 292 (2008)
- 30 5. A. Lazar, in *ICMLA* (Citeseer, 2004), pp. 143–149
- 31 6. D. Webbink, J. Hartog, Economics of Education Review
 32 **23**(2), 103 (2004)
- 33 7. B.J. Lee, M.J. Lee, The journal of the Korean economy
 34 **7**(1), 1 (2006)
- 35 8. P. Oehrlein, (2009)
- 36 9. M. Strand, T. Truong. Predicting Student Earnings After
 37 College
- 38 10. M.d.C. Figueiredo, E. Fontainha, (2015)
- 39 11. R. Sharath, K.N. Nirupam, B.J. Sowmya, K.G. Srinivasa,
 40 in *2016 International Conference on Computation Sys-
 41 tem and Information Technology for Sustainable Solu-
 42 tions (CSITSS)* (IEEE, 2016), pp. 249–254
- 43 12. P. Khongchai, P. Songmuang, in *2016 13th International
 44 Joint Conference on Computer Science and Software En-
 45 gineering, JCSSE 2016* (2016). DOI 10.1109/JCSSE.
 46 2016.7748896
- 47 13. L. Chen, Y. Sun, P. Thakuriah, in *International Confer-
 48 ence on Hybrid Intelligent Systems* (Springer, 2018), pp.
 49 61–74
- 50 14. D.T. Larose, *Discovering Knowledge in Data: An In-
 51 troduction to Data Mining* (2005). DOI 10.1002/
 52 0471687545
- 53 15. B. Leventhal, Journal of Direct, Data and Digital Mar-
 54 keting Practice **12**(2), 137 (2010)
- 55 16. R.W. Koenker, *Quantile Regression (Economet-
 56 ric Society Monographs)* (2005). DOI 10.1017/
 57 CBO9781107415324.004. URL [http://www.amazon.com/
 58 Quantile-Regression-Econometric-Society-Monographs/
 59 dp/0521608279](http://www.amazon.com/Quantile-Regression-Econometric-Society-Monographs/dp/0521608279)
- 60 17. A. Géron, *Hands-on machine learning with Scikit-Learn,
 61 Keras, and TensorFlow: Concepts, tools, and techniques
 62 to build intelligent systems* (O'Reilly Media, 2019)
- 63 18. S.B. Kotsiantis, Artificial Intelligence Review **39**(4), 261
 64 (2013)
- 65 19. A. Burkov, *The hundred-page machine learning book*,
 vol. 1 (Andriy Burkov Canada, 2019)
20. J. Friedman, T. Hastie, R. Tibshirani, *The elements of
 statistical learning*, vol. 1 (Springer series in statistics
 New York, 2001)
21. B. Efron, T. Hastie, *Computer age statistical inference*,
 vol. 5 (Cambridge University Press, 2016)
22. O. Troyanskaya, M. Cantor, G. Sherlock, P. Brown,
 T. Hastie, R. Tibshirani, D. Botstein, R.B. Altman,
 Bioinformatics **17**(6), 520 (2001)
23. W.J. Dixon, K.K. Yuen, Statistische Hefte **15**(2-3), 157
 (1974)
24. P. Embrechts, H. Schmidli, Zeitschrift für Operations Re-
 search **39**(1), 1 (1994)
25. R. Ibragimov, The New Palgrave Dictionary of Eco-
 nomics Online (2009)
26. A.J. McNeil, R. Frey, P. Embrechts, *Quantitative risk
 management: concepts, techniques and tools-revised edi-
 tion* (Princeton university press, 2015)
27. R. Ibragimov, A. Prokhorov, *Heavy tails and copulas:
 topics in dependence modelling in economics and finance*
 (World Scientific, 2017)
28. R.J. Hyndman, G. Athanasopoulos, *Forecasting: princi-
 ples and practice* (OTexts, 2018)
29. G.J. Myatt, W.P. Johnson, *Making sense of data II: A
 practical guide to data visualization, advanced data min-
 ing methods, and applications* (Wiley Online Library,
 2009)
30. M.A. Efroymson, Mathematical methods for digital com-
 puters pp. 191–203 (1960)
31. I. Guyon, J. Weston, S. Barnhill, V. Vapnik, Machine
 learning **46**(1-3), 389 (2002)
32. M. Kuhn, K. Johnson, *Feature engineering and selection:
 A practical approach for predictive models* (CRC Press,
 2019)
33. J. Miles, Wiley StatsRef: Statistics Reference Online
 (2014)
34. C.C. Aggarwal, *Data mining: the textbook* (Springer,
 2015)
35. R. Koenker, J.A.F. Machado, Journal of the american
 statistical association **94**(448), 1296 (1999)
36. I. Kurzawa, J. Lira, Metody Ilościowe w Badaniach Eko-
 nomicznych **16**(2), 33 (2015)
37. I. Steinwart, A. Christmann, Bernoulli **17**(1), 211 (2011)
38. F. Zhang, X. Fan, H. Xu, P. Zhou, Y. He, J. Liu, arXiv
 preprint arXiv:1911.05441 (2019)
39. A. Rai, Journal of the Academy of Marketing Science
48(1), 137 (2020)
40. L.S. Shapley, Contributions to the Theory of Games
2(28), 307 (1953)
41. J.H. Friedman, Annals of statistics pp. 1189–1232 (2001)
42. O. Loyola-González, M.A. Medina-Pérez, J.F. Martínez-
 Trinidad, J.A. Carrasco-Ochoa, R. Monroy, M. García-
 Borroto, Knowledge-Based Systems **115**, 100 (2017)
43. K.S. Bordens, B.B. Abbott, *Research design and meth-
 ods: A process approach* (McGraw-Hill, 2002)