# ITALIAN GRADUATES AND INTERNATIONAL MOBILITY: A POTENTIAL OUTCOME MODEL APPLIED TO ALMALAUREA DATA

**Furio Camillo**

*University of Bologna, Italy*

**Giorgio Vittadini**

*University of Milan Bicocca, Italy; CRISP, Milan, Italy*

**Sara Binassi**

*Interuniversity Consortium AlmaLaurea, Italy*

**Abstract**. *In order to investigate and correctly measure the monetary benefits obtained when human capital migrates, we have to control for individuals' characteristics and as many aspects of their "life history" as possible, regardless of the selected statistical approach. Individuals' backgrounds and experiences will affect their propensity to migrate; as such, only by comparing similar individuals is it possible to estimate the monetary effect of choosing to work abroad after graduation. Socio-economic backgrounds, training choices, learning experiences, and more general life stories all contribute to the propensity of migrating, as well as the possible outcomes. This is a well-known bias problem generated by a self-selection mechanism. We present a method using a statistical adjustment of the self-selection problem applied to a database of life history records. The AlmaLaurea database used in this paper collects thorough data from graduates of Italian universities, including information on graduates' previous studies and social backgrounds. We analyse this data to measure the impact of the emigration choice in the context of human capital theory as applied to the modern view of a global labour market. We perform a statistical multivariate analysis to study migrants vs. non-migrants, using an innovative data-driven potential outcome approach that corrects for self-selection bias. We continue with a cluster analysis to identify homogeneous groups based on a multiple correspondence analysis and testing the self-selection bias reduction using a global imbalance measure.*

**Keywords**: *Human capital theory, migration, higher skilled, selection bias treatment, wage premium.*

## 1. HUMAN CAPITAL THEORY AND HIGHER-SKILLED MIGRATION

A large subsection of the literature on the university-to-job-market transition deals specifically with education as a potential means for ensuring greater earnings and

the key role of the human capital (HC) theory as developed by the University of Chicago (Schultz, 1959, 1961; Becker, 1962, 1964; Mincer, 1958, 1974). In this approach, the "rational choice" of investment strategies in HC is similar to that of investment in physical capital; the difference is that the effect of education and training depends on specific characteristics of each individual. Mincer's pioneering model (Mincer, 1958) used school years as an educational measure to predict annual earnings or hourly wage rates. In a subsequent extension, Mincer (1974) and other authors (Griliches, 1977; James et al. 1989; Lorenz and Wagner, 1990; Rumberger and Thomas, 1993; Cohn and Huches, 1994; Belzil and Hansen, 2002; Thomas, 2003) investigated the causes of the heterogeneous effect of education on earnings (i.e., occupation, university attended, or academic major), providing the basis for estimating a rate of return to schooling. Other authors studied the effects of self-selection and ability bias on Mincerian equations (Garen, 1984; Card, 2001; Heckman, 2000; Heckman et al., 2003, 2005, 2008), as well as general problems of the basic Mincerian wage equation (see Le et al., 2003; Oxley et al., 2008; Folloni and Vittadini, 2010).

In this context, our paper investigates the phenomenon of graduates moving abroad to work to improve their income and HC. Some papers connect the mobility of higher-skilled labour to economic development of the countries involved, business cycles, redistribution across countries' labour markets, temporary or long-term phenomenon of "brain drain" (Hatton and Jeffrey,2006; Sciulli and Signorelli 2011; Triandafyllidou 2012, 2015).

However how to measure if potential income and HC for new graduates are raised moving abroad?

HC theory gives some interesting suggestions for answering this question. HC was traditionally estimated by different methods. The retrospective method (Kendrick, 1976; Eisner, 1985) relies on a microeconomic approach, measuring the cost of "producing" a graduate for the labour market. However, this method does not consider the actual effects of HC investment on the income and wealth of households.

In the prospective method (Petty, 1690; Farr, 1853; Jorgenson and Fraumeni 1989, 1992), HC can be defined from a macroeconomic perspective as the present actuarial value of an individual's expected income weighted by the survival probability and net of maintenance costs. However, the prospective method reduces the relevance of the HC stake in education, job training and other investments and gives only macroeconomic estimations.

The educational-attainment macroeconomic method (World Bank, 1995; United Nations, 2002; Wöβmann, 2003) measures HC through educational

attainment variables (i.e., schooling, educational investment costs), non-cognitive skills, and macroeconomic educational investments such as educational infrastructures and ratio of government spending on education to GDP (Barro and Lee, 1996; Hanushek, 1996; OECD, 1998; Wöβmann, 2003). However, this method gives the same weight to each year of schooling regardless of level, the quality of educational institution, or personal cognitive skills (Wöβmann, 2003) and gives only macroeconomic estimations.

The OECD (1998) report summarizes these different methods defining HC as "the knowledge, skills, competencies, and attributes embodied in individuals that are relevant to economic activity", suggesting that it is a complex, multifaceted phenomenon not directly measurable by a set of personal attributes. This definition implies that HC must be measured at a microeconomic level using variables both concerning investment in education and monetary income. Therefore in order to evaluate if higher skilled migrants improve their economic situation and HC moving abroad it is not enough to evaluate their monetary returns. We have to control these monetary benefits for the individuals' characteristics, socio-economic background, educational history.

However, all these personal aspects not only determine the propensity to migrate but, simultaneously, influence the entity of the impact on the outcome variable, the income. In literature it has been widely demonstrated and discussed that in cases like these the impact evaluation of a treatment (migrate or not migrate) on an outcome variable is biased (Rosenbaum and Rubin, 1983). We need statistical adjustments to correctly compare groups of higher–skilled migrants and not migrants.

Given that approach, in Section 2 we describe AlmaLaurea database that contains important information on graduates' past educational paths and social backgrounds proposing some descriptive results. In Section 3 we control monetary benefits of migrants and not migrants graduates for many characteristics. We perform an innovative data-driven potential outcome statistical multivariate approach that corrects for self-selection bias the differences between migrants and not migrants. In Section 4 we use a cluster analysis to identify homogeneous groups based on a multiple correspondence analysis (MCA) and test for self-selection bias reduction using a global imbalance measure. In Section 5 we summarize the main results of our paper suggesting further remarks.

## 2. ALMALAUREA DATA COLLECTION SYSTEM AND THE DATA SET USED

Our investigation studies 2008 Italian master graduates employed in Italy and abroad. They were interviewed five years after graduation.

The data comes from a survey conducted by AlmaLaurea, an inter-university consortium founded in Italy in 1994 at the University of Bologna. The consortium began as an initiative that has currently grown to include 73 universities, as well as 91% of graduates from these universities. As such, AlmaLaurea could be considered as representative of all Italian graduates.

The AlmaLaurea project consists of three interlinked components:

1) An annual survey of graduates' profiles (based on the previous year's data) and reports on the internal effectiveness of the higher education system primarily based on administrative data. The response rate of the 2015 survey was around 90%.

2) An annual survey on graduates' occupational statuses conducted one, three, and five years after graduation with reports on the external effectiveness of the higher education system (AlmaLaurea, 2014). The response rates of the 2015 survey were 82%, 75%, and 72%, respectively, for the three groups of graduates.

3) An online CV-database with over 2,200,000 curricula. This is a tool aimed at improving the match between the supply and demand of university graduates, their international mobility, and the transparency of the labour market.

One notable feature of the annual survey is that, because of the completeness and high response rate, data on university graduates are available at a single-course level, which enables comparisons of students across institutions, departments, and so forth. The data collected and the documentation resulting from analysis of this data are provided to university governing bodies that are part of the consortium and to committees involved with teaching activities and career guidance[1].

AlmaLaurea resamples the data using a well-known Raking Ratio iterative procedure (Ardilly, 2006, pg. 283-295) that weights the sampled graduates to ensure that the data set is representative of the overall population of Italian graduates. The variables used in the re-weighting process are gender, degree subject, geographical area of university, and area of graduate-residence at the time of graduation. In order to obtain results that are as close as possible to the

---

[1] The data is also made available to all the stakeholders involved in higher education – families, students, companies and policy makers – as a solid basis for fostering all decision-making processes, activity planning, and the transparency of the labour market.

composition of the population, in the present work all the analyses have been carried out considering the weighting system.

While we focus on those who graduated in 2008, we note a significant increase in the number of graduates working abroad registered after 2008 – the year the economic crisis began – especially for graduates in the first and third years subsequent to attaining a degree.

## 2.1 DESCRIPTIVE RESULTS

Table 1 and Table 2 present descriptive statistics for graduates. Five years after graduation, 5.3% of 2008 graduates were employed abroad. Migrants are mainly male (55%) and come from the fields of engineering, political and social sciences, and economics/statistics. Graduates from the fields of education, medicine, agriculture and veterinary science, psychology, and law are less likely to migrate.

**Tab. 1: Master Graduates of the 2008 Cohort Interviewed Five Years after Graduation: Main Descriptive Statistics of Migrants vs. Non-migrants**

| | | Migrants (n = 1,484) | Non-migrants (n = 28,052) | Total (n = 29,536) |
|---|---|---|---|---|
| Gender (%) | Male | 54,6 | 39,8 | 40,6 |
| Degree subject grouping (%) [*] | Agriculture, veterinary | 1,8 | 2,5 | 2,5 |
| | Architecture | 5,6 | 5,6 | 5,6 |
| | Chemistry, pharmacy | 4,2 | 4,9 | 4,9 |
| | Economics, statistics | 11,1 | 15,2 | 15,0 |
| | Geology, biology, geography | 7,3 | 4,3 | 4,4 |
| | Law | 3,2 | 9,5 | 9,2 |
| | Engineering | 25,9 | 14,9 | 15,5 |
| | Education | 0,4 | 7,2 | 6,9 |
| | Humanities | 5,7 | 6,0 | 6,0 |
| | Foreign languages | 7,4 | 3,1 | 3,3 |
| | Medicine | 1,8 | 5,5 | 5,3 |
| | Politics, social sciences | 15,8 | 11,3 | 11,5 |
| | Psychology | 2,2 | 7,2 | 6,9 |
| | Mathematics, physics, natural sciences | 7,4 | 2,7 | 2,9 |
| Residential geographic area (%) | North of Italy | 54,0 | 47,2 | 47,5 |
| Educational qualification of mother (%) | Higher education degree | 27,8 | 17,0 | 17,6 |
| Educational qualification of father (%) | Higher education degree | 31,6 | 19,8 | 20,4 |
| Earned study abroad experiences (%) | Erasmus or other European Union programme experiences | 25,7 | 6,7 | 7,7 |
| Prepared a significant part of dissertation abroad (%) | | 29,9 | 6,4 | 7,6 |
| Current job: work geographic area (%) | North | - | 54,3 | 51,6 |
| | Center | - | 24,1 | 22,9 |
| | South and Isles | - | 21,5 | 20,5 |
| | Abroad | 100,0 | - | 5,0 |
| Current job: effectiveness of the degree (%) | Very effective/effective | 63,9 | 61,6 | 61,8 |

[*] Defence and security and Physical education are not considered
**Data source**: AlmaLaurea.

**Tab. 2: Master Graduates of the 2008 Cohort Interviewed Five Years after Graduation:
Main Continuous Descriptive Statistics of Migrants vs. Non-migrants**

|  | Migrants (n = 1,484) | | Non-migrants (n = 28,052) | | Total (n = 29,536) | |
| --- | --- | --- | --- | --- | --- | --- |
|  | Mean | Standard deviation | Mean | Standard deviation | Mean | Standard deviation |
| Age at graduation [1] | 25.7 | 2.2 | 27.3 | 5.3 | 27.2 | 5.2 |
| Graduation Mark [2] | 109.3 | 5.3 | 107.8 | 6.1 | 107.9 | 6.1 |
| Delay in degree completion times (averages, in years) | 0.3 | 0.7 | 0.4 | 1.0 | 0.4 | 1.0 |
| Average monthly net earnings (in euro) | 2,225 | 807.6 | 1,310 | 533.0 | 1,356 | 585.6 |
| Satisfaction with current job (average value on a 1 to 10 scale) | 8.0 | 1.6 | 7.5 | 1.8 | 7.6 | 1.7 |

[1]   Calculated on the basis of the age, which is considered as an entire number, the date of birth,
       and the graduation date.

[2]   Expressed on a scale of 110/110. For calculating the average marks, it has been established
       that the mark 110/110 with honours corresponds to 113/110.

**Data source**: AlmaLaurea.

More than half of the graduates who moved abroad for work reasons come from the north of Italy. Almost all of the graduates relocated to other countries within Europe (mainly the U.K., France, and Germany), with a small percentage moving to North America and other continents.

In term of educational careers, migrants seem to be more successful. They reach higher degree marks with 109.2/110 points on average compared to 108 among non-migrants, and complete their degrees before non-migrants (25.7 years old vs. 27.3 for non-migrants). The average duration of studies are lower than the levels registered among non-migrants.

Migrants participated in more study-abroad experiences, such as Erasmus Mundus or other EU mobility programmes, during their time in school. A student's willingness to go abroad, participating in foreign study programs, or completing a thesis abroad is a strong indicator for the likelihood of finding employment outside of Italy.

Other significant factors that influence successful professional establishment abroad are social class and parental education. The data shows that migrants have parents who are more likely to hold degrees, and their average social state is higher than non-migrants.

The vast majority of migrants are satisfied with the decision to leave Italy less than 2% regret their decision, and over 80% indicate high satisfaction. More than 70% of respondents consider themselves unlikely to return to Italy. The difference in average net monthly earnings after five years is significant: €1,310 for graduates who work in Italy as compared to € 2,225 for those who work abroad, and this gap is even larger for graduates in medicine and engineering.

The high rates of Italian graduates leaving Italy seem to contradict a cliché

often repeated in the Italian media: that is, the Italian university system is of poor quality. In fact, more than half of the respondents who decide to work abroad consider their educational preparation to have been very effective (64% vs. 62% of non-migrants). In addition, those who go abroad want to continue professional training: 74% of graduates who migrate participated in at least one postgraduate program.

The descriptive statistics presented here make an important contribution to evaluating the amount of HC generated by Italian universities and suggest that the problem lies with the Italian labour market, which is unable to fully recognise and enhance the HC value generated by Italian universities. The large wage gap between those who work in Italy and those who find employment abroad indicates that a standardised underpayment exists for new hires from university, which encourages graduates to leave the country.

This evidence suggests that the idea that the brain drain generated by low levels of financial rewards for Italian labour-market graduates is not in accordance with HC theory. Young Italians are responding to this dilemma with entirely rational behaviour seeking better pay elsewhere even when it means leaving their families and home country.

## 3. MODELLING APPROACH

To investigate and correctly measure the monetary benefits obtained by a graduate when he or she migrates using HC methods, we must control for the personal characteristics of each individual, using as rich a data set as possible.

In order to overcome self-selection problems, we need an approach that minimises bias from selection on observables. We add two features to a standard approach previously described in D'Attoma (2009) and Camillo and D'Attoma (2010) and initially applied to a real-life example in Peck, Camillo and D'Attoma (2010). The first distinct feature is our use of a cluster analysis, which enhances the possibility of finding local spaces of the multidimensional data-space in which covariate variables are balanced, according to the potential outcome approach (Rubin, 2005).

A second feature is our use of non-parametric methods, such as the Global Imbalance (GI) method (D'Attoma and Camillo, 2011), to base statistical adjustments of the self-selection problem on previously observed evidence. The GI method collects data from a set of key points in the lives of previously observed individuals, starting with basic characteristics about family structure, social class, and behaviours.

It is based on the concept of inertia as a measure of dependence between categorical covariates and the treatment-assignment indicator. Whereas it is common practice to assess the extent to which comparison and treatment groups are matched variable by variable, the GI measure allows for an overall assessment of how well individual cases match. For example, while the "matching technique" considers one-by-one differences between the two groups, the GI measure assesses the comparability between groups, taking into account the variation in all baseline covariates simultaneously.

The GI method adopts a three-step approach for estimating treatment effects in non-experimental data in which bias from selection on observables is minimised as follows:

1) The balance is measured and tested. The GI measurement is calculated on the whole sample; its statistical significance is then tested to isolate the statistical effects on outcomes and the net of the person's experiences.

2) In cases of an imbalance – which is probable in all non-experimental data – it is necessary to move from the global predictor space to local predictor spaces. Intuitively, this is done by means of a cluster analysis to estimate treatment effects in which bias due to selection on observables is minimised. This approach can be considered a general subgroup analysis in which the primary goal is to obtain a cluster partition that generates balanced groups thereby minimising selection bias in generating impact estimates.

The GI performs a subgroup analysis in which a multiple correspondence analysis (MCA) is used to obtain a continuous and a low-dimensional representation of the X-space generated by covariates (conditioning variables). The MCA coordinates are then entered into a cluster analysis to identify homogeneous groups.

3) It is necessary to test the balance within each cluster and to compute the local treatment effects within the balanced clusters, pruning observations in unbalanced clusters.

This method can also be used to separate classes of graduates with different results for emigration that were not a-priori observable. One can divide the graduates into homogeneous classes with respect to the migration results on the basis of other probabilities for each class.

Finally, because we only study the effects of migration in the initial post-degree years, it is not necessary to utilise the actualised lifecycle income, as we are assuming that migration is only temporary.

### 3.1 THE GI MEASURE

Camillo and D'Attoma (2009, 2010) reported that the between-group inertia of a set of units leads to a GI measure that can be expressed as:

$$GI = I_b = \frac{1}{Q} \sum_{t=1}^{T} \sum_{j=1}^{J_Q} \frac{b_{tj}^2}{k_{t.} k_{.j}} - 1 \tag{1}$$

Where

$Q$   denotes the number of baseline covariates introduced in the analysis;

$T$   denotes the number of treatment levels;

$J_Q$ denotes the set of all categories of the Q variables considered;

$b_{tj}$ is the number of units with categoryn $j \in J_Q$ in the treatment group $t \in T$;

$k_{t.}$ is the group size $t \in T$;

$k_{.j}$ is the number of units with category $j \in J_Q$.

The GI measure results from using a conditional MCA (Escofier et al., 1988) to quantify between-group inertia. This measure originates from noting that the dependence between X (covariates matrix) and T (treatment vector) is not directly observable and that displaying the relationship among them on a factorial space represents a first step to discover the hidden relationship. If dependence between X and T exists, any descriptive factorial analysis may exhibit this link.

MCA is a conventional method often used to deal with the factorial variance decomposition related to a juxtaposition of the X covariance matrix and the T treatment vector. Lebart, Morineau, and Warwick (1984) give a comprehensive description of this method, as well as computational details and applications. Camillo and D'Attoma (2010) and D'Attoma and Camillo (2011) are useful references for cases in which a conditioning variable is present.

Given that the variability (inertia) of the covariance matrix X can be decomposed into eigenvalues and eigenvectors, the presence of a conditioning treatment vector T would strongly influence the structure of the matrix-decomposition process.

Hence, a conditional analysis can help to isolate the part of the variability of the X-space that is caused by the assignment mechanism. Escofier and Pagés (1988) first studied the conditioning applied to problems arising from the dependence between categorical covariates and an external categorical variable (1988) in order to produce a method of conditional multiple-correspondence analysis (MCA_cond).

Using the terminology of Huygens–Steiner theorem (Haas, 1928), in which total inertia (I_T) can be decomposed into within-groups inertia (I_W) and between-groups inertia (I_B), MCA_cond produces a factorial decomposition of

the within-group inertia. Furthermore, MCA_cond could also be considered an intra-analysis since the inertia resulting from the conditioning variable (T) is not taken into account. Specifically, an inter-group analysis considers the relative position of groups, whereas an intra-group analysis detects and describes differences between units within each group independent of the effect caused by the partition's structure. In the context of evaluation, this structure is induced by the non-random selection mechanism. An intra-analysis allows the influence of conditioning to be measured, and this produces, as reported in Camillo and D'Attoma (2010), a comparative measure between treatment groups.

This method is particularly useful in the presence of categorical covariates, and as reported in Wermuth and Cox (1998), background knowledge tends to be qualitative in the social sciences so a frequent need exists to work with categorical covariates. It is also common for continuous variables to be split into classes in databases that are accessible to researchers.

The key result of using MCA_cond is represented by the quantified "between-group inertia" (I_B). The no omitted variable bias assumption underlying the approach then assumes a crucial role and must be emphasised. The assignment mechanism is assumed to be known, which means that the X matrix includes all baseline variables associated with both the treatment assignment and the observed outcome.

### 3.2 THE IMBALANCE TEST

We perform an imbalance test to determine the significance of any detected imbalance. We specify the null hypothesis of no dependence between X and T as:

$$\mathbf{H_0 : I_W = I_T} \tag{2}$$

where $I_w$ is the intragroup inertia and $I_t$ the total inertia. To establish an interval of plausible values for $\mathbf{I_B}$ under the null hypothesis, we use results obtained by Daunis-i-Estadella, Aluja-Banet, and Thiò-Henestrosa (2005), who studied the asymptotic distribution function of the intergroup inertia $\mathbf{I_B}$

$$\mathbf{I_B} \approx \frac{\chi^2_{(T-1)(J_Q-1)}}{nQ} . \tag{3}$$

Once the above distribution of the between-group inertia is derived, the interval of plausible values for GI is defined as:

$$GI \in \left( 0, \frac{\chi^2_{(T-1)(J_Q-1),\alpha}}{nQ} \right).$$     (4)

Specifically, if the GI calculated on the specific data set is outside the interval, the null hypothesis of no dependence between X and T is then rejected and the data are deemed unbalanced. Simulation results show that unbiased estimates of the average treatment effect (ATE) are obtained (Camillo and D'Attoma 2009, 2010) when the test detects balance.

The GI method is particularly useful in our case in which we implement a semi-automatic system on a large database to monitor the performance of Italian graduates in the job market. In this type of data system, we expect to examine many subsets of matched treatment and comparative cases. A variable-by-variable assessment of the balance of these cases is likely to not only be tedious, but also to reveal differences that exist purely by chance – some of which may be real and some of which may be random – with an unknown distinction between them.

The AlmaLaurea database structure allows the reconstruction of life stories of those interviewed with approximately 90 covariates, generated by the recognition of about 30 logical constructs and behavioural statements and opinions before going abroad. A portion of these covariates refers to constructs generated after the decision to work abroad, mostly describing the current work of the graduates (Tables 1 and 2).

Data show the existence of a conditioning-selection mechanism among those who have decided to emigrate and those who stay to work in Italy. Through a set of discriminant models and estimates on 100 sub-random samples among non-migrants, the average confusion matrix and the distribution of the average correct classification rate can be constructed. Confusion, or error, matrix is a 2x2 table layout to visualize an algorithm performance by crossing the criterion variable with the frequency distribution estimated through that algorithm. Each row of the matrix represents the instances in a predicted class, while each column represents the instances in an actual class (or vice versa). We obtain this information by applying a discriminant analysis to coordinates of a multiple correspondence analysis on pre-treatment covariates according to the Disqual approach (Saporta, 1977).

The average confusion matrix shows that the average correct classification rate (78.2%) is high enough to exclude that, conditional on the available covariates (the life stories), the decision to go abroad is not random.

In order to complete the abovementioned three-step analysis, we first compute

the GI measure for the entire sample. The resulting value of 0.0513 can be interpreted as a sign of the presence of imbalance in the data for any confidence interval between 0.001 and 0.10. This is unsurprising given the results of the discriminant models previously described in which it was possible to predict the treatment assignment using available covariates.

The second step involves a cluster analysis to identify homogeneous groups using MCA coordinates. The MCA has been implemented on all the variables available in the AlmaLaurea database that represent the life history of respondents: socio-demographic variables, variables that describe the individual's training path, and variables that constitute the respondents' socio-attitudinal profiles. Some of these variables are shown in Tables 1 and 2. The results are a set of new variables (factorial coordinates) that are continuous and orthogonal to one another. Based on these new MCA coordinates, we perform a cluster analysis to generate homogeneous groups. There is theoretical justification to expect that the variables used in the MCA are useful predictors of working abroad. While others have used cluster analyses to understand how various treated subgroups are impacted (e.g., Peck 2005; Yoshikawa, Rosman, and Hsueh 2001), our approach differs for two reasons: (1) the MCA application in creating cluster-based, treatment-comparison group impact analyses; and (2) the introduction of the GI measure to conclude whether a particular cluster has balanced treatment-comparison cases. In accordance with the tandem strategy applied in previous research (Peck, Camillo and D'Attoma 2010), we also use a Ward algorithm on the MCA coordinates to generate clusters.

A tree diagram for documenting the clustering process, known as a dendrogram, can illustrate the result of the Ward clustering method. Based on the structure of the dendrogram, we primarily focus on the 8-, 11-, 18-, 22-, 34- and 41-cluster solutions. We retain the 11-cluster solution because it provides balance within a suitable number of clusters compared to other examined cluster solutions. With the 11-cluster solution, we measure and test the balance within each group using the GI measure, which provides a clear advantage over the more traditional variable-by-variable assessment of group similarity. Using the GI approach, the number of clusters should be chosen by maximizing the generation of the so-called "common support" through which it is possible to compare treated and non-treated, but also by minimizing the possibility of non-balancing within the groups, in order to define the greatest number of equivalent groups.

Individuals within each cluster are probabilistically equivalent (and thus comparable) between treated (graduates working abroad) and untreated (graduates remaining in Italy to work). Table 3 shows the results of this cluster analysis in terms of balance, including the number of treatment and comparison cases in each cluster.

The only group defined as not balanced by the GI test is Group 4. Group 9 and Group 10, that are close to the confidence limit of the GI test, are considered balanced.

Table 4 shows computations regarding the local group effects measured using

**Tab. 3: Net Monthly Earnings: Selection Bias Cleaning**

| equivalent groups | N_treated | N_untreated | Interval for GI (alpha=0.01) | GI index | balanced |
|---|---|---|---|---|---|
| group1 | 379 | 4,836 | (0;0.11) | 0.9000 | y |
| group2 | 234 | 6,944 | (0;0.36) | 0.2200 | y |
| group3 | 112 | 984 | (0;0.03) | 0.0010 | y |
| group4 | 6 | 413 | (0;0.17) | 0.1900 | n |
| group5 | 28 | 2,052 | (0;0.08) | 0.0030 | y |
| group6 | 68 | 2,499 | (0;0.10) | 0.0080 | y |
| group7 | 18 | 784 | (0;0.11) | 0.1000 | y |
| group8 | 89 | 3,426 | (0;0.13) | 0.0600 | y |
| group9 | 237 | 1,131 | (0;0.06) | 0.0650 | y (very close to the confidence limit) |
| group10 | 297 | 1,679 | (0;0.07) | 0.0780 | y (very close to the confidence limit) |
| group11 | 67 | 2,409 | (0;0.008) | 0.0001 | y |

the real values of income. Groups are listed according to the local estimated effect in absolute value (delta real_income). In determining such effects, we used an adjusted t-test to signal groups with no significant difference of real income between treated (abroad) and non-treated (no-abroad) workers. The final average effect –1,110.80 real euros per month, the net of the local cost of living – is the result of different effects across equivalent groups. The local average effect ranges from 1,501.60 to 794.6 real euros per month.

**Tab. 4: Net Monthly Earnings in Real Terms: Comparison with Monetary Values in Equivalent Groups**

**ATE (Average Treatment Effect, weighted using only treated)**

| equivalent groups | adj_tValue | p-value | N_treated | Mean_treated | N_untreated | Mean_untreated | ATE in real EURO | ATE % |
|---|---|---|---|---|---|---|---|---|
| group3 | 4.46 | 0.0009480 | 112 | 3547.1 | 984 | 2040.5 | 1,501.6 | 73.6 |
| group4 | 1.63 | 0.1632800 | 6 | 3275.1 | 413 | 1817.5 | 1,457.7 | 80.2 |
| group9 | 2.57 | 0.0000000 | 237 | 3337.9 | 1,131 | 2094.5 | 1,243.4 | 59.4 |
| group11 | 6.19 | 0.0000000 | 67 | 3239.7 | 2,409 | 2096.4 | 1,143.3 | 54.5 |
| group1 | 16.9 | 0.0000000 | 379 | 3194.9 | 4,836 | 2088.9 | 1,106.1 | 52.9 |
| group10 | 3.63 | 0.0000000 | 297 | 3138.3 | 1,679 | 2039.6 | 1,098.7 | 53.9 |
| group6 | 6.36 | 0.0000000 | 68 | 2814.9 | 2,499 | 1826.3 | 988.7 | 54.1 |
| group7 | 3.46 | 0.0029600 | 18 | 2788.3 | 784 | 1803 | 985.3 | 54.6 |
| group2 | 0.74 | 0.0000000 | 234 | 2984.3 | 6,944 | 2010 | 974.3 | 48.5 |
| group8 | 6.76 | 0.0000000 | 89 | 2626.2 | 3,426 | 1771.3 | 854.9 | 48.3 |
| group5 | 2.95 | 0.0063890 | 28 | 2473.8 | 2,052 | 1679.2 | 794.6 | 47.3 |
| Total | | | 1,535 | 3133.8 | 27,157 | 1957.2 | 1,110.8 | 54.8 |

## 4. EQUIVALENT GROUPS DESCRIPTION

As policy and program evaluators begin to more systematically consider the heterogeneity of effects that accrue across various subpopulations, methods such as the cluster-based approach presented here will become more widely used (Peck, Camillo and D'Attoma, 2010). A better understanding of treatment effects by subgroup can help to improve policy and program targeting, making interventions both more effective and efficient. In particular, the welfare policy arena might consider the results of such data-driven approaches to research and measure impacts using a modern system of big data generated by workers and their decisional systems. The analysis of this type of data would allow to segment the motivations of people to migrate, starting from their stories and their training paths and considering their initial socio-economic status.

In our case, it therefore seems difficult to trace the determinants to migrate for Italian graduates to few objectively detectable variables, but rather the personal paths of sedimentation of the experiences made seem decisive. If our policy makers wish to manage and influence this propensity, they have to improve the quality of the Italian labour market in order to reduce migration of the best graduates.

With this in mind, it is helpful to provide an "equivalent group" description, using all available information regarding units, and taking the nature of generated clusters into account. As such, we describe, in descending rank according to the absolute value of the effect on income (ATE in real Euro in Table 4), the main statistical significant characteristics of each equivalent group. Please note that the average general ATE is 1100.8 euros.

The quality of the achieved higher education is discriminant between migrating and non-migrating people. Hence, final mark and distinctive discipline(s) are to be considered as distinctive traits of each group (provided the cluster analytic method works), in particular for Group 3 which is the group with the largest absolute and relative income.

Group 3. (ATE in real Euro=1501.6) The graduates of this group are open to work in any kind of company. Their average age is higher than that of the other graduates. They come primarily from southern Italy. Their parents have limited educational background. Their academic performances were good but not excellent, and during college they held temporary jobs. After graduation, their most common occupation is teacher in technical high school.

Group 4. (ATE in real Euro=1457.7) This group was created without balancing because it is composed of a number of students not sufficiently high.

Group 9. (ATE in real Euro=1243.4) This group is comprised mostly of subjects who have completed their final thesis abroad. They are students involved

in university life, except for international mobility programs such as for example Erasmus. Their family educational background is medium/high and their academic performances had some minimal delays. Compared to the average graduate, they highly value jobs that guarantee stability both in terms of their employment contracts and protection from economic cycles. During their studies, these students worked sporadically.

Group 11. (ATE in real Euro=1143.3) The main characteristic of this group is that the subjects graduated at a late age. This group held jobs outside of school during their studies, and they had poor academic performance. Their social class is medium-low. A significant number of these graduates studied Medicine, and a few pursued degrees in Chemistry-Pharmacology.

Group 1. (ATE in real Euro=1106.1) This group consists of very young students from Italian universities in the northwest, who spent time studying abroad during college. These students graduated on time and likely with honours, after having performed at a high level in high school. Students in this group stated very strongly that they would accept a job in the north of Italy or in Europe. Their social class is generally not high, and their parents have modest backgrounds.

Group 10. (ATE in real Euro=1098.7) These young people completed tests abroad and studied abroad. Their father and probably their mother earned a university degree, and they worked outside of school occasionally during their studies. It is very probable that they attended high school abroad. A large proportion of this group are engineers.

Group 6. (ATE in real Euro=988.7) This is a group without any experience abroad during their studies. The group includes many architects, and in particular, people who completed degree programs in Rome. The academic performance of this group, assessed at intermediate tests or final exams, is not particularly high. Their fathers generally hold a college degree, are often self-employed or work in management. Their mothers likely graduated from college and have high-paying jobs.

Group 7. (ATE in real Euro=985.3) The students in this group don't have any study or work experience abroad. University of Perugia graduates and those who pursued a degree in Agriculture comprise a large proportion of this group, although a significant number of students who studied Political or Social Sciences is also included here. Generally, this group represents youth of central Italy, who currently work in consulting areas or electronics industries. These young people have an above-average social background and parents with university degrees.

Group 2. (ATE in real Euro=974.3) They are graduates from North-Eastern universities, such as the University of Padua. The students in this group had no experiences abroad, and good but not excellent academic performances. Enginee-

ring is the most prevalent discipline. These students have a technical high school diploma obtained with very high grades. In this group, the average parental education level is middle school. Their parents generally have low-paying jobs.

Group 8. (ATE in real Euro=854.9) The key characteristic of this group is that graduates are primarily from universities in the centre of Italy and in particular from schools in Rome and Florence. These students have high academic performances, a high school diploma, graduated on time, and did not travel abroad during their studies. Students in this group pursued degrees mostly in the Humanities, and, less frequently, in Architecture and Economics-Statistics. Members of this group are primarily women, and tend to have part-time jobs outside of school during their studies.

Group 5. (ATE in real Euro=794.6) Students of this group did not work during their studies nor studied abroad. A very important feature of this group seems to be the students' provenance: although these are young graduates from southern Italy who attended universities in the south or in central Italy, they most likely lived away from their family. Output expectations are another interesting element, because this group, in its formative path, has really achieved a great professional growth. Their parents are generally middle class.


## 5. CONCLUSION

Using the potential outcomes approach, we performed a comparison between the treatment group, migrant graduates, and the control group, non-migrants. The isolated covariates demonstrated the effects on outcomes and net of the person's background and experiences.

We used a non-parametric method of GI that takes into account variation in all baseline covariates, simultaneously minimises bias from selection on observable variability, and corrects imbalance problems among groups by means of transition-generating balanced groups.

Overall, the GI method not only allowed a general comparison among migrant-graduates and non-migrant graduates, it also facilitated division into homogeneous classes on the basis of conditional probabilities for a given class. We utilized this method to analyse the migration of Italian graduates.

In the past several years, the percentage of Italian graduates who work abroad has increased dramatically; emigration is no longer limited to southern Italians or graduates of low social and economic classes. This "brain drain" is mainly due to the fact that in the first five years after graduation migrating graduates income in real values is greater than the income of graduates working in Italy, although this

difference varies across destination countries.

Moreover, the average income, as well as the variance of income, is greater abroad than in Italy not only in general, but also for each cluster of graduates. We argue that that escaping from Italy is a rational choice of Italian graduates. In fact, usually, to work abroad when their qualifications and competencies are more highly rewarded is the final choice after some "tasting" of the Italian labour market. In this sense their decision is comparative, hence rational. But this decision, sometimes far from the origin family, and in some cases also from the seeds of own family, is painful, a wound. The average income between different clusters differs greatly. The level of education and parental socio-economic status, the kind of institution attended to earn a degree, and individual characteristics all affect graduates' income.

It would be interesting to apply the GI method to graduates' outcomes of different countries to see if our results persist globally. It would also be helpful to have administrative data similar to that prepared by the AlmaLaurea databases to carry out comparative studies in various European and worldwide university systems.

Our methodological suggestions are based on the Neyman potential outcomes framework and the assignment mechanism: every unit has different potential outcomes depending on their "assignment" to a condition.

To measure the causal effect of treatment 1 versus treatment 2, the investigator should look at the outcome for the same individual in both alternative futures. Since it is impossible to see both potential outcomes at once, one of the potential outcomes is always missing.

The approach used in this paper does not raise the problem of using covariates to describe cause and effect relationships, but above all, they are used to place individuals "on equal terms".

The granularity with which the categories of qualitative variables used as "covariates" were then treated makes it possible to evaluate better and more accurately outcome differences.

One of main contribution of our research is to have tested and then proposed an innovative approach to the extent of the impact of going to work abroad for Italian graduates. This approach is completely data-driven and above all does not require the specification and estimation of any model, if not subsequently within each "equivalent group", as well as the modern causal inference based on the potential outcome approach in order to show how the groups themselves are different. The fact of being data-driven makes it possible to implement this type of study as a semi-automatic process of monitoring causal impacts over time.

A longitudinal analysis could be fruitful for understanding the perspectives of migrants and in that case a proper GI method is to be devised.

## REFERENCES

AlmaLaurea (2014). *The 16th Survey on Graduates' Employment Conditions.* Bologna, Il Mulino. Available on http://www.almalaurea.it/en/universita/occupazione/occupazione12.

Ardilly, P. (2006). *Les techniques de sondage.* Editions TECHNIP, Paris.

Barro, R. and Lee, J.W. (1996). International measures of schooling years and schooling quality. *American Economic Review*. 86(2): 218-223.

Becker, G.S. (1962). Investment in human capital: A theoretical analysis. *The Journal of Political Economy.* 70: 9-49.

Becker, G.S. (1964). *Human Capital. Columbia University Press for the National Bureau of Economic Research,* New York.

Belzil, C. and Hansen, J. (2002). Unobserved ability and the return to schooling. *Econometrica.* 70(5): 2075- 2091

Camillo, F. and D'Attoma, I. (2009). A multivariate approach to assess balance of categorical covariates in observational studies. In Book of short papers, Seventh Scientific Meeting of the Classification and Data Analysis Group of the Italian Statistical Society. CLEUP, Padova: 35-38.

Camillo, F. and D'Attoma, I. (2010). A new data mining approach to estimate causal effects of policy interventions. Expert Systems with Applications. 37(1): 171–181.

Card, D. (2001). Immigrant inflows, native outflows, and the local market impacts of higher immigration. *Journal of Labor Economics.* 19(1): 22-64.

Cohn, E. and Huches, W.W. (1994). A benefit-cost analysis of investment in college education in the United States: 1969-1985. *Economics of Education Review.* 13:109-123.

D'Attoma, I. and Camillo, F. (2011). A multivariate strategy to measure and test global imbalance in observational studies. Expert Systems with Applications. 3: 3451–3460.

Daunis-i-Estadella,, J., Aluja-Banet, T. and Thió-Henestrosa, S. (2005). Distribution of the inter and intra inertia in conditional MCA. Computational Statistics. 20 (3): 449-463.

Eisner, R. (1985). The total incomes system of accounts. Survey of Current Business, 65(1): 24–48.

Escofier, B. and Pages, J. (1988). Analyses factorielles simples et multiples; objectifs, méthodes et interprétation. Dunod. Paris.

Farr, W. (1853). Equitable taxation of property. *Journal of the Royal Statistical Society.* XVI: 1–45.

Folloni, G. and Vittadini, G. (2010). Human capital measurement: A survey. *Journal of Economic Surveys.* 24:248-279.

Garen, J. (1984). The returns to schooling: A selectivity bias approach with a continuous choice variable. *Econometrica.* 52(5): 1199-1218.

Griliches, Z. (1977). Estimating the returns to schooling. Some econometric problems. *Econometrica,* 45:1-22.

Haas, A.E. (1928). Introduction to Theoretical Physics. Constable & Company ltd., New York.

Hanushek, E.A. (1996). Measuring investment in education. *The Journal of Economic Perspectives* 10:9-30.

Hatton T J. and Williamson J.G. (2006). Global Migration and the World Economy: Two Centuries of Policy and Performance. MIT Press, Cambridge, (Massachusetts).

Heckman J. (2000). Policies to foster human capital. *Research in Economics.* 54:3-56.

Heckman, J., Lochner, L. and Todd, P. (2003). Fifty Years of Mincer Earnings Regressions. Technical report. *National Bureau of Economic Research,* Cambridge (Massachusetts).

Heckman, J., Lochner, L. and Todd, P. (2005). Earnings Functions, Rates of Returns and Treatment Effects: The Mincer Equation and Beyond. Working Paper Series N¡ 11544. National Bureau of Economic Research, Cambridge (Massachusetts).

Heckman, J., Lochner, L. and Todd, P. (2008). Earnings Functions and Rates of Returns. Working Paper Series N¡ 13780. National Bureau of Economic Research, Cambridge (Massachusetts).

James, E., Alsalam, N., Conaty, J.C. and To, D.L. (1989). College quality and future earnings: Where should you send your child to college? *The American Economic Review,* 79:247-252.

Jorgenson, D.W. and Fraumeni, B.M. (1989). The accumulation of human and non-human capital, 1948-84. In R. E. Lipsey and H. Stone Tice, editors, The Measurement of Saving, Investment, and Wealth. University of Chicago Press, Chicago: 227-286

Jorgenson, D.W. and Fraumeni, B.M. (1992). The output of the education sector. In Z. Griliches editor, Output Measurement in the Service Sectors. University of Chicago Press, Chicago: 303-341.

Kendrick, J.W. (1976). The Formation and Stocks of Total Capital. Columbia University Press, New York.

Le, T., Gibson, J. and Oxley, L. (2003). Cost- and income-based measures of human capital. J*ournal of Economic Surveys,* 17:271-307.

Lebart, L., Morineau, A. and Warwick, K.M. (1984) Multivariate Descriptive Statistical Analysis. John Wiley & Sons, New York.

Lorenz, W. and Wagner, J. (1990). A Note on Returns to Human Capital in the Eighties: Evidence from Twelve Countries. Luxembourg Income Study, Luxembourg.

Mincer, J.A. (1958). Investment in human capital and personal income distribution. *The Journal of Political Economy* 66:281-302.

Mincer, J.A. (1974). Schooling and earnings. In J.A. Mincer, editor, Schooling, Experience and Earnings. National Bureau of Economic Research. Cambridge (Massachusetts): 41-63.

OECD (1998). Human capital investment: An international comparison. *Economics of Education Review* 20:93-94.

Oxley, L., Le, T. and Gibson, J. (2008). Measuring human capital: Alternative method and international evidence. *Korean Economic Review.* 24:283-344.

Peck, L.R., Camillo, F. and D'Attoma, I. (2010). A promising new approach to eliminating selection bias. *The Canadian Journal of Program Evaluation.* 24(2): 31–56.

Peck, L.R. (2005). Using cluster analysis in program evaluation. *Evaluation Review.* 29(2): 178-196.

Petty, W. (1690). Political Arithmetick..Reprinted in: C.H. Hull, The Economic Writings of Sir William Petty (1899).

Rosenbaum, P.R. and Rubin, D.B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1): 41-55.

Rubin, D.B. (2005). Causal inference using potential outcomes: Design, modeling, decisions. *Journal of The American Statistical Association.* 100: 322-331.

Rumberger, R.W. and Thomas, S.L. (1993). The economic returns to college major, quality and performance: A multilevel analysis of recent graduates. *Economic of Education Review.* 12:1-19.

Saporta G. (1977). Une méthode et un programme d'analyse discriminante sur variables qualitatives, in E.Diday, editor, Analyse des Données et Informatique. INRIA, Paris: 201-210.

Schultz, T.W. (1959). Investment in man: An economist's vie. *The Social Service Review.* 33:109-117.

Schultz, T.W. (1961). Investment in human capital. *The American Economic Review.* 51(1): 1-17.

Sciulli, D. and Signorelli, M. (2011). University-to-work transitions: An empirical analysis on Perugia graduates. *European Journal of Higher Education* 1:39-65.

Thomas, S.L. (2003). Longer-term economic effects of college selectivity and control. Research in Higher Education. 44:263-299.

Triandafyllidou, A. (2012). Intra EU-Mobility and EU Citizenship in Times of Crisis. https://www.eui.eu/Projects/EUDO/Documents/powerpoints/TriandafyllidouAnna.pdf

Triandafyllidou, A. and Gropas, R. (2015). What is Europe. Palgrave, London.

United Nations (2002). Human Development Report 2000: Human Rights and Human Development. Oxford University Press, New York.

Yoshikawa, H., Rosman, E.A. and Hsueh, J. (2001). Variation in teenage mothers' experiences of child care and other components of welfare reform: Selection processes and developmental consequences. Child Development. 72(1): 299-317.

Wermuth, N. and Cox, D.R. (1998). On the application of conditional independence to ordinal data. *International Statistical Review.* 66, 181-199.

Wöβmann, L. (2003). Specifying human capital. *Journal of Economic Surveys.* 17(3): 239–270.

World Bank (1995). World Development Report. Oxford University Press, Oxford.