

Better information on regression coefficients in predictive capacity studies with missing outcomes

Inés M. Varas ^{(0), †‡} Eduardo Alarcón-Bustamante ^{(0), ¶§‡||} and Jorge González ^{(0)†‡||}

†Department of Statistics, Pontificia Universidad Católica de Chile, Santiago, Chile

‡Interdisciplinary Laboratory of Social Statistics, Santiago de Chile, Chile

¶Departamento de Evaluación Medición y Registro Educacional (DEMRE), Universidad de Chile, Santiago de Chile, Chile §Escuela de Psicología, Pontificia Universidad Católica de Chile, Santiago de Chile, Chile

||Millennium Nucleus on Intergenerational Mobility: From Modelling to Policy (MOVI), Santiago de Chile, Chile *Corresponding author. Email: imvaras@uc.cl

Abstract

Tests play a central role in psychology and education, particularly in selection processes such as university admissions, where their predictive capacity—how well they forecast future performance—is typically assessed using regression models. However, a critical challenge arises in such analyses: while test scores are available for all applicants, outcome data (e.g., grade point averages) are only observed for selected individuals. Consequently, neither the regression model nor its parameters are identifiable from the sampling process.

Traditional approaches to address this problem rely on strong, often unrealistic assumptions about the missing data mechanism. Recent work has turned to partial identification, which provides bounds for regression functions and their parameters rather than point estimates. Building on Stoye (2007), this paper advances this approach by applying narrower and more informative identification bounds. Using real data from the Chilean university admissions system and contextually grounded assumptions, we demonstrate how these refined bounds enable modeling of selection process complexities. Our results underscore the value of partial identification in enhancing predictive capacity studies, offering a data-driven methodology to better understand the relationship between test scores and performance in selection settings.

Keywords: predictive capacity, ignorability, identification bounds

1. Introduction

Tests play a fundamental role in psychology and education, serving as tools for assessment, diagnosis, decision-making, and research. Their use is grounded in psychometric principles, ensuring validity, reliability, and fairness (Borsboom et al., 2004; Embretson & Reise, 2000; Messick, 1989). Particularly in selection processes such as in employment or university admission processes, they serve as critical tools for distinguishing individuals who are likely to succeed from those who may not. The selection decision is based on the assumption that test scores are indicative of future performance (Lord, 1980). This assumption underscores the need for tests to have strong predictive capacity, which is the ability of test scores to accurately forecast future outcomes.

Studies on the predictive capacity of selection tests are typically conducted through correlation analyses (Grassau, 1956; Lawley, 1943; Makransky et al., 2017; Pearson, 2013; K. Pearson, 1903; Thorndike, 1949) or regression models (e.g., Alarcón-Bustamante et al., 2021; Ayers & Peters, 1977; Manzi & Carrasco, 2021). In selection settings, positive regression coefficients are expected, indicating that higher test scores are linked to better performance. However, regression coefficients can only be estimated when the conditional distribution of the outcome given the explanatory variables is fully observed, or instead when suitable structural assumptions are imposed on the unobserved data, i.e., when the conditional distribution is identified (Koopmans, 1949). When assessing the predictive capacity of selection tests, this is not the case as the performance is only observed for those who were selected. In university admissions, for example, test scores are available for all applicants, but their academic performance (e.g., the graded point average (GPA) at the first year) is only recorded for those who gained admission. In studies on predictive capacity, this issue—referred to as the selection problem or range restriction—occurs when the performance of non-selected applicants is not observed (Manski, 1993).

To tackle the problem of missing outcomes in this type of studies, researchers have explored various approaches. A common practice in the literature is to restrict the analysis to students who were selected, and to draw inferences about the predictive capacity of selection tests based solely on this subset of the population (see, for instance, Alarcón-Bustamante et al., 2021; Geiser & Studley, 2002). This practice is justified using the assumption that the data are Missing At Random (MAR), which is also called *ignorability* (Florens & Mouchart, 1982; Hirano & Imbens, 2004; Imbens, 2000; Manski, 2013; Rosenmbaum & Rubin, 1983). Under this assumption, it can be shown that the distribution of the outcome of interest conditional on test scores in the full population of applicants is equal to the corresponding distribution among those who were selected. As a consequence of this equality the parameters indexing both distributions are also equal, allowing to estimate the population-level conditional expectation-expressed in terms of regression coefficients- using only the data from selected individuals. However, this approach is often inconsistent with the data-generating process observed in university admissions, as it contradicts the fundamental premise of selection: non-selected individuals were excluded precisely because their expected performance differed, presumably being lower (Alarcón-Bustamante et al., 2025; Grassau, 1956).

A more recent approach is the use of partial identification, which does not aim to estimate an exact parameter value but rather it provides a range of plausible values, usually an interval, containing information about the parameter of interest (Tamer, 2010). The interval contains all the values for the parameters that are compatible with the researcher's believe and the available data (Manski, 1989). Thus, the wider (narrow) the identification interval, the lesser (larger) is the information about the parameter of interest.

Applications of partial identification analysis in several contexts can be reviewed in Diemer et al., 2024; Manski, 2016; Pepper, 2000; San Martín and González, 2022; San Martín et al., 2024; San Martín and Alarcón-Bustamante, 2022; Stoye, 2011; Stoye, 2007, among others. One important advantage of this approach is that it allows researchers to incorporate milder, contextually grounded assumptions about the non-observed population, leading to more credible results (Manski, 2003), in contrast to strategies based on assumptions on the probability distribution for the missing outcome data (e.g., ignorability), which are rarely justified in empirical research.

Stoye, 2007 derived identification intervals for regression coefficients in the case of incomplete outcome data. These intervals are based on identification bounds for the whole linear regression model. Using these results, and in the context of predictive capacity of tests in a university selection process, Alarcón-Bustamante et al., 2023 derived identification intervals for the regression coefficients based on the widest (i.e., less informative) identification bounds of the linear regression. On the other hand, Alarcón-Bustamante et al., 2025 derived more informative identification bounds for a regression function based on different beliefs about the selection process.

In this paper, we aim to improve the results in Alarcón-Bustamante et al., 2023 narrowing the identification intervals for regression parameters by using the more informative bounds of the regressions derived in Alarcón-Bustamante et al., 2025.

The paper is organized as follows. In Section 2 the partial identification approach is introduced to obtain identification bounds for both the regression functions and for regression coefficients,

when data contains missing outcomes. Building on Alarcón-Bustamante et al., 2025, Section 3 describes the main characteristics of the Chilean university admission process and, using real data and different assumptions based on the Chilean context, informative identification bounds for the regression function are presented. In Section 4, the results of Section 3 are used to obtain informative identification intervals for the regression coefficients in a linear regression model. The paper ends with a discussion in Section 5.

2. Partial identification approach

In this section, we first make explicit the identification problem in regression models, and thus also in regression coefficients, when missing outcomes are present in the data. Next, we introduce the partial identification approach to obtain bounds on the regression function which will then be used to obtain identification intervals for the regression coefficients.

2.1 Unidentifiability of the regression function

Consider regression data (X, Y) where Y is an outcome variable and X is a vector of regression covariates or explanatory variables. Let Z be a binary random variable such that Z = 1 if the outcome is observed and 0 otherwise. By the law of iterated expectations (or law of total probability Kolmogorov, 1950) the conditional expectation E(Y | X), i.e., the regression function, decomposes as

$$E(Y \mid X) = E(Y \mid X, Z = 1)P(Z = 1 \mid X) + E(Y \mid X, Z = 0)P(Z = 0 \mid X).$$
(1)

To simplify notation, in what follows we will use $\mu_Y(X) = E(Y \mid X)$, $\mu_Y(X, Z = z) = E(Y \mid X, Z = z)$ and $m(X) = P(Z = 0 \mid X)$.

In (1), the regression function $\mu_Y(X, 0)$ is unknown and cannot be estimated from the available data because the outcome is not observed in this case. Consequently, neither the complete regression function $\mu_Y(X)$, nor the parameters that might characterize it (e.g., the regression coefficients in a linear model) can be known and estimated either, i.e., they are unidentified (Koopmans, 1949). However, assumptions can be imposed on the unobserved regression so that it becomes identified (Manski, 1993). In particular, if it is assumed that the conditional expectations when Y is unobserved and observed are equal, i.e., $\mu_Y(X, 0) = \mu_Y(X, 1)$ (known as the weak ignorability assumption and also referred to as mean missing at random in the terminology of Manski, 2003), then the regression $\mu_Y(X)$ can be identified and estimated using only the information from the observed outcomes, i.e., $\mu_Y(X) = \mu_Y(X, Z = 1)$. In this case, the regression is said to be point-identified. The ignorability assumption serves as an identification restriction, allowing one to learn about the regression model and, in particular, its regression coefficients. It underlies the missing-at-random (MAR) framework and certain conditional imputation methods (Little & Rubin, 2002). However, such assumption may not be appropriate in certain contexts as it will be seen later.

In the following section, we elaborate on alternatives that allow learning about the regression function using the available data and milder assumptions.

2.2 Partial identification of the regression function

Although the regression function, $\mu_Y(X)$, is not identified, it can be partially identified in the following way (Manski, 1989). Assume that $\mu_Y(X, 0)$ is bounded by the functions $[h_0(X), h_1(X)]$, i.e. $h_0(X) \leq \mu_Y(X, 0) \leq h_1(X)$. Then, the regression $\mu_Y(X)$ in (1) can be bounded as

$$\mu_Y(X, Z = 1)(1 - m(X)) + h_0(X)m(X) \le \mu_Y(X) \le \mu_Y(X, Z = 1)(1 - m(X)) + h_1(X)m(X)$$
(2)

Bounds in (2) define a partial identification region for the regression function $\mu_Y(X)$. Thus, all possible regression functions which are consistent with the observed data lie in this region. For a

fixed X = x, this region becomes an interval, and its width is given by $m(x)(h_1(x) - h_0(x))$. As a consequence, the smaller the difference $h_1(x) - h_0(x)$ is, the narrower and thus more informative is the identification interval.

The bounds in (2) gives information about all the plausible regression functions under the assumption that the non-observed regression is bounded. However, although it is possible to obtain bounds for the regression function $\mu_Y(X)$, these are not bounds for the parameters that might characterize it. For instance, for X = x if a linear regression model $\mu_Y(x) = x^{\top}\beta$ is assumed, as it is usual in predictive capacity of test studies, then identification bounds on the regression coefficients, β , are of interest. The following section elaborates on how these bounds can be obtained.

2.3 Partial identification of the regression coefficients

In the context of missing outcome data when conditional distributions, and therefore conditional expectations, are partially identified, Stoye, 2007 proposed bounds for general linear predictors of the form $K(\mathbf{x}^{\top}\boldsymbol{\beta})$, where K is a strictly increasing function. In particular when K is taken to be the identity function, and the linear regression $\mu_Y(\mathbf{x})$ is bounded by $\left[\underline{\mu}_Y(\mathbf{x}); \overline{\mu}_Y(\mathbf{x})\right]$, then identification bounds for a linear combination of the regression coefficients, namely $c\boldsymbol{\beta}$ where $c \in \mathbb{R}^J$ is a row vector, are given by:

$$\boldsymbol{c} \Big[E(\boldsymbol{x}^{\top} \boldsymbol{x}) \Big]^{-1} E(\boldsymbol{x}^{\top} \underline{k}(\boldsymbol{x})) \leq \boldsymbol{c} \boldsymbol{\beta} \leq \boldsymbol{c} \Big[E(\boldsymbol{x}^{\top} \boldsymbol{x}) \Big]^{-1} E(\boldsymbol{x}^{\top} \overline{k}(\boldsymbol{x})).$$
(3)

Here, $E(\cdot)$ denotes the expectation over the distribution F_X of X (i.e., $E(\mathbf{x}^{\top}\mathbf{x}) = \int \mathbf{x}^{\top}\mathbf{x}dF_X$), and $\underline{k}(\mathbf{x})$ and $\overline{k}(\mathbf{x})$ are defined such that

• if
$$cE(\mathbf{x}^{\top}\mathbf{x})^{-1}\mathbf{x}^{\top} > 0$$
, then $\underline{k}(\mathbf{x}) = \underline{\mu}_{Y}(\mathbf{x})$ and $\overline{k}(\mathbf{x}) = \overline{\mu}_{Y}(\mathbf{x})$; and
• if $cE(\mathbf{x}^{\top}\mathbf{x})^{-1}\mathbf{x}^{\top} \leq 0$, then $\underline{k}(\mathbf{x}) = \overline{\mu}_{Y}(\mathbf{x})$ and $\overline{k}(\mathbf{x}) = \underline{\mu}_{Y}(\mathbf{x})$.

Since (3) holds for a linear combination of the regression coefficients, identification bounds for β_l , $l \in \{1, ..., J\}$, can be obtained choosing *c* as the *l*-th canonical vector in \mathbb{R}^J . For further details on the derivation of these results, see Stoye, 2007.

Alarcón-Bustamante et al., 2023 used this approach to study the predictive capacity of test scores in the context of university admissions. These authors found identification bounds for the regression coefficients assuming that $\underline{\mu}_Y(\mathbf{x})$ and $\overline{\mu}_Y(\mathbf{x})$ were the minimum and maximum values attainable of the GPA, respectively. In this paper, we extend their work by incorporating additional assumptions based on the selection process, which enable the derivation of more informative bounds for the regression function and, consequently, for the regression coefficients. Before presenting these assumptions, it is essential to understand how the selection process operates in practice. In the following section, we describe the Chilean university admission process in detail and outline how this knowledge can be used to formulate more informative assumptions for the analysis.

3. The Chilean university admission process and an empirical study

The derivation of more informative bounds for the partially identified regression coefficients will be illustrated with an empirical study that aims to examine the predictive capacity of selection factors in the Chilean university admission process.

In Chile, the process begins with applicants taking a battery of standardized selection tests, from which test scores are obtained. Other selection factors considered are the high school GPA (HS-GPA), and a variable called Ranking, which reflects the applicant's relative position at the end of high school. All selection factors are measured on a scale from 150 to 850 and are combined using predetermined

weights to obtain a single composite score for the admissions process. Applicants whose composite score meets or exceeds the university's established cutoff are admitted. Consequently, even if two applicants obtained identical scores on a particular test, their overall selection status (selected or non-selected) may differ due to differences in the values of the remaining selection factors.

Based on the composite score and their personal preferences, applicants then decide to apply to one or more undergraduate programs. Subsequently, a selection process is carried out in which applicants are evaluated and a proportion of them is admitted to a single undergraduate program among those they applied to.

The empirical study uses data of applicants to three programs at the School of Biology in a university in Chile: Marine Biology (MB), Biology (B), and Biochemistry (BC). The selection factors we consider are the standardized test scores in Mathematics and Language, high school GPA (HS-GPA), and the applicant's ranking. The outcome variable corresponds to the academic performance as measured through the first-year university GPA on a 1.0 to 7.0 scale, where 4.0 is the minimum passing grade.

In the selection process, the School of Biology uses predetermined cut-off scores to filter applicants. To be selected for programs MB, B, or BC, an applicant must achieve a minimum score of $\tau_1 = 630.8$, $\tau_2 = 643.6$, or $\tau_3 = 701.5$, respectively. This information will be used to formulate different beliefs about the selection process.

3.1 Mathematical formulation

Building on Alarcón-Bustamante et al., 2025, each applicant will be characterized by (Y, X, G, U), where Y is the GPA at the first year at university, X is a vector of selection factors, G denotes the selection status with levels 0 = non-selected, 1 = selected in MB, 2 = selected in B, and 3 = selected in BC, and U is a vector with entries either 1 (apply) or 0 (did not apply) denoting the application status. For our study, $U = (u_1, u_2, u_3)$, where u_1, u_2 , and u_3 indicate whether applicant applied to MB, B, or BC, respectively. In addition, all individuals applied to at least one of the three specified undergraduate programs so the status U = (0, 0, 0) is not considered.

By the law of iterated expectations the regression function can be written as

$$\mu_Y(X) = \sum_{u \in \mathcal{U}} \mu_Y(X, G = 0, U = u) P(G = 0, U = u \mid X) + \sum_{g=1}^3 \mu_Y(X, G = g) P(G = g \mid X),$$
(4)

where \mathcal{U} is the set of all possible application statuses. To simplify notation, let us denote by $P_{G,U}(X, g, u) = P(G = g, U = u | X)$ and $P_G(X, g) = P(G = g | X)$. Note that all terms in (4) can be identified from the sampling process, with the exception of the regression corresponding to non-selected applicants, i.e., $\mu_Y(X, G = 0, U = u)$ for all $u \in \mathcal{U}$. As bounds for this regression, Alarcón-Bustamante et al., 2025 considered $h_0(X) = y_0$ and $h_1(X) = y_1$, the minimum and maximum possible value of the GPA respectively, obtaining the widest interval in (2) for the regression model. Assuming that the non-observed regression is bounded by the minimum and maximum possible GPA was called the Weakly Informative Assumption (WIA). The authors also proposed using other assumptions for bounding the regression function considering beliefs about the selection process. In this study, we apply these assumptions to define bounds $[\underline{\mu}(X); \overline{\mu}(X)]$ for the regression function and then, by evaluating the results in (3) for a linear regression model, we obtain more informative bounds for the regression coefficients.

3.2 More informative assumptions

To obtain more informative bounds for the regression function, Alarcón-Bustamante et al., 2025 introduced assumptions about the selection process. These assumptions can be written in terms

of both the application status and the most attainable undergraduate programs that non-selected applicants would have been admitted to.

For non-selected applicants, the most attainable program depends on their application status. For instance, for an application status U = (1, 1, 1), the most attainable program is MB; whereas for U = (0, 1, 1) would be B. Thus, the set U of all possible application statuses can be partitioned into subsets $U_1 = \{(1, 0, 0), (1, 1, 0), (1, 0, 1), (1, 1, 1)\}, U_2 = \{(0, 1, 0), (0, 1, 1)\}, and U_3 = \{(0, 0, 1)\}$ where U_1, U_2 , and U_3 collect the application statuses for which MB, B, and BC are the most attainable programs, respectively. This partition will be useful for the specification of the assumptions considered about the selection process.

Specifically, applicants in U_1 could only have been admitted to Marine Biology, as their scores were below the cut-off $\tau_1 = 630.8$. Similarly, those in U_2 could only have been admitted to Biology, given their scores were below $\tau_2 = 643.6$, and applicants in U_3 could only have been admitted to Biochemistry, as their scores were below $\tau_3 = 701.5$. This information is instrumental in representing different beliefs about the selection process, as it allows us to define alternative bounds for the regression model that go beyond those derived under the WIA assumption.

Building on this framework, in what follows we briefly revisit and summarize three assumptions that lead to more informative bounds for the regression function: Perfect Selection Assumption (PSA), Worst Selection Assumption (WSA), and Fallible Selection Assumption (FSA). A more detailed description of the formulation of these assumptions, as well as the mathematical derivations can be seen in Alarcón-Bustamante et al., 2025.

3.2.1 Perfect Selection Assumption (PSA)

Under this scenario, it is assumed that the tests perfectly select applicants in the sense that, given the application status, the expected performance of non-selected applicants would be worse than the observed for those who were selected in the most attainable program the non-selected applicant would have been admitted to.

Let m_{0g} represent the minimum observed GPA among selected applicants in program g, where g = 1, 2, 3 corresponds to MB, B, and BC, respectively. Under the PSA, the expected performance of non-selected applicants with an application status $u \in U_1$ would be at most m_{01} . Similarly, for non-selected applicants with an application status $u \in U_2$, the mean GPA would be at most m_{02} , and for those with an application status $u \in U_3$, the expected GPA would be at most m_{03} . These conditions can be expressed in terms of the regression function as follows: for all $g \in \{1, 2, 3\}$,

$$\mu_Y(X, G = 0, U = u) \leqslant m_{0g} \quad \text{for } u \in \mathcal{U}_g.$$

In terms of lower bounds, the PSA does not restrict the expected mean of non-selected applicants. Since the GPA is bounded below by y_0 , it follows that for all application statuses of non-selected applicants, the condition

$$\gamma_0 \leq \mu_Y(X, G = 0, U = u)$$

holds. This ensures that the expected GPA of non-selected applicants is at least γ_0 , regardless of their application status. Thus, under the PSA assumption the regression function in (4) can be bounded by $\left[\underline{\mu}_V(X), \overline{\mu}_Y(X)\right]$, where

$$\begin{split} \underline{\mu}_{Y}(X) &= \gamma_{0} P_{G}(X,0) + \sum_{g=1}^{3} \mu_{Y}(X,G=g) P_{G}(X,g), \\ \overline{\mu}_{Y}(X) &= \sum_{g=1}^{3} \left\{ m_{0g} \sum_{u \in \mathcal{U}_{g}} P_{G,U}(X,0,U=u) \right\} + \sum_{g=1}^{3} \mu_{Y}(X,G=g) P_{G}(X,g). \end{split}$$

3.2.2 Worst Selection Assumption (WSA)

This scenario is entirely contrary to the optimistic viewpoint of the PSA. It is assumed that nonselected applicants would perform better than those who were selected, suggesting that the admission process may have selected poorly.

Under this assumption, the expected GPA of non-selected applicants who applied to the MB, B, and BC programs, would be higher than that of those who were selected. Formally, for $g \in \{1, 2, 3\}$ and $u \in U_q$,

$$\mu_Y(X, G = g) \leq \mu_Y(X, G = 0, U = u).$$

Additionally, if γ_1 denotes the maximum value of the GPA range, the potential GPA of all non-selected applicants, given their application status, would be at most the upper limit of the GPA. That is, for all $u \in U$,

$$\mu_Y(X, G = 0, U = u) \leq \gamma_1.$$

Thus, by considering the WSA, the bounds for the regression model in (4) are given by

$$\underline{\mu}_{Y}(X) = \sum_{g=1}^{3} \mu_{Y}(X, G = g) \left\{ P_{G}(X, g) + \sum_{u \in \mathcal{U}_{g}} P_{G, U}(X, g, u) \right\}$$

$$\overline{\mu}_{Y}(X) = \gamma_{1} P_{G}(X, 0) + \sum_{g=1}^{3} \mu_{Y}(X, G = g) P_{G}(X, g)$$

3.2.3 Fallible Selection Assumption (FSA)

Noting that both the PSA and the WSA are extreme assumptions regarding the selection process, which may even be unrealistic for a selection process, Alarcón-Bustamante et al., 2025 proposed the *Fallible Selection Assumption (FSA)*, under which it is assumed that the selection process might make mistakes in selecting applicants but it is not completely deficient.

Under this alternative assumption, given the selection status, the potential performance of nonselected applicants is bounded below by the minimum GPA observed among selected applicants in their most attainable program. Specifically, for $g \in \{1, 2, 3\}$, the expected GPA of non-selected applicants with an application status to program g, i.e. for all $u \in U_g$, satisfies

$$m_{0\sigma} \leq \mu_Y(X, G = 0, U = u).$$

These lower bounds reflect a potential imperfection in the admission process. Assuming that the selection process is not entirely flawed, the expected potential performance of non-selected applicants with a specific application status would not exceed the mean performance of selected applicants in their most attainable program. Formally, it holds:

$$\mu_Y(X, G = 0, U = u) \leq \mu_Y(X, G = g),$$
 for all $u \in \mathcal{U}_g$,

where $g \in \{1, 2, 3\}$ denotes the MB, B, and BC programs, respectively. Thus, by considering the lower and upper bounds for the regression function of non-selected applicants, the bounds of the complete regression model under the FSA are given by the terms $\mu(X)$ and $\overline{\mu}(X)$:

$$\underline{\mu}_{Y}(X) = \sum_{g=1}^{3} \left\{ m_{0g} \sum_{\boldsymbol{u} \in \mathcal{U}_{g}} P_{G,\boldsymbol{U}}(\boldsymbol{X}, \boldsymbol{0}, \boldsymbol{u}) \right\} + \sum_{g=1}^{3} \mu_{Y}(\boldsymbol{X}, \boldsymbol{G} = g) P_{G}(\boldsymbol{X}, g),$$

$$\overline{\mu}_{Y}(\boldsymbol{X}) = \sum_{g=1}^{3} \mu_{Y}(\boldsymbol{X}, \boldsymbol{G} = g) \left\{ P_{G}(\boldsymbol{X}, g) + \sum_{\boldsymbol{u} \in \mathcal{U}_{g}} P_{G,\boldsymbol{U}}(\boldsymbol{X}, g, \boldsymbol{u}) \right\}.$$

Although the FSA relaxes the extremes of perfect and worst-case selection, alternative trade-off assumptions can be considered to characterize the selection process in ways that differ from those previously analyzed. In a more conservative approach, distinct assumptions regarding the upper bound of performance among non-selected applicants can be formulated. We present two of such alternatives.

3.2.4 Mean Most Attainable Assumption (MMAA)

Suppose that the selection process guarantees that on average, non-selected applicants cannot outperform those who were admitted, while placing no restriction on their minimum potential performance.

Under this approach, the mean performance of the non-selected applicants in a specific application status u is bounded by the mean performance of those selected applicants in their most attainable program, i.e.

$$\mu_Y(X, G = 0, U = u) \leq \mu_Y(X, G = g), \quad \text{for all } u \in \mathcal{U}_g,$$

where $g \in \{1, 2, 3\}$ denotes the MB, B and BC programs, respectively. In the words of Manski and Pepper (2000), the Mean Most Attainable assumption asserts that selected applicants in a specific program have a weakly higher mean GPA than those who applied but were not selected. In addition, since the GPA is bounded by γ_0 , it holds that

$$\gamma_0 \leq \mu_Y(X, G = 0, U = u).$$

Accordingly, under the MMAA, the complete regression model $\mu_Y(X)$ is bounded by

$$\underline{\mu}_{Y}(\boldsymbol{X}) = y_{0}P_{G}(\boldsymbol{X}, 0) + \sum_{g=1}^{3} \mu_{Y}(\boldsymbol{X}, G = g)P_{G}(\boldsymbol{X}, g),$$

$$\overline{\mu}_{Y}(\boldsymbol{X}) = \sum_{g=1}^{3} \mu_{Y}(\boldsymbol{X}, G = g) \left\{ P_{G}(\boldsymbol{X}, g) + \sum_{\boldsymbol{u} \in \mathcal{U}_{g}} P_{G, \boldsymbol{U}}(\boldsymbol{X}, g, \boldsymbol{u}) \right\}.$$

3.2.5 Conservative Most Attainable Assumption (CMAA)

A more conservative approach regarding the performance of those non-selected applicants is to consider that the mean performance of those non-selected applicants in a specified application status cannot exceed the maximum observed performance of those selected applicants in their most attainable program. We call this approach the *Conservative Most Attainable Assumption (CMAA)*. Under this assumption, the expected GPA of non-selected applicants who applied to the MB, B, and BC programs, would be at most equal to the maximum observed GPA in the most attainable program. This is,

$$\mu_Y(X, G = 0, U = u) \leqslant m_{1g}, \quad \text{for all } u \in \mathcal{U}_g,$$

where $g \in \{1, 2, 3\}$ denotes the MB, B and BC program, respectively and m_{1g} denotes the maximum observed GPA of the selected applicants under the respective program. In addition, because the potential performance of those non-selected applicants is at least γ_0 for all application status $u \in U_g$, for $g \in \{1, 2, 3\}$ it holds that

$$\gamma_0 \leqslant \mu_Y(X, G = 0, U = u)$$
 for all $u \in \mathcal{U}_{\varphi}$.

Consequently, the lower and upper bounds for the complete regression model in (4) are defined as:

$$\begin{split} \underline{\mu}_Y(X) &= \gamma_0 P_G(X,0) + \sum_{g=1}^3 \mu_Y(X,G=g) P_G(X,g), \\ \overline{\mu}_Y(X) &= \sum_{g=1}^3 m_{1g} \left\{ P_G(X,g) + \sum_{\mathbf{u} \in \mathcal{U}_g} P_{G,\mathbf{U}}(X,g,\mathbf{u}) \right\}. \end{split}$$

The assumptions discussed—PSA, WSA, FSA, MMCA and CMMA—provide a framework for understanding the potential performance of non-selected applicants under different scenarios of the selection process. As shown in Alarcón-Bustamante et al., 2025, these assumptions yield tighter (and therefore more informative) bounds $[\mu_V(X); \overline{\mu}_Y(X)]$ than those obtained under the WIA assumption. Using these results in (3), we evaluate these refined bounds for the regression parameters, offering insights into the predictive capacity of admission tests under varying conditions. To illustrate the practical implications of these theoretical results, we apply them to the dataset from the Chilean university admission process described above. The identification bounds will be obtained by considering linear regression models for selected applicants and multinomial models for the conditional probabilities of application and selection status.

4. Results

The Chilean university admission process incorporates test scores from mandatory tests (Language and Mathematics) and non-mandatory tests (History or Science), as well as high school performance factors. To evaluate the effect of these selection factors on GPA, a multiple linear regression model is typically employed, which includes mandatory test scores and high school performance metrics (Manzi et al., 2008). The regression model is specified as follows:

$$\mu_Y(\mathbf{x}) = E(Y \mid \mathbf{X} = \mathbf{x}) = \beta_0 + \beta_1 \text{Math} + \beta_2 \text{Lang} + \beta_3 \text{Ranking} + \beta_4 \text{GPA}_{\text{HS}}, \tag{5}$$

where

- the outcome variable *Y* represents the first-year university GPA,
- · Math and Lang denote the test scores on the mandatory Mathematics and Language tests, respectively,
- GPA_{HS} denotes the GPA at the end of high school, and
- Ranking reflects the applicant's relative position at the end of high school.

All selection factors included in the model are defined on a scale ranging from 150 to 850 points, while the outcome GPA is defined on a scale ranging from 1.0 to 7.0.

The analysis of selection factors in Chile is typically conducted using complete information from the selection factors and the selected applicants. This approach relies on the assumption of a MAR selection process, expressed as:

$$\mu_Y(\mathbf{x}) = \mu_Y(\mathbf{x}, Z = 1).$$

In our data, 56.8% (164) of the applicants were non-selected. Among the selected applicants, 47.2% were admitted to Biology and 32.8% to Biochemistry. Using this data, we assess the selection factors in the Chilean admission system by fitting the linear regression model (5) under the MAR assumption. As previously discussed, the advantage of this approach is that the regression model—and consequently, the parameters of the linear model—are point-identified. Table 1 shows the ordinary least squares (OLS) estimates (Rao, 1973) for this model.

Although the estimated coefficients in Table 1 are small, this is justified by the difference in scale between the selection factors (150–850) and the outcome GPA (1.0–7.0). Additionally, because no

9

Coefficient	β_1	β_2	β_3	β_4	
Estimation	0.00151	0.00814	-0.00092	0.00351	

Table 1. Coefficients of the regression model for each selection factor under MAR assumption

distributional assumption is considered for estimating the parameters of the regression model, no confidence intervals are displayed. Under the MAR assumption, the sign of the effect of each selection factor can be determined. Notably, most coefficients are positive, suggesting—preliminarily—a positive impact of the selection factors on first-year university GPA. Specifically, higher scores in Mathematics, Language, and GPA_{HS} are associated with better academic performance. However, this pattern does not hold for the coefficient associated with the applicant's Ranking, which does not exhibit a similar positive relationship.

The identification bounds for the regression coefficients under the different assumptions about the selection process (FSA, WSA, PSA, MMAA, CMAA, WIA) are presented in Table 2 showing the range of all possible values for the selection factor coefficients that are compatible with each assumption.

Table 2. Identification bounds (lower bound (LB) and upper bound (UB)) for regression coefficients under different beliefs about the selection process

Assumption	β_1		β ₂		β ₃		β_4	
	LB	UB	LB	UB	LB	UB	LB	UB
FSA	0.00535	0.01009	-0.00067	0.00317	-0.00175	0.00376	-0.00059	0.00633
WSA	-0.00712	0.01501	-0.00879	0.00828	-0.01745	0.01622	-0.02211	0.02090
PSA	-0.00543	0.02526	-0.00842	0.01601	-0.02129	0.02543	-0.02303	0.03727
MMAA	-0.00797	0.02746	-0.01024	0.01804	-0.02564	0.02659	-0.02653	0.04069
CMAA	-0.01864	0.03268	-0.01760	0.02289	-0.03757	0.03923	-0.04486	0.05378
WIA	-0.05420	0.05420	-0.04289	0.04289	-0.07672	0.07672	-0.10106	0.10106

From the results in Table 2, it is evident that the identification of the regression parameters varies significantly under different assumptions about the selection process. Specifically, under the WIA, PSA, MMAA, CMAA and WSA assumptions, the sign of the regression parameters for all selection factors—Mathematics, Language, Ranking, and HS-GPA—remains unidentified, as the bounds for these parameters include both positive and negative values. Consequently, under these assumptions, we cannot draw conclusions about the predictive capacity of the selection factors. However, when the Fallible Selection Assumption (FSA) is considered, the sign of the parameter associated with Mathematics is identified, as the bounds for this factor lie entirely above zero. This suggests that, if we accept the possibility of errors in the admission process, Mathematics has a positive predictive capacity over GPA in the analyzed School of Biology. Moreover, the bounds for the other selection factors under the FSA still include zero, leaving their signs unidentified and preventing us from drawing conclusions about their predictive capacity over GPA.

Results on Tables 1 and 2 are graphically displayed in Figure 1. The dashed line represents the zero line, serving as a reference to evaluate whether the sign of the parameter can be identified under each assumption. Additionally, the red dot for each selection factor indicates the estimated value of the regression parameter under the MAR assumption.

In Figure 1 we can see that, with the exception of the Ranking factor, the regression parameter estimations obtained under the MAR assumption fall outside the bounds derived under the FSA, indicating that these estimates are incompatible with the fallible assumption. These findings highlight the sensitivity of the results to the underlying assumptions about the selection process and underscore



Figure 1. Identification bounds for regression parameters under different assumptions about the selection process

the importance of carefully considering these assumptions in the analysis. The results have practical implications for decision-making in university admissions, as they demonstrate how different assumptions can lead to varying conclusions about the predictive validity of selection factors.

5. Conclusion

This study explored the predictive capacity of selection factors within the context of the Chilean admission process under the presence of missing outcomes. By proposing a linear regression model relating GPA to these factors, we examined how different assumptions about the selection process influence the identification of regression parameters and, consequently, the interpretation of predictive capacity. These assumptions include the Perfect Selection Assumption, Worst Selection Assumption and Fallible Selection Assumption proposed in Alarcón-Bustamante et al., 2025 and two additional assumptions: the Mean Most Attainable Assumption and the Conservative Most Attainable Assumption.

Without any assumptions regarding the performance of non-selected applicants, the predictive capacity of the selection factors—measured through the regression parameters—remains unidentified. Even under conservative assumptions such as the CMAA and MMAA or informative assumptions, such as the PSA and the WSA, the sign of the regression parameters for most selection factors cannot be determined due to the missing outcome data. However, under the FSA, the sign of the regression parameter for the mathematics test is identified, providing evidence of its positive predictive capacity. This finding underscores the importance of mathematics as a key predictor of academic performance in the Chilean university context under a fallible scenario.

Notably, the results obtained under the MAR assumption are compatible with those derived under the WSA, PSA, CMAA, CMAA and WIA for almost all selection factors. This suggests that, while the MAR assumption simplifies the analysis, it may not fully capture the complexities of the selection process. The partial identification approach, on the other hand, offers a more flexible and realistic framework for estimating regression parameters by imposing justifiable restrictions on unobserved values.

The derived partial identification bounds in this paper, allow to illustrate how conclusions about the predictive capacity of selection factors change under different assumptions about the selection mechanism. However, it should be noted that they are not confidence intervals as they do not account for sampling variability. To address sampling variability, Imbens and Manski, 2004; Kaido et al., 2019; Stoye, 2009 introduced confidence intervals for partially identified parameters. Illustrations of these approaches in the context of predictive capacity studies is a topic for future research.

Alternative approaches to tackle the missing data problem include statistical techniques such as Heckman's correction and related selection models, which aim to account for the effects of missing data by incorporating additional information about the selection process (Heckman, 1976, 1979; Hsu, 1995; Kennet-Cohen et al., 1999; Marchenko & Genton, 2012). Additionally, methods for correcting correlation coefficients, such as range restriction corrections, attempt to adjust for the bias introduced by selection (Koretz et al., 2016; Linn, 1983; Mendoza & Mumford, 1987; Zimmermann et al., 2017). Although these methods intent to improve the accuracy of academic performance predictions despite missing data, we did not include them in our analyses as they often rely on strong assumptions about the performance given the test scores that may not be justified in practice (Manski, 2013).

This study highlights the value of the partial identification approach in analyzing the predictive capacity of selection factors through regression parameters. By addressing the challenges posed by missing outcome data and incorporating reasonable assumptions, this approach provides a more detailed understanding of the relationship between selection factors and academic performance.

Funding Statement Inés M. Varas was partially funded by the Agencia Nacional de Investigación y Desarrollo (ANID) grant FONDECYT 11251581. Eduardo Alarcón-Bustamante was partially funded by the Postdoctoral FONDECYT grant 3220422 and by the FONDEF IDeA I+D ID22I10228 project. Jorge González was partially supported by ANID grant FONDECYT 1230968.

Competing Interests None

References

- Alarcón-Bustamante, E., González, J., Torres Irribarra, D., & San Martín, E. (2025). From missing data to informative GPA predictions: Navigating selection process beliefs with the partial identifiability approach. *British Journal of Mathematical and Statistical Psychology*, 78(2), 647–671.
- Alarcón-Bustamante, E., San Martín, E., & González, J. (2021). On the marginal effect under partitioned populations: Definition and interpretation. In M. Wiberg, D. Molenaar, J. González, U. Böckenholt, & J.-S. Kim (Eds.), Quantitative psychology. Springer International Publishing.
- Alarcón-Bustamante, E., Varas, I. M., & San Martín, E. (2023). On the impact of missing outcomes in linear regression. Chilean Journal of Statistics, 14(1), 26–36.
- Ayers, J. B., & Peters, M. (1977). Predictive validity of the test of english as foreign language for asian graduate students in engineering, chemistry, or mathematics. *Educational and Psychological Measurement*, (37), 461–46.
- Borsboom, D., Mellenbergh, G. J., & van Heerden, J. (2004). The concept of validity. Psychological Review, 111(4), 1061–1071.
- Diemer, E., Shi, J., & Swanson, S. (2024). Partial identification of the effects of sustained treatment strategies. Epidemiology.
- Embretson, S. E., & Reise, S. P. (2000). Item response theory for psychologists. Lawrence Erlbaum Associates.
- Florens, J., & Mouchart, M. (1982). A note on noncausality. Econometrica, 50(3), 583-592.
- Geiser, S., & Studley, R. (2002). UC and the SAT: Predictive Validity and Differential Impact of the SAT I and SAT II at the University of California. *Educational Assessment*, 8(1), 1–26.
- Grassau, E. (1956). Análisis estadístico de las pruebas de bachillerato. Anales de la Universidad de Chile, (102), 77-93.
- Heckman, J. (1976). The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. *The Annals of Economic and Social Measurement*, 46, 931–961.
- Heckman, J. (1979). Sample selection bias as a specification error. *Econometrica*, 47(1), 153–161.
- Hirano, K., & Imbens, G. W. (2004). The propensity score with continuous treatments. In W. A. Shewhart & S. S. Wilks (Eds.), Applied bayesian modeling and causal inference from incomplete-data perspectives: An essential journey with donald rubin's statistical family (pp. 73–84). Wiley Series in Probability; Statistics.
- Hsu, J.-W. Y. (1995). Sampling behaviour in estimating predictive validity in the context of selection and latent variable modelling: A monte carlo study. *British Journal of Mathematical and Statistical Psychology*, 48(1), 75–97.
- Imbens, G. (2000). The role of the propensity score in estimating dose-response functions. *Biometrika*, 87(3), 706–710.
- Imbens, G., & Manski, C. (2004). Confidence intervals for partially identified parameters. *Econometrica*, 72(6), 1845–1857.
- Kaido, H., Molinari, F., & Stoye, J. (2019). Confidence intervals for projections of partially identified parameters. *Econometrica*, 87(4), 1397–1432.

- Kennet-Cohen, T., Bronner, S., & Oren, C. (1999). The predictive validity of the components of the process of selection of candidates for higher education in Israel. National Institute for Testing & evaluation.
- Kolmogorov, A. N. (1950). Foundations of the theory of probability. New York: Chelsea Pub. Co.
- Koopmans, T. C. (1949). Identification problems in economic model construction. Econometrica, 17(2), 125-144.
- Koretz, D., Yu, C., Mbekeani, P. P., Langi, M., Dhaliwal, T., & Braslow, D. (2016). Predicting freshman grade point average from college admissions test scores and state high school test scores. AERA Open, 2(4), 1–13.
- Lawley, D. (1943). Iv.-a note on karl pearson's selection formulae. Proceedings of the Royal Society of Edinburgh. Section A. Mathematical and Physical Science, 62(1), 28–30.
- Linn, R. L. (1983). Pearson selection formulas: Implications for studies of predictive bias and estimates of educational effects in selected samples. *Journal of Educational Measurement*, 20(1), 1–15.
- Little, R. J. A., & Rubin, D. B. (2002). Statistical analysis with missing data (2nd ed.). John Wiley & Sons.
- Lord, F. (1980). Applications of item response theory to practical testing problems. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Makransky, G., Havmose, P., Vang, M. L., Andersen, T. E., & Nielsen, T. (2017). The predictive validity of using admissions testing and multiple mini-interviews in undergraduate university admissions. *Higher Education Research & Development*, 36(5), 1003–1016.
- Manski, C. (2016). Credible interval estimates for official statistic with survey nonresponse. Journal of Econometrics, 191, 293–301.
- Manski, C., & Pepper, J. V. (2000). Monotone instrumental variables: With an application to return to schooling. *Econometrica*, 68(4), 997–1010.
- Manski, C. (1989). Anatomy of the selection problem. The Journal of Human Resources, 24(3), 343-360.
- Manski, C. (1993). Identification problems in the social sciences. Sociological Methodology, 23, 1-56.
- Manski, C. (2003). Partial identification of probability distributions. New York: Springer.
- Manski, C. (2013). Public policy in an uncertain world: Analysis and decisions. Harvard University Press.
- Manzi, J., Bravo, D., del Pino, G., Donoso, G., Martínez, M., & Pizarro, R. (2008, July). Estudio de la validez predictiva de los factores de selección a las universidades del consejo de rectores, admisiones 2003 al 2006 (tech. rep.). Comité Técnico Asesor, Honorable Consejo de Rectores de las Universidades Chilenas.
- Manzi, J., & Carrasco, D. (2021). Validity evidence of the university admission tests in chile: Prueba de Selección Universitaria (PSU). In J. Manzi, M. R. García, & S. Taut (Eds.), Validity of educational assessments in Chile and Latin America (pp. 331–351). Springer.
- Marchenko, Y. V., & Genton, M. G. (2012). A Heckman Selection-t Model. Journal of the American Statistical Association, 107(497), 304–317.
- Mendoza, J., & Mumford, M. (1987). Corrections for attenuation and range restriction on the predictor. Journal of Educational Statistics, 12(3), 282–293.
- Messick, S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3rd, pp. 13–103). American Council on Education; Macmillan.
- Pearson. (2013, January). Evaluation of the chile psu (tech. rep.). Pearson.
- Pearson, K. (1903). Mathematical contribution to the theory of evolution-xi on the influence of natural selection on the variability and correlation of organs. *Philosophical Transactions of the Royal Society of London*, 200(Ser. A), 1–66.
- Pepper, J. (2000). The Intergenerational Transmission of Welfare Receipt: A Nonparametric Bounds Analysis. *The Review of Economics and Statistics*, 82, 472–488.
- Rao, C. R. (1973). Linear statistical inference and its applications (Vol. 2). Wiley New York.
- Rosenmbaum, P., & Rubin, D. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1), 41–55.
- San Martín, E., & González, J. (2022). A Critical View on the NEAT Equating Design: Statistical Modelling and Identifiability Problems. *Journal of Educational and Behavioral Statistics*, 47(4), 406–437.
- San Martín, E., Perticará, M., Varas, I. M., Asahi, K., & González, J. (2024). The role of identifiability in empirical research. In W. P. Fisher Jr., L. R. Pendrill, K.-D. Sommer, & T. Fröhlich (Eds.), Models, measurement, and metrology extending the si: Trust and quality assured knowledge infrastructures (pp. 133–158). De Gruyter Oldenbourg.
- San Martín, E., & Alarcón-Bustamante, E. (2022). Dissecting Chilean surveys: the case of missing outcomes. Chilean Journal of Statistics, 13(1), 17–45.
- Stoye, J. (2011). Minimax regret treatment choice with covariates or with limited validity of experiments. *Journal of Econometrics*, *166*, 138–156.
- Stoye, J. (2007). Bounds on generalized linear predictors with incomplete outcome data. Reliable Computing, 13, 293–302.
- Stoye, J. (2009). More on confidence intervals for partially identified parameters. *Econometrica*, 77(4), 1299–1315. https: //doi.org/10.3982/ECTA7347
- Tamer, E. (2010). Partial identification in econometrics. Annual Review of Economics, 2(1), 167-195.
- Thorndike, R. (1949). Personnel selection: Test and measurement techniques. New York: Wiley.
- Zimmermann, S., Klusmann, D., & Hampe, W. (2017). Correcting the predictive validity of a selection test for the effect of indirect range restriction. BMC Medical Education, 17(1), 246.