

IEB Working Paper 2023/13

**GENDER DIFFERENCES IN HIGH-STAKES PERFORMANCE AND
COLLEGE ADMISSION POLICIES**

Andreu Arenas, Caterina Calsamiglia

Version December 2023

Public Policies

IEB Working Paper

**GENDER DIFFERENCES IN HIGH-STAKES PERFORMANCE
AND COLLEGE ADMISSION POLICIES**

Andreu Arenas, Caterina Calsamiglia

The **Barcelona Institute of Economics (IEB)** is a research centre at the University of Barcelona (UB) which specializes in the field of applied economics. The IEB is a foundation funded by the following institutions: “La Caixa” Foundation, Saba, the Barcelona City Hall, the Barcelona Metropolitan Area, the University of Barcelona, the Autonomous University of Barcelona, the Barcelona Provincial Council, Agbar, Cuatrecasas and Consorci Zona Franca Barcelona.

The **IEB** research program in **Public Policies** aims at promoting research related to the design, implementation and evaluation of public policies that meet social needs and enhance social and individual welfare. Special emphasis is put on applied research and on work that tries to shed light on the Economics of Education, Health Economics, Innovation, Labour Markets and Security Policies. Disseminating research findings in these topics to a broader audience is also an aim of the program.

Postal Address:

Institut d’Economia de Barcelona

Facultat d’Economia i Empresa

Universitat de Barcelona

C/ John M. Keynes, 1-11

(08034) Barcelona, Spain

Tel.: + 34 93 403 46 46

ieb@ub.edu

<http://www.ieb.ub.edu>

The IEB working papers represent ongoing research that is circulated to encourage discussion and has not undergone a peer review process. Any opinions expressed here are those of the author(s) and not those of IEB.

**GENDER DIFFERENCES IN HIGH-STAKES PERFORMANCE
AND COLLEGE ADMISSION POLICIES***

Andreu Arenas, Caterina Calsamiglia

ABSTRACT: The Gale-Shapley algorithm is one of the most popular college allocation mechanism around the world. A crucial policy question in its setting is designing admission priorities for students, understanding how they disadvantage certain demographic groups, and whether these differences are related to differences in college performance potential (i.e., whether these differences are fair). Studying a policy change in Spain, we find a negative effect of increasing the weight of standardized high-stakes exams on female college admission scores, driven by students expected to be at the top. The effect on admission scores does not affect enrolment, but the percentage of female students in the most selective degrees declines, along with their career prospects. Using data on college performance of pre-reform cohorts, we find that female students most likely to lose from the reform tend to do better in college than male students expected to benefit from the reform. The results show that rewarding high-stakes performance in selection processes may come along with gender differences unrelated to the determinants of subsequent performance.

JEL Codes: J16; I23; I24

Keywords: College Admissions, High-stakes Exams, Algorithmic Fairness, Gender Gaps

Andreu Arenas
Universitat de Barcelona & IEBCaterina Calsamiglia
ICREA & IPEG & IZA

*We are grateful to seminar participants at Bar-Ilan, Bologna, EIEF, Columbia, Florida, Kadir-Has, Maryland, Oxford, Toulouse, UAB, UB, UC3M, UPF; the Catalan Economics Society Conference, EUI Alumni Conference, IZA Workshop on Economics of Education, Nuremberg conference on Gender Economics and the Workplace, Simposio of the Spanish Economic Society, SOLE, Workshop on Public Policies: Inequality of Opportunity, ZEW Mannheim Workshop on the Economics of Higher Education, the 16th Matching in Practice Workshop; and to Sule Alan, Ghazala Azmat, Juan Dolado, Andrea Ichino, Raquel Fernández, Rosa Ferrer, Naomi Friedman-Sokuler, Libertad González, Nagore Iriberrí, Hannes Mueller, Daniele Paserman, Perihan Saygin and Miguel Urquiola for helpful comments. Arenas acknowledges funding from the Spanish Ministry of Science PID2020-120359RA-I00.

1 Introduction

In many countries across the world, such as Brasil, Chile, China, Croatia, France, Germany, Hungary, Ireland, Israel, Korea, Norway, Spain, Sweden, Portugal, Russia, Turkey, or Ukraine, students are allocated to college through algorithms.¹ The most common is the Gale-Shapley algorithm (Gale and Shapley, 1962; Roth and Sotomayor, 1992), where both colleges and students are able to find the match they most prefer from among those who prefer them (i.e., stable matchings). In this process, there are two crucial inputs to the algorithm: colleges' preferences for students, and students' preferences for colleges. Colleges' preferences for students are typically summarized by an admission score which takes into account students' performance in standardized high stakes exams (taking place in a given day) as well as their high school GPA (which averages multiple exam results over a longer time period). For instance, in Australia, both school evaluation grades and final standardized exams determine college admissions, with weights ranging from 75%-25% to 50-50%. Likewise, in Chile (with weights around 60%-40%), in Ecuador (60%-40%), in Portugal (50-50%), or in Turkey (50%-50%). In Canada or Sweden, admissions are based on high school grades, complemented with standardized exams in the most selective programs. In Brasil and South Africa, admissions are fully based on standardized exams.

A crucial policy question concerns whether the weights to the inputs of the algorithm to determine admission grades disadvantage certain demographic groups, and whether these differences are related to differences in college performance potential (i.e., whether these differences are fair). A very salient case concerns gender differences. In lab experiments, men's performance tends to be more elastic to the competitiveness of the environment than women's. Examples include solving mazes in tournaments (Gneezy *et al.*, 2003) or running in a physical education class (Gneezy and Rustichini, 2004). Females are also less likely to

¹Source: matching-in-practice.eu

self-select into competitive tournaments, even after controlling for performance, confidence and risk aversion (Niederle, 2015; Niederle and Vesterlund, 2011, 2007). These gender differences in lab experiments do not necessarily relate to relevant differences in qualifications or subsequent performance: for instance, Balafoutas and Sutter (2012) and Niederle *et al.* (2013) find that affirmative action interventions encourage women to enter competitions more often, and performance is at least equally good, both during and after the competition.

The effects changing the inputs to the college allocation algorithm on gender differences in the students' allocation to college are however difficult to predict. First, because students may react to policy changes, partially offsetting any effects that we may expect based on baseline differences in grades. Second, because the effects will depend on who are the most affected students, and whether they are competing for the same programs of admission. For instance, if the most affected students are at the margin of being admitted to college, changes in the inputs to the college algorithm will lead to changes in college enrolment. Instead, if the most affected students tend to be high performing students, there will be a reshuffling of students across academic programs. Likewise, whether these differences lead to changes in college quality or field of study will depend on the interaction between students' preferences and their gains from any policy change. For instance, in a world with fully segregated gender differences in preferences for academic programmes, changes in the inputs to the college allocation algorithm would have little effect on the college allocation. To sum up, understanding and quantifying those effects is an empirical question.

In this paper, we study the effect of a policy change which increased the weight of the high-stakes standardized exam for (centralized) college admissions in Spain from 40% to 57%, using administrative data on college applications and college performance in the region of Catalonia, which hosts some of the best universities of Spain. First, we study the effect of the reform on gender differences in admission scores. Second, we quantify the effect of the reform on gender differences in college enrolment, college selectivity and career prospects. Last but

not least, we study the relationship between gender differences in high-stakes performance and college performance skills, by studying what type of students (based on their potential for college performance) are most affected by the reform.

This is an important question, because the number of students attending higher education has more than doubled in the last decades (UNESCO, 2017), and a large share of them are allocated to college through similar algorithms. The field and the institution of enrolment have been shown to have a large impact on life prospects such as earnings (Kirkebøen *et al.*, 2016), whom one marries (Kirkebøen *et al.*, 2021) and even the well-being of potential children (Kaufmann *et al.*, 2021), particularly so for women. Hence, the algorithms determining who has access to higher education and where do have a large impact on society. Furthermore, in countries with decentralized college admissions, like the US or Italy, admissions also frequently rely on high school grades or high-stakes exams.

The three main results of the paper are the following. First, we find a negative effect of increasing the weight of the high-stakes exam on female admission scores. The size of the effect is similar to the date of birth effect in our sample (i.e., the effect of being born in January rather than in December); to 15% of the parental college education gradient in admission scores in our sample; or to the effect of taking an exam in a day with high pollution (Ebenstein *et al.*, 2016). The effect is slightly larger than the effect of re-weighting high school grades (where females largely outperform males) and high-stakes grades (where there are smaller differences in performance) differently. This suggests that students' reaction to the policy is small compared to its mechanical effects, although they both go in the same direction.

Second, we study the effect of the reform on students' allocation to college. This effect depends on who are the most affected students and whether they are competing for the same academic programmes. We find no effect of the reform on college enrolment, because the effect on admission scores is driven by students expected to be top performers. This

is consistent with previous evidence finding that performance gaps at high percentiles are related to the differential manner in which men and women respond to competitive test-taking environments (Niederle and Vesterlund, 2010). However, we do find that female students become significantly less likely to attend the most selective programmes. Enrolment in programs above the median level of selectivity declines by 3pp, compared to enrolment in programs below the median. We also estimate that this change in the allocation to college leads to worse career prospects for female students: a 2pp points increase in the expected gender wage gap, on top of a pre-reform gender wage gap of 9.3% (within-field) and 20% (unconditional) four years after graduation.

Third, we study the correlation between gains from the reform and college performance skills. To this aim, we identify the types of students who are most likely to benefit from the reform, based on a large set of pre-determined covariates. Then, we study whether students predicted to win from were underplaced or overplaced before the reform (i.e., whether their college performance was above or below that of classmates with the same admission grades). Within gender, we find that students expected to win from the reform used to be underplaced: they were performing better in college than comparable students with the same admission grade and enrolled in the same college-major and pre-reform cohort. This suggests a stronger relationship between high-stakes performance and college-performance skills than between high-school performance and college-performance skills. However, across genders, the sign of this relationship is the opposite. Females predicted to lose from the reform were underplaced before the reform (i.e., doing better in college), compared to similar male students predicted to win. Hence, this suggests that the gender differences in high-stakes performance are not related to differences in college performance potential.

We make three novel contributions. First, our results emphasize that choosing different inputs for the algorithm has significantly different effects across groups and that these do not necessarily relate to group-differences in college performance potential. Our results are

relevant for a large number of countries which make use of very similar centralized college allocation mechanisms, but which differ in the weights given to high school and high-stakes GPAs, such as Brasil, Chile, China, Croatia, France, Germany, Hungary, Ireland, Israel, Korea, Norway, Sweden, Portugal, Russia, Turkey, or Ukraine.² They are also broadly relevant for designing algorithms in admission or selection processes, in educational settings and in the labor market.

Second, we contribute to the literature on gender-differences in high-stakes exams by estimating its policy consequences. Gender differences in high-stakes exams in educational competitive settings have been found in various countries (Jurajda and München, 2011; Saygin, 2018; Montolio and Taberner, 2018; Iriberry and Rey-Biel, 2019; Arenas *et al.*, 2021; Graetz and Karimi, 2022). For instance, Schlosser *et al.* (2019) and Cai *et al.* (2018) find significant gender differences in performance between mock and actual GRE and Gaokao (the Chinese college admission exam) tests. Azmat *et al.* (2016) find that throughout secondary and high school, girls always outperform boys, but especially in lower-stakes exams. Duckworth and Seligman (2006) find that self-discipline is an important driver of these differences. Ors *et al.* (2013) find that male students outperform female students in admission exams of the most selective French Business School, but not in first-year courses nor in high school. Morin (2015) finds that male average grades and the proportion of male students graduating on time in college increased relative to females within a cohort of students in Canada which was exceptionally large, which increased competition for grades.

Third, we characterise the compliers' profile and relate it to a relevant trait that policy-makers would like to select for (in this case, college performance skills). Our results directly speak to this question, and jointly evaluate the distributional and efficiency implications of policies that put more or less weight on high-stakes performance. The results suggest that changing the input to the college choice algorithm does not only have unequal admission

²Source: matching-in-practice.eu

effects across demographic groups, but also that these may not be related to subsequent performance. This complements existing studies which have shown that self-discipline is an important driver of gender differences in high school vs. high-stakes exams (Duckworth and Seligman, 2006). Graetz and Karimi (2022) document that that cognitive skills, motivation, and effort are strongly positively related to high school grades in Sweden, while SAT scores are more informative about cognitive skills, but they show no correlation with motivation or effort. Our results are compatible with these findings, for instance if self-discipline, motivation and effort are more important determinants of high school grades and college performance, and higher among females.

Overall, our paper provides new evidence on the effects of admission criteria on schools and colleges. This complements a large literature on how the type of mechanism used affects admissions (i.e., deferred vs. immediate acceptance, the unequal effects of non-strategy-proofness, reviewed in Abdulkadiroglu and Andersson (2022)) and the literature on the effects of other admission criteria in earlier education stages, such as catchment areas or families' proximity (Söderström and Uusitalo, 2010; Calsamiglia and Güell, 2018; Gortázar *et al.*, 2023). We contribute to the question of how admission criteria affect college admissions by studying their effects on gender inequalities, a question that has attracted less attention before, despite a large literature showing gender differences in performance in high-stake exams.

2 Context, policy change, and conceptual framework

The college allocation process starts with students listing their preferences in an application form. Then, they are allocated to academic programmes (i.e., pairs major \times university) based solely on their admission grades, which are a weighted average of high school grades and grades in a comprehensive high-stakes exam at the end of high school, namely the PAU

(Proves d'Accés a la Universitat), which covers the contents of high school. High school lasts for two years, and students specialize in one of five possible specialities: arts, humanities, social science, science, or technology. The high-stakes exam (PAU) includes exams on core subjects common for all high school students (namely Catalan, Spanish, English, and Philosophy or History), and on three field subjects corresponding to the students' specialization in high school. Table A23 in the Appendix C shows that the structure of each subject exam's remained very similar (in many cases, identical), mostly consisting of open ended questions and with some degree of choice between questions for all subjects. Students are then allocated into academic programmes (i.e. pairs college-major), which are capacity constrained, using a Gale-Shapley mechanism (Gale and Shapley, 1962). Every year, the admission grade of the last student admitted into an academic programme becomes public and it is known as the *threshold grade*. The allocation is managed by regions, and it follows the same standard Gale-Shapley mechanism for the slots of public universities in every region.³

Before the 2010 reform, the admission grade was computed as follows:

■ Before 2010:

$$\text{Admission Grade} = \frac{60 \times \text{High School GPA} + 40 \times \text{high-stakes GPA}}{100}$$

The aim of the reform was to give more points to field subjects that are related to the degrees students' want to pursue. However, the way this was implemented ended up giving a lot more weight to the high stakes exam overall.⁴ After the 2010 reform, the admission grade was computed as follows:

³Every student has to fill out an application form for every region where she is applying to college.

⁴The pre-reform weights were set in 2000 (before that, they were 50%-50%), to give more significance to the *years of hard work* and minimize that students would put everything at risk in a single test.

■ After 2010:

$$\text{Admission Grade} = \frac{60 \times (\text{High School GPA}) + 40 \times (\text{high-stakes GPA, Core} + \text{Field Subject A})}{140} \\ + \frac{W_B \times (\text{high-stakes GPA, Field Subject B}) + W_C \times (\text{high-stakes GPA, Field Subject C})}{140}$$

Where W_B , W_C can be 10 or 20 depending on the subject relevance for the degree where the student is applying. These weights could be zero if a student does not take the subject's exam, which is only possible for these subjects (taking the rest, including field subject A, is compulsory). Students may not take the exam if they prefer to relocate their effort to some of the other exams.⁵ This means that the post-reform high-stakes exam amounts to up to $\frac{80}{140} \approx 57\%$ of the admission grade, a substantial increase from the pre-reform weight (40%).

Besides increasing the weight of the high-stakes exam, the reform comes along with two additional relevant changes, which we will also study to understand whether they could be confounding any effects driven by the change in the weight of the high school versus the high-stakes exam. First, there are changes in the relative weights of subjects within the high-stakes exam GPA. Indeed, the main reason for the reform was to increase the weight of field subjects for college admissions. However, this was done in a way that led to a quite large change in the overall weight of the high-stakes exam compared to the high school GPA. Field subjects account for up to 60% of the high-stakes GPA after the reform, compared to 50% before the reform.⁶ If there are systematic gender differences in performance in field vs. core subjects, this could have an effect on admission grades beyond the change in the weights of high school and high-stakes GPAs.

Second, after the reform, the weight of two field subjects may change depending on

⁵58.7% of females and 58.2% of males take both; 27.1% of females and 26.2% males take only one; 14.1% of females and 15.6% of males take none.

⁶Before the reform, each core subject in the high-stakes exam counted for 12.5%; one field subject for 10%, and two field subjects for 20%. After the reform, each core subject counts for 10%, one field subject for 10%, and two field subjects for up to 25%.

whether the student is being considered for enrolment in a related field. In practice, because students tend to apply and enrol into programmes related to their high school studies, W_B and W_C are on average 19, conditional on taking the field exams. Table A1 in the Appendix reports the distribution of weights by gender, and table A2 shows that there are no significant gender differences neither in taking field subjects exams nor on the average weights. As a benchmark, we use the admission grade with W_B and W_C from students' program of enrolment, but as a robustness check, we will also provide estimates using $W_B = W_C = 19$ for all students, as well as controlling for them.

Overall, conceptually, these are the relevant variables, functions, and policy aims:

- Variables:

- School performance skills: the skills and knowledge acquired by students during their high school education
- High-stakes exam performance skills: the skills and knowledge tested by the high-stakes exam taken at the end of high school
- College performance skills: the skills and knowledge required for success in college
- Admission grades: grade that determines students' ability of enrolling into the academic program of their choice, which is based on a combination of their school performance and high-stakes performance skills
- Students' allocation to college: student enrolment into different academic programmes and universities
- Demographic divisions: differences in the previous variables by gender

- Functions:

- Admission score formula: a weighted average of high school grades and PAU exam scores, used to determine a student's admission grade

- Choice algorithm: a process for allocating students to academic programs based on their admission scores and program preferences
- Objective function of the policy maker:
 - To allocate students to academic programs in a way that gives the highest priority to students with the highest college performance skills, while also promoting equal opportunity for all students.

This paper studies (1) how a change in the admission score formula changed admission grades by gender (which significantly differ in their school vs. high-stakes performance skills), and (2) subsequently how this affected the students' allocation to college. It furthermore investigates (3) whether the gender difference in high-school vs. high-stakes performance is related to gender differences in college performance skills.

2.1 Data

The main data source for this paper consists of administrative records on enrolment applications to public universities in Catalonia, a large region of Spain with some of the best universities in the country (for instance, according to the 2018 Times Higher Education Ranking, five out of the seven best Spanish Universities are in Catalonia).⁷ Cross-region mobility for undergraduate studies in Spain is low, such that 85% of students stay in their region.⁸ In the period of analysis, 90% of students in Catalonia attend public universities, where tuition fees are highly subsidized.⁹ In 2018, Catalonia's PPS GDP per capita was €33200, slightly above Spain (€28100) and the EU (€31000) (Eurostat, 2020).

⁷Universitat Pompeu Fabra (1st), Universitat Autònoma de Barcelona (2nd), Universitat de Barcelona (3rd), Universitat Politècnica de Catalunya (6th) and Universitat Rovira i Virgili (7th).

⁸Source: El Mundo. According to Eurostat, Spain is one of the EU countries where young people live with their parents for longer, leaving at age 29.5, compared to an EU average of 26.

⁹Source: <https://www.idescat.cat/pub/?id=aec&n=753&t=2010>

We use administrative data on all applicants to Catalan universities, on the regular track (high school + PAU), who took the high-stakes exam and applied to college every year between 2006 and 2012. The main outcome variables in the sample are the students' Admission Grades and their Academic Programme (degree \times university) of admission. The main pre-determined covariates in the sample are parental and maternal education and occupation, postal code of residence, and high school. For every programme, we compute, every year, the threshold grade of admission, which is the lowest admission grade of a student that managed to enrol into that program, given the capacity constraints. For every programme, we observe the field of study, the faculty and the municipality where it is taught. We refer to these data as the *Selectivitat* dataset.

We combine these data with three additional datasets. First, with an administrative dataset of all students enrolled in public high schools for the post-reform period, including detailed information on their high-school grades.

Second, with a survey dataset of a sample of pre-reform students of Catalan universities, with information on their earnings four years after graduation, to compute the career prospects associated with each academic program.

Third, with an administrative dataset on college performance of the three main public universities in Catalonia (Universitat de Barcelona and Universitat Pompeu Fabra for all cohorts, and Universitat Autònoma de Barcelona for pre-treatment cohorts), which enrol more than 60% of students in Catalan Public Universities.

How do these data compare to the ideal data? First, we would like to observe the high school performance and the choices of high school students who do not apply to college, to understand whether the reform does affect selection into application. While we do not observe these students, we do observe very detailed information about students' characteristics, which allows us to understand whether changes in students' characteristics (such as maternal and parental education and occupation, postal code, high school, or month of birth) could

be driving the results. On the other hand, we observe whether students most likely to be affected by the reform are at the enrolment or application margin. We discuss this further in section 3.

Second, we would like to observe the high school and high-stakes grades of all pre-reform students. This would be useful to better understand what share of the effect of the reform is due to mechanical re-weighting. However, we only observe them separately post-reform. There are two symmetric ways of decomposing the effect into a mechanical and a behavioral component, which involve studying the effect of the reform holding constant test scores. On the one hand, holding pre-reform test scores constant (which we do not observe), applying post-reform weights. On the other hand, holding post-reform test scores constant, applying pre-reform weights (for which we present results). The intuition behind either decomposition is rather similar, although the results may differ. We discuss this further in section 3.1.

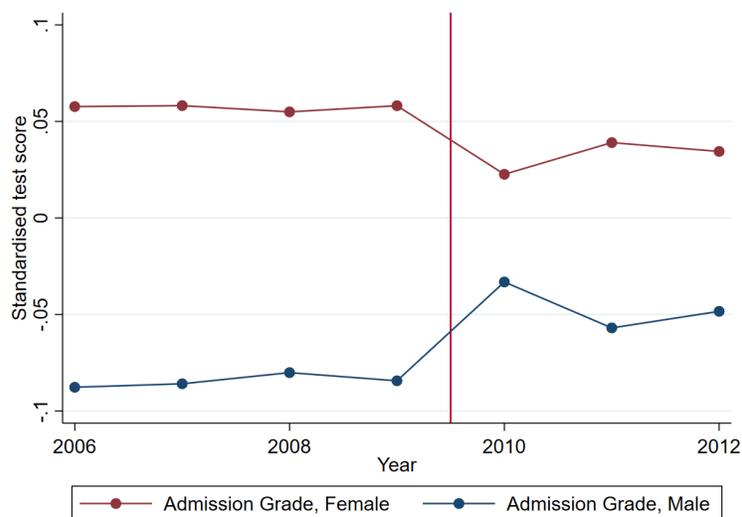
Finally, a very comprehensive evaluation of the reform would involve understanding its effects on labor market outcomes. Our focus, however, is to examine the effects on college admissions and performance. While we acknowledge that the effect on labor market outcomes is beyond the scope of our available data, we provide results about the effects on career prospects in the Appendix. Additionally, we highlight the effects on post-reform college performance in section 5.

3 Admission Grades

Figure 1 displays standardized admission grades by gender over time. Pre-reform, females' admission grades were around 0.14 standard deviations (s.d.) higher than males' admission grades, and this difference was stable over time. Post-reform, this difference shrinks to around 0.08 s.d.. Hence, the reform had a negative effect on females' admission grades.¹⁰

¹⁰Figure A1 in the Appendix displays female-by-year coefficients, where the baseline year is 2009 (the year before the reform), showing that these differences are statistically significant.

Figure 1: Effect of the reform on admission grades by gender



We also estimate differences-in-differences regressions:

$$Admission\ Grade_{it} = \alpha_t + \beta Female_i + \gamma (Female_i \times Post_t) + \epsilon_{it}$$

Where we regress the admission grade of student i in year t on year fixed effects α_t , a female indicator, and a post-reform indicator (year = 2010, 2011, 2012) interacted with a female indicator. Table 1 reports point estimates. As suggested by figure 1, the reform had a significant negative effect on females' admission grades. Adding gender-specific time trends and controls (parental and maternal education and occupation dummies, high school, postal code, nationality), the estimates show a very similar picture. This suggests that the results are not driven by differential trends by gender nor by changes in socio-economic characteristics by gender over time. Quantitatively, the magnitude of this effect is similar to the effect of taking an exam on a high pollution day (Ebenstein *et al.*, 2016), or to the date of birth effect (January-December) in our sample, as reported in table A3 in the Appendix; or to 15% of the parental college education gradient in admission scores, as reported in table A4 in the Appendix. Throughout the paper, we use robust standard errors, since we focus on two groups and cluster-robust standard errors do not perform well in that case. We cluster our standard errors only whenever our explanatory variable features a clustered pattern (Abadie *et al.*, 2023).

Table 1: Dependent Variable: Admission Grade

	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.0636*** (0.00953)	-0.0741*** (0.0189)	-0.0739*** (0.00872)	-0.0662*** (0.0172)
Female	0.142*** (0.00645)		0.110*** (0.00601)	
Female	✓	✓	✓	✓
Year FE	✓	✓	✓	✓
Gender-specific trends		✓		✓
Controls			✓	✓
Mean Dep. Var	7.08e-08	7.08e-08	-0.00240	-0.00240
N	183451	183451	182259	182259

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Controls: Mother and father education and occupation, year and month of birth, nationality, high school, postal code.

It is also interesting to measure the effect in terms of students' percentile rank (i.e., the percentage of scores that are equal to or lower than the student's admission grade), which is closely linked to college admissions. Table A5 in the Appendix reports point estimates of the effect on the admission grade rank, where the rank is equal to one for the highest admission grade, and zero for the lowest admission grade. Pre-reform, females were ranked 4% higher, on average, and post-reform this declines to around 2%. Again, adding gender-specific time trends or controls does not substantially change the point estimates.

Regarding the role of W_B and W_C , table A2 in the Appendix show that there are no gender differences, and table A6 in the Appendix reports estimates setting $W_B = W_C = 19$ (i.e., the average weight in the sample) for all students (columns 1 and 2) and controlling for individuals' W_B and W_C (columns 3 and 4), with very similar results.

3.1 Students' response

The reform has a significant effect on gender differences in admission grades, which are a weighted average of high-school and high-stakes grades. There are three possible sources for the estimated effect. First, an equilibrium effect of the reform on gender differences in high school vs. high-stakes performance. This will happen if students' behavior or effort reacts differently to the increased importance of the high-stakes exam. Second, a re-weighting of baseline gender differences in performance between high school and the high-stakes exam. If female and male students tend to perform relatively differently in high school compared to the high-stakes exam, we would expect the reform to affect admission grades via re-weighting. Third, a re-weighting effect due to differential performance across field and core subjects in the high-stakes exam. As explained in the previous section, the reform also changes the relative weight of core and field subjects in the high-stakes exam. If female and male students tend to perform differently in field subjects compared to core subjects, we would expect this to affect the admission grade as well.

We examine these alternative mechanisms by combining the Selectivitat dataset with administrative data on admission grades, high school grades and high-stakes grades for all post-reform students in Catalan public high schools. We weight the sample of public high schools so that it matches the full sample in terms of average admission grades by year and gender, using entropy balancing (Hainmueller, 2012).¹¹ Using the weighted sample, we study gender differences in the different components of the admission grades. This procedure is useful to interpret the results with respect to the full sample, since they both feature the same effect of the reform on female admission grades (and the same averages by year-gender).

The top panel in figure 2 displays gender differences in standardized high school grades and standardized high-stakes grades. First, it shows a very large gender difference in high

¹¹Weights are chosen by minimizing the entropy distance metric: $\min_{w_i} H(w) = \sum_{i \in Public\ Schools} w_i \log(w_i)$; subject to the balance constraint that the 1st and 2nd moment of the admission grade by year and gender of the re-weighted public school sample is equal to the population one.

school grades and a very small difference in high-stakes grades. This suggests that the re-weighting between baseline high school and high-stakes grades may have played an important role in the effect of the reform. Second, the figure shows both the post-treatment high-stakes GPA and a high-stakes GPA based on the pre-treatment formula, where core subjects have a 50% weight (as opposed to 60% under the new formula). The figure shows that in both cases, the gender differences in high-stakes performance remain almost identical. This suggests that the change in weights across field and core subjects in the high-stakes exam does not play an important role in the effect of the reform.

The bottom panel in figure 2 displays gender differences in standardized admission grades based on the pre-treatment formula, such that the high-stakes GPA has a weight of 43% for the admission grade, and the high school GPA a weight of 57%. It shows that if high-stakes and high school GPAs were to be re-weighted according to the pre-reform weights, the effect of the reform would have been smaller (and similarly smaller regardless of whether core subjects count for 60% or 50% within the high-stakes exam). Table 2 reports the corresponding figures, comparing (1st row) the gender difference in admission grades before the reform, (2nd row) after the reform but using the pre-reform weights, and (3rd row) after the reform. It shows that if weights had stayed the same, the gender difference in admission grades would have barely changed. Denoting by w the vector of weights on high school and high stake exam grades, and by g the vector of grades, we can distinguish a mechanical and a behavioral effect:

$$\text{Total Effect} = \underbrace{(w_{post} - w_{pre}) \times g_{post}}_{\text{Mechanical Effect}} + \underbrace{w_{pre} \times (g_{post} - g_{pre})}_{\text{Behavioral Effect}}$$

The mechanical effect is computed by evaluating post-reform grades under the different weights, and corresponds to the difference between rows 2 and 3 in table 2. The behavioral effect is the change in performance weighted by the pre-reform weights. In the simple comparison above, the reform has an effect of 0.64, out of which 0.45 are mechanically due to

the re-weighting of high school and high stakes grades.¹²

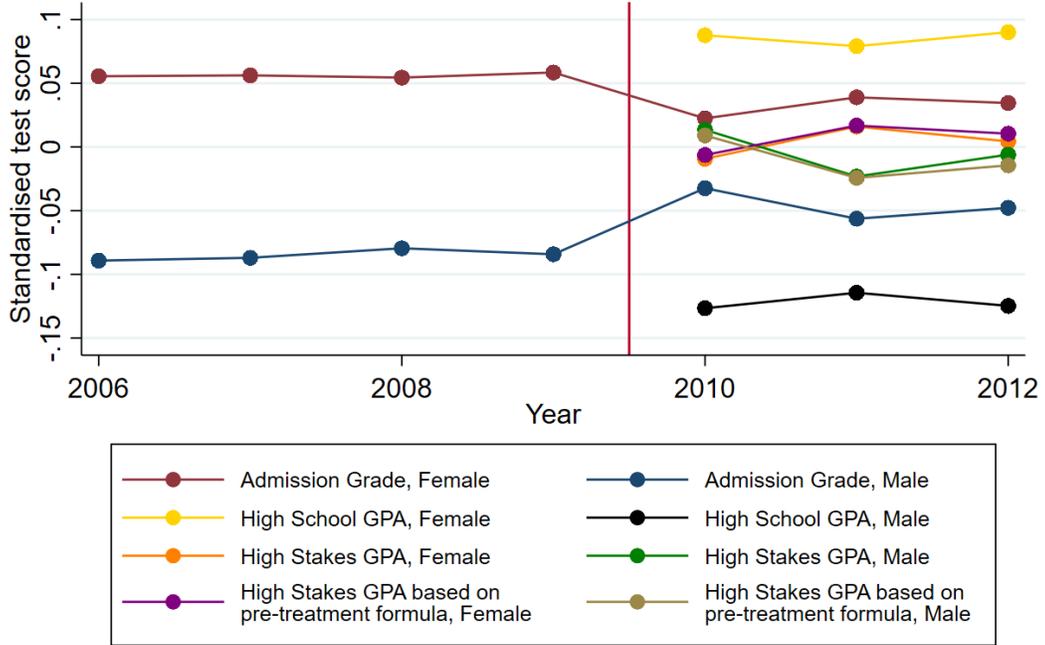
An alternative symmetric decomposition (which we cannot compute since we do not observe high stakes and high school GPA separately for the pre-reform period) is given by:

$$\text{Total Effect} = \underbrace{(w_{post} - w_{pre}) \times g_{pre}}_{\text{Mechanical Effect}} + \underbrace{w_{post} \times (g_{post} - g_{pre})}_{\text{Behavioral Effect}}$$

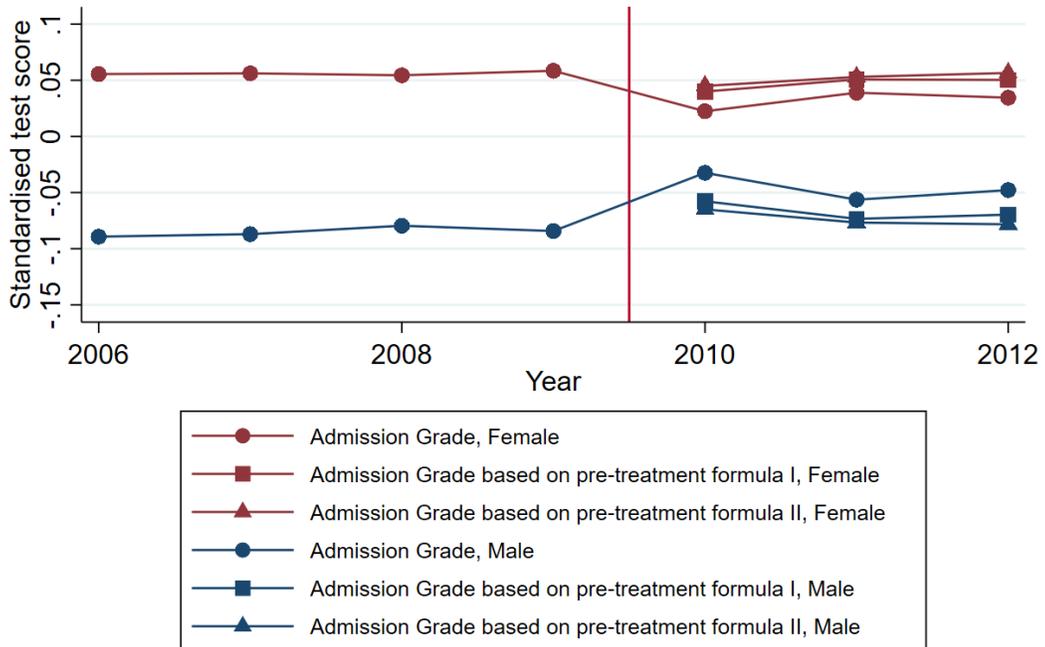
Where now the mechanical effect is computed by evaluating pre-reform grades under the different weights, and the behavioral effect is the change in performance weighted by the post-reform weights. This gives us the mechanical effect of moving from the post-reform to the pre-reform admission policy, while the previous decomposition gives us the symmetrical mechanical effect of moving from the post-reform to the pre-reform admission policy. Our results show that the behavioral effect, under the first decomposition, is small. If the behavioral effect would be zero, both decompositions would give the same result. When the behavioral effect is not zero ($g_{post} \neq g_{pre}$), then the second decomposition would deliver a smaller mechanical effect and a larger behavioral effect. For instance, if the reform had increased the gender difference between high school and high stake grades (that we do not observe because we observe them separately only post-reform) by 30%, the alternative decomposition would imply a mechanical effect of 55%, which is smaller than the 70% implied by our decomposition.

¹²In table A7 in the Appendix, we report estimates of the re-weighting effect of the reform with controls. The results likewise show that the re-weighting effect explains around 76% of the total effect of the reform on admission grades. We also report results for public schools, without weighting to match admission grades in the population, in table A8 in the Appendix. In this case, the effect also seems largely driven by re-weighting, although the effect of the reform on admission is smaller, which would suggest a slight behavioral reaction of the opposite sign.

Figure 2: Re-weighting and the effect of the reform



Sample: public schools, weighted to match admission scores by gender in the population via entropy balancing. High stakes GPA based on pre-treatment formula gives 50% weight to core subjects.



Sample: public schools, weighted to match admission scores by gender in the population via entropy balancing. Pre-treatment formula I: pre-treatment weights for high school vs. high stakes. Pre-treatment formula II: pre-treatment weights for high school vs. high stakes + pre-treatment weights for core vs. field subjects within high stakes GPA.

Table 2: Mechanical effects of the reform

Variable	(1)
	Difference
	Female - Male
Admission Grade, Pre-reform	0.141*** (0.011)
Admission Grade, Post-reform with Pre-reform weights	0.122*** (0.013)
Admission Grade, Post-reform	0.077*** (0.013)
High School GPA, Post-reform	0.208*** (0.013)
High Stakes GPA, Post-reform	0.016 (0.012)
High Stakes GPA, Post-reform with Pre-reform weights	0.021* (0.012)

Robust standard errors in parentheses. Pre-treatment weights: pre-treatment weights for high school vs. high stakes + pre-treatment weights for core vs. field subjects within high stakes GPA. Sample: public schools, weighted to match admission scores by gender in the population via entropy balancing. N pre-reform: 40295. N post-reform: 29933.

4 Students' allocation to college

In this section, we quantify the consequences of the gender differences in admission grades induced by the reform. The consequences will crucially depend on who are the most affected students, and whether they are competing for the same programmes. For instance, in an extreme case where female and male students' preferences were completely segregated, any effects on gender differences in admission grades would not affect the college allocation. We study three outcomes related to the allocation of students to college: enrolment, selectivity of the program of attendance, and career prospects.

Figure 3 displays admission grades over time across the predicted admission grades' distribution. In a first step, we regress admission grades on a vector of pre-determined covariates (namely parental and maternal education and occupation dummies, high school, postal code, month of birth, and nationality), for the pre-treatment sample. Then, we split

the sample according to whether students are predicted to be in different quartiles of the admission grade distribution. Figure 3 displays the effect of the reform across these groups (table A9 reports the corresponding point estimates). The main takeaway is that the most affected students are those expected to be top performers. For those expected to have lower grades, instead, the differences are small. This is consistent with the findings in Niederle and Vesterlund (2010) that performance gender gaps at high percentiles can partially be explained by the differential manner in which men and women respond to competitive test-taking environments. This also shows that the most affected students are not competing for enrolment into college, but for enrolment into rather selective programmes, and makes it less likely that our results are affected by idiosyncratic changes in enrolment/application behavior.

Figure 3: Effect of the reform on admission grades, along the performance distribution

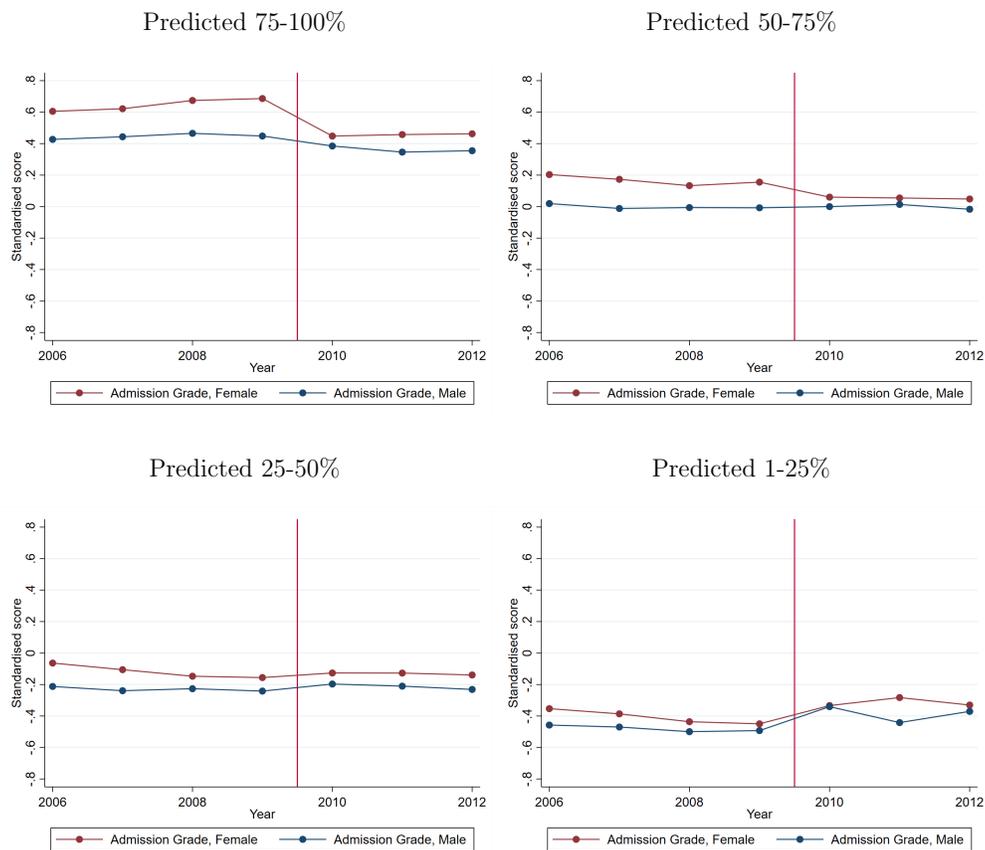
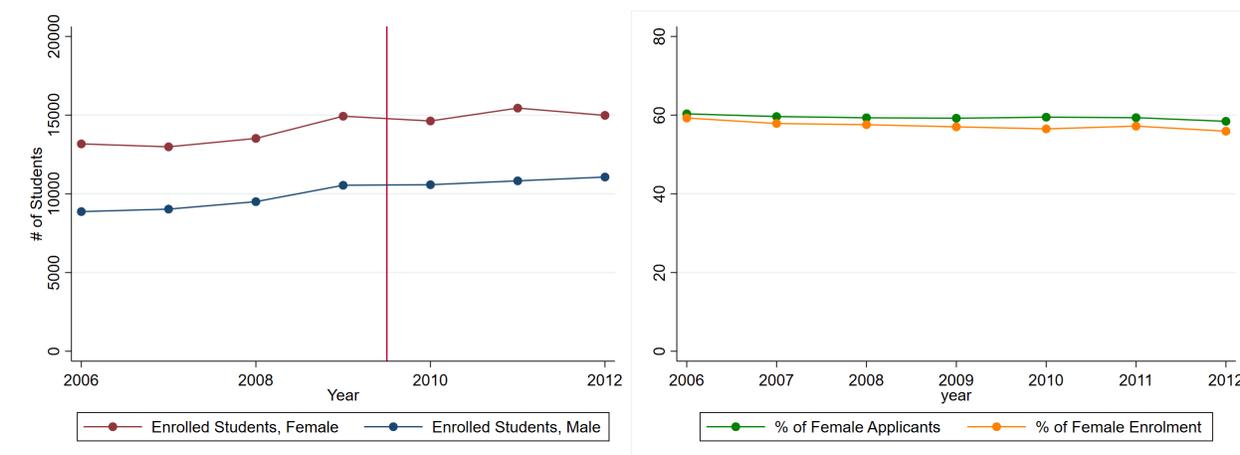


Figure 4 displays the number of enrolled students by year and gender, showing no differences due to the reform, as one would expect from figure 3.¹³

Figure 4: College enrolment



We next study the effect of the reform on another margin, namely the selectivity of the academic program attended. The Spanish setting provides a straightforward measure of access to more or less preferred or selective programmes, which is the *threshold grade* of the program of enrolment. The threshold grade is the admission grade of the student with the lowest admission grade who is admitted into a program. It is a measure of how selective is a programme, it is public information and strongly serially correlated. It is also a measure of peer quality and reputation: MacLeod *et al.* (2017) find that in Colombia, programmes' average admission grades across programmes causally matter for labour market outcomes.

To study the effect of the reform on gender differences in the selectivity of the program of enrolment, we rearrange the data and take academic programmes p as the unit of analysis and look at how the reform changes their gender composition depending on pre-reform threshold grades. Studying differences in the allocation according to pre-treatment threshold grades

¹³Enrolment is increasing during the period of analysis, which includes the great recession, in line with the literature on the counter-cyclicity of education (Arenas and Malgouyres, 2018). Spanish regions most affected by the crisis saw gender differences in educational attainment because of diminished blue-collar labor market opportunities in the construction sector (Aparicio-Fenoll, 2016), but these compliers are unlikely to be at the high school-college enrolment margin.

is useful because it keeps the measure of selectivity constant. Threshold Grades themselves are likely to be affected by the reform, and differently depending on the typical gender composition of academic programmes. An extreme case would be a scenario of full gender segregation across programmes: the reform would not change the students' allocation, but it would change average threshold grades by gender. Hence, we estimate the regression:

$$\%Females_{pt} = \alpha_p + \pi_t + \beta (\text{Pre-reform Thresh.Grade}_p \times Post_t) + \epsilon_{pt}$$

The outcome is the % of females in programme p in year t , we control for program fixed effects α_p and year fixed effects π_t , and the estimates are weighted by the number of students in each programme. Since the coding of academic programmes is fuzzy, with frequent changes that are difficult to track, we take university \times faculty \times municipality \times field of study as the unit of analysis, for which we obtain a more balanced panel. We obtain 210 units (on average, every unit offers 2.5 programmes per year). Table A10 in the Appendix shows that this is a meaningful grouping since there is a high serial correlation within this unit of observation in outcomes such as threshold grades, the number of enrolled students or the fraction of females.

Figure 5 displays the fraction of female students in programs above and below the median pre-reform level of selectivity. It shows that the reform affected the students' allocation, such that the percentage of female students in the most selective programs declines after the reform. The figure suggests that the percentage of female students in the most selective programmes declined by 3pp, compared to the percentage of females in the least selective programmes.

Figure 5: Fraction of Female Students by pre-reform Threshold Grade

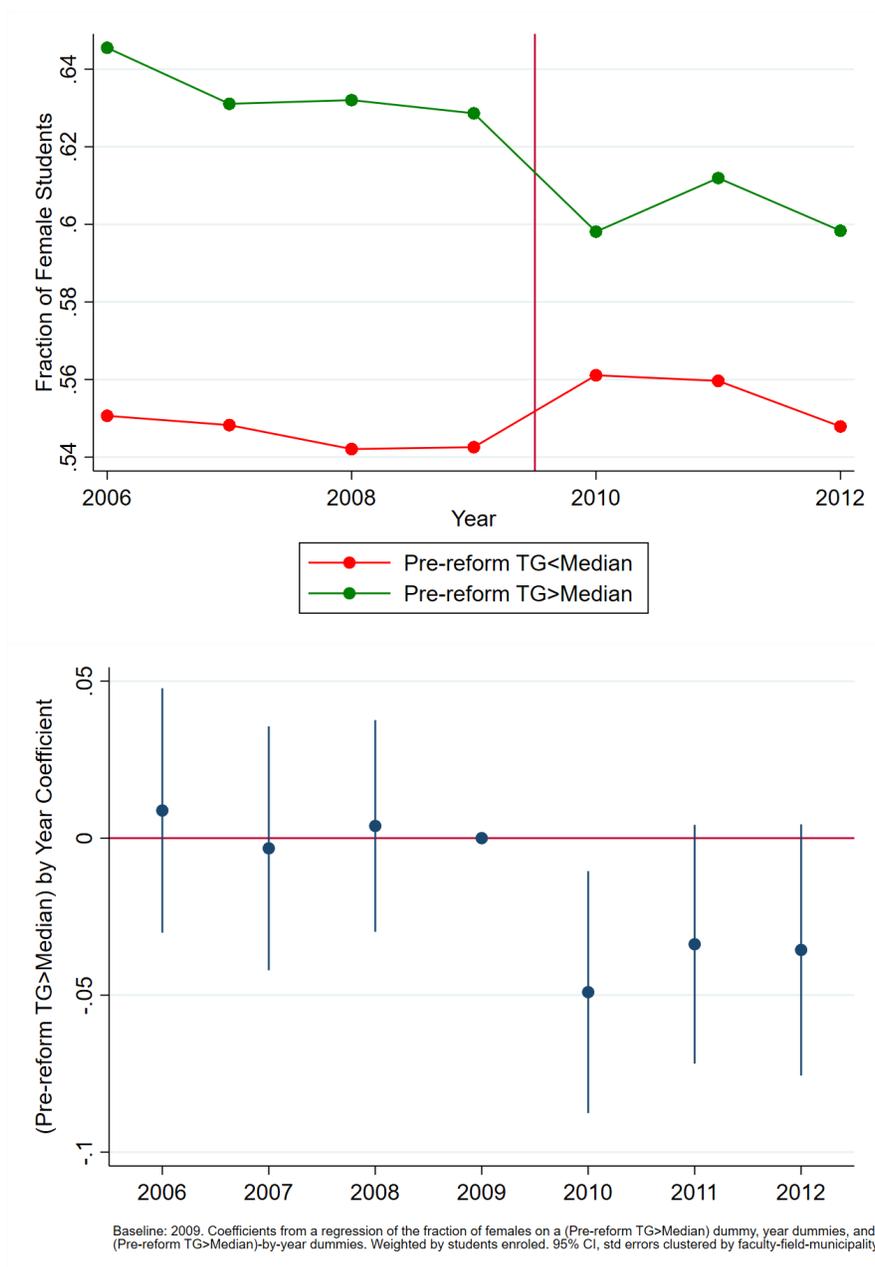


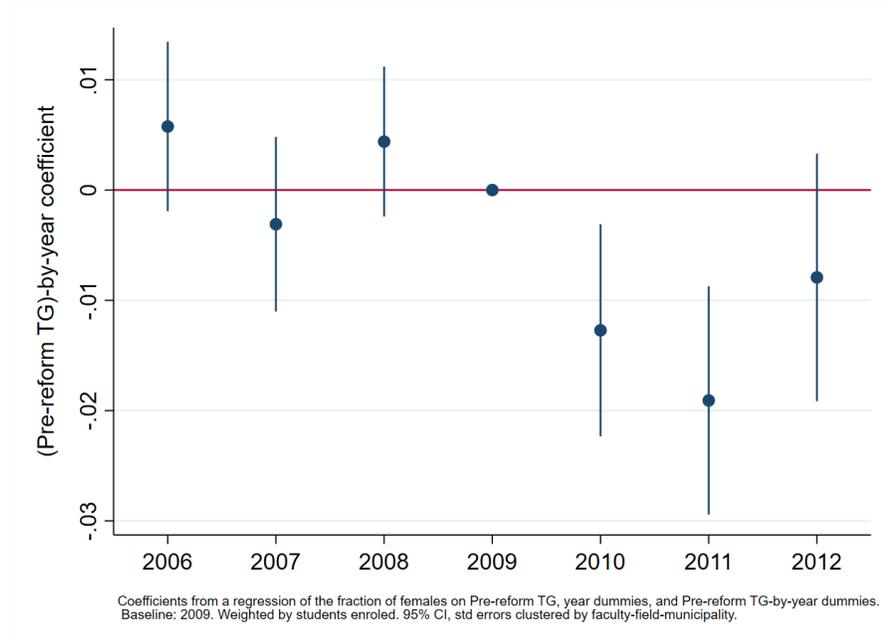
Table 3 reports differences-in-differences estimates with a continuous treatment measure.¹⁴ In a similar vein, the results show that the reform significantly decreased the percentage of female students in the most selective programmes. Compared to a program in the 25th percentile of selectivity, the percentage of female students in a program in the 75th percentile of selectivity declines by around 1.5 pp.

Table 3: Enrolment in selective programs

Dependent variable: fraction of female students		
	(1)	(2)
Post \times Pre-Reform T.Grade	-0.0146*** (0.00418)	-0.0167** (0.00775)
Faculty-Field-Municipality FE	✓	✓
Year FE	✓	✓
Faculty-Field-Municipality trends		✓
Mean Dep. Var	0.588	0.588
N	1018	1018

Standard errors clustered at the panel unit faculty-field-municipality in parentheses. Estimates weighted by the number of enrolled students.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$



¹⁴This includes Faculty-Field-Municipality FE, where pre-reform threshold grades are averaged in every cell.

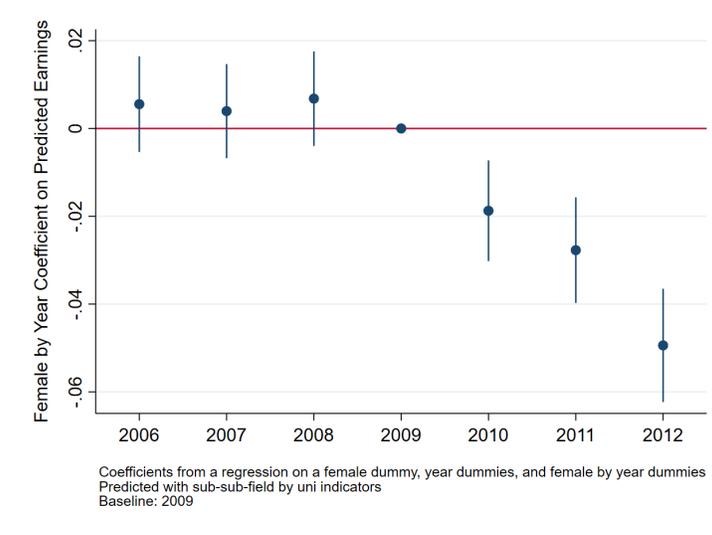
We also report estimates from individual-level regressions for the effect of the reform on the threshold grade of the program of enrolment in table A11 in the Appendix. In the first two columns, the dependent variable is the average pre-treatment threshold grade of the program of enrolment (again, at the level of university \times faculty \times municipality \times field of study), which keeps the selectivity measure constant. In the third and fourth columns, the dependent variable is the average pre or post-treatment threshold grade of the program of enrolment (again, at the level of university \times faculty \times municipality \times field of study). The results show a negative effect of the reform on the threshold grade of the program of enrolment when keeping its selectivity measure constant (columns 1 and 2), and an even larger effect on actual post-reform threshold grades (columns 3 and 4), which could be due to a decrease in the threshold grades of programs with a large percentage of female students. Finally, we also examine whether threshold grades change because of changes in enrolment across fields or within field, in table A12 in the Appendix. The results show that the effects on selectivity arise because female students move to similar but less selective programs, rather than across fields. This suggests that the effect is not driven by changes in choices or preferences for different fields over time.

Hence, overall, gender differences in admission grades due to the reform translate into significant changes in the colleges' allocation. The magnitude of the effect is again comparable to the date of birth effect on threshold grades in our sample (i.e., the effect of being born in January rather than in December); and to around 15% of the parental college education gradient in threshold grades, as reported by tables A13 and A14 in the Appendix.

In the Appendix B, we further study whether this effect on the selectivity of the program of enrolment is associated with changes in career prospects. This is interesting because threshold grades and wages are only positively correlated within field of study, and because female students tend to sort into fields and academic programmes with worse career (wage, employment) prospects. Hence, the effect will depend on whether students very much sub-

stitute their most preferred programmes for less selective programmes within the same field. Using a survey of pre-treatment college graduates to compute expected wages and employment by academic program, we estimate that the effect on the college allocation comes along with an increase of 2% in the expected gender wage gap four years after graduation (on top of a 20% wage gap) and with a small but significant effect on expected employment as well. Figure 6 displays female by year coefficients, where the baseline year is 2009, the last pre-reform year.

Figure 6: Career prospects: predicted log(wages)



5 Match quality

The reform has a significant effect on gender differences in admission grades and on the allocation of students to academic programmes, because of gender differences in high school vs. high-stakes performance. However, an open and very policy-relevant question is whether there is a trade-off between gender inequality and the quality of the match between students to college. To address this question, we study how gender differences in high-stakes performance in college admissions relate to college performance skills.

To this aim, we proceed in two steps. First, using machine learning techniques, we identify the types of students who are most likely to benefit from the reform (i.e., predicted winners and losers), based on a large set of detailed pre-determined student characteristics. Then, focusing on pre-treatment cohorts, we compare the college performance of students with the same admission grade and enrolled in the same program, college and (pre-reform) cohort, based on whether they are predicted to be winners or losers from the reform. The aim is to understand whether students who pre-reform were doing better in college (beyond what one would expect given their admission grades) are those most likely to gain from the reform and whether there are gender differences.

More precisely, in our first step, we estimate a prediction model for the heterogeneous effect of the reform across students, based on individual pre-determined covariates. An important concern about this type of prediction exercise is over-fitting. Over-fitting is a concern because, for instance, OLS coefficient estimates of the heterogeneous effects of the reform maximize the in-sample fit. Instead, Machine Learning methods, such as Lasso (least absolute shrinkage and selection operator), are estimated to maximize their out-of-sample predictive power, although the coefficient estimates cannot be interpreted as indicating any meaningful structure (Mullainathan and Spiess, 2017). Given the large set of covariates at hand and that we are interested in predicting the effect of the reform on admission scores, this is a suitable approach.

Lasso regressions are a form of penalized regression, with a penalty for each non-zero coefficient, that overcome over-fitting via cross-validation: slicing the sample into different parts, a training sample and a testing sample, and delivering estimates that maximize the predictive power of the training samples on the testing samples (Athey and Imbens, 2019).¹⁵ Lasso's $\hat{\beta}$ are the solution to: $\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{x}_i \beta')^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$, where

¹⁵In this case, the pre-treatment sample is sliced into ten different parts, as suggested by Kuhn and Johnson (2013) and Kohavi (1995).

$\lambda > 0$ is the Lasso penalty parameter that is chosen through cross-validation to maximize the out-of-sample performance of the training sample on the testing sample and p is the number of covariates.

In this case, we fit the Lasso models separately for the pre and post-reform periods, to obtain $\hat{\beta}(X)$, the predicted gain of the reform as a function of covariates X , where X is a vector of parental and maternal education and occupation dummies, postal code, high school, and month of birth dummies, all of them interacted with a gender indicator.

$$\hat{\beta}(X) = \widehat{Admission\ Grade}(X)^{Post} - \widehat{Admission\ Grade}(X)^{Pre}$$

Figure 7 plots the distribution of predicted effects of the reform $\hat{\beta}(X)$ by gender, where on average $\hat{\beta}(X, Female) = -0.03$ and $\hat{\beta}(X, Male) = 0.045$ (note that this is not symmetric because there are 60% of female students).

Figure 7: Distribution of expected gains from the reform

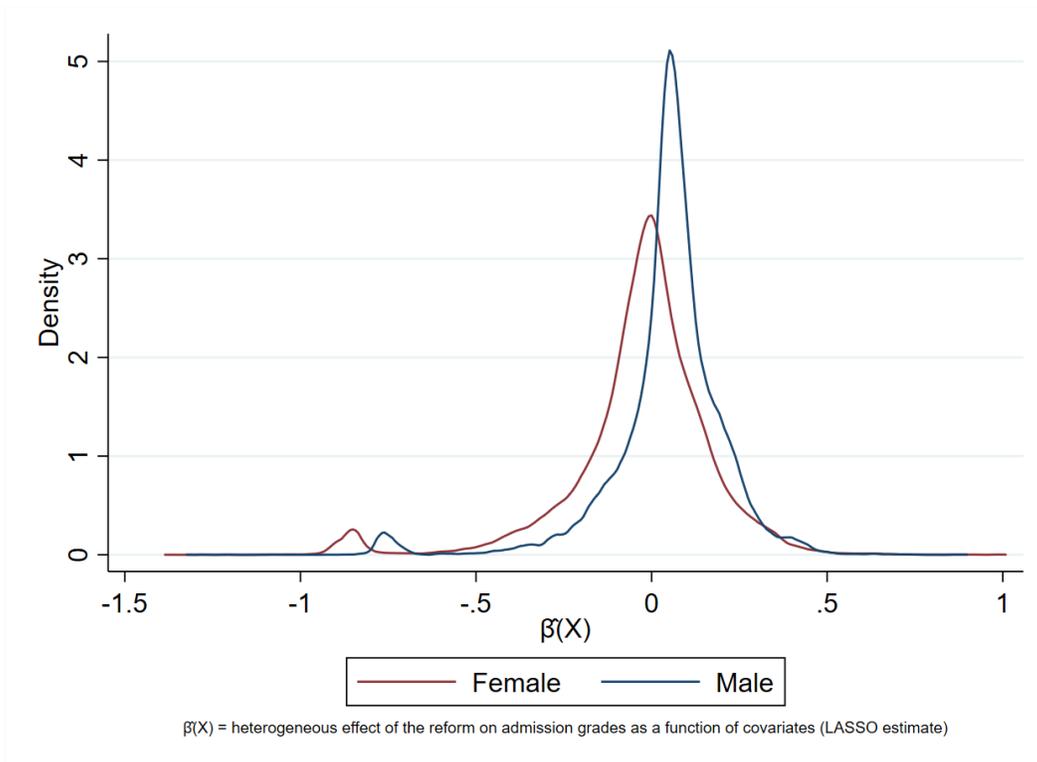


Table A16 in the Appendix shows that (1) $\widehat{\beta}(X)$ is more negative for individuals with high pre-reform admission grades and (2) females with higher pre-reform admission grades are those with the largest negative effect of the reform, which is consistent with figure 3 on the effect of the reform along the performance distribution. We further validate this measure by looking at its correlation with high school performance for the post-treatment cohorts (sample of public schools). We would expect that students predicted to gain from the reform are those doing relatively worse in high school. Figure A2 in the Appendix shows that within students with similar admission grades, those predicted to benefit from the reform are indeed those with worse high school grades (relative to their high-stakes performance).

Once we have obtained an individual-level measure of the predicted effects of the reform ($\widehat{\beta}(X)$), the second step is to relate it to college performance skills. The data on college performance by pre-treatment students enrolled comes from UB (Universitat de Barcelona), UAB (Universitat Autònoma de Barcelona) and UPF (Universitat Pompeu Fabra), which enrol around 61% of students in Catalan Public Universities.¹⁶ We merge these data with the college applications data (i.e., the *selectivitat* dataset).¹⁷ European undergraduate degrees are structured into subjects. Subjects have a number of credits (usually around 6 per subject, where a credit represents a certain amount of coursework time, which is standardized across all EU countries), and completion of an undergraduate degree typically requires passing 180 credits.

For UB, we observe, for all students in the 2006 to 2009 enrolling cohorts, for every year they are enrolled, the number of subjects (credits) they enrol, the number of credits they pass, and the average GPA in the passed subjects. For UAB, for all students in the 2006 to 2009 enrolling cohorts, the number of credits they enrol and pass, for the academic

¹⁶UB: 29.5%, UAB: 22.5%, UPF: 9%.

¹⁷We match the main college applications dataset with the college performance datasets, which are provided by universities, based on detailed demographics, matching 72.3% of students.

years 2008 to 2012.¹⁸ For UPF, for all students in the 2006 to 2009 enrolling cohorts, the yearly number of credits they enrol and pass. Hence, we use as the main measure of college performance the fraction of credits that a student passes out of the credits she enrolls during her time in college. We also present results with students' college GPA (average GPA in the completed subjects, unconditional on graduation, available for UB) in the Appendix.

We measure college performance with the residuals of a regression of the raw measure of college performance (fraction of credits passed out of credits enrolled and GPA, both standardized by cohort by academic programme) on admission grades: $\widetilde{CP}_i = CP_i - \widehat{CP}_i$. We weight the observations so that the college performance sample matches the population in admission grades by gender and cohort, using entropy balancing (Hainmueller, 2012), but report unweighted results in the Appendix as well.

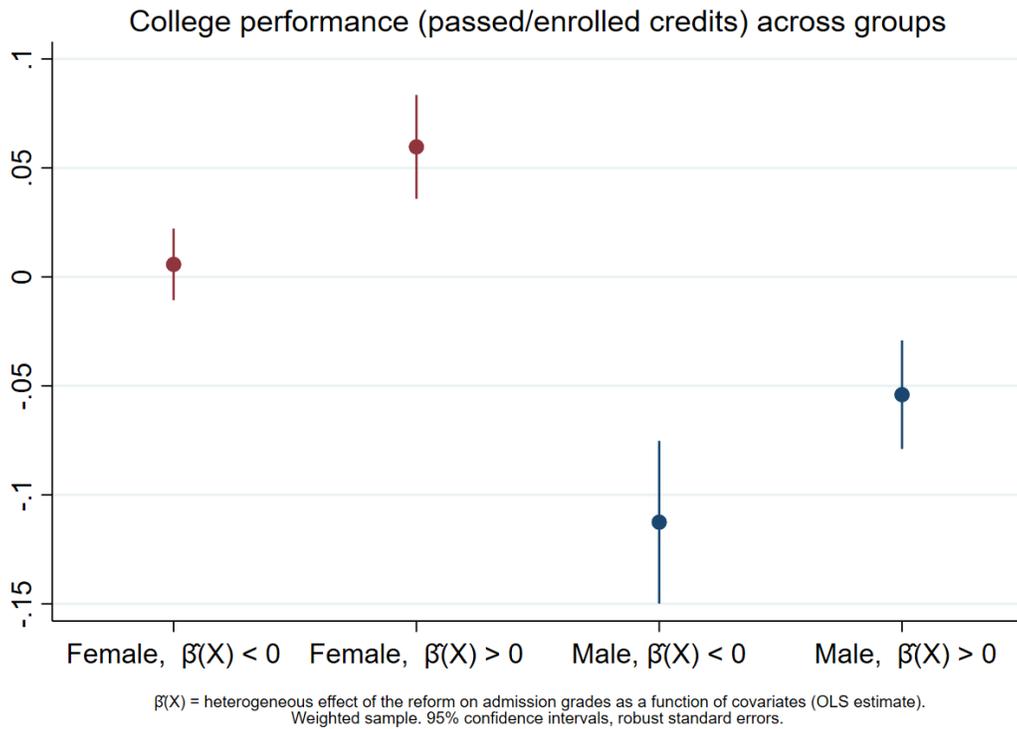
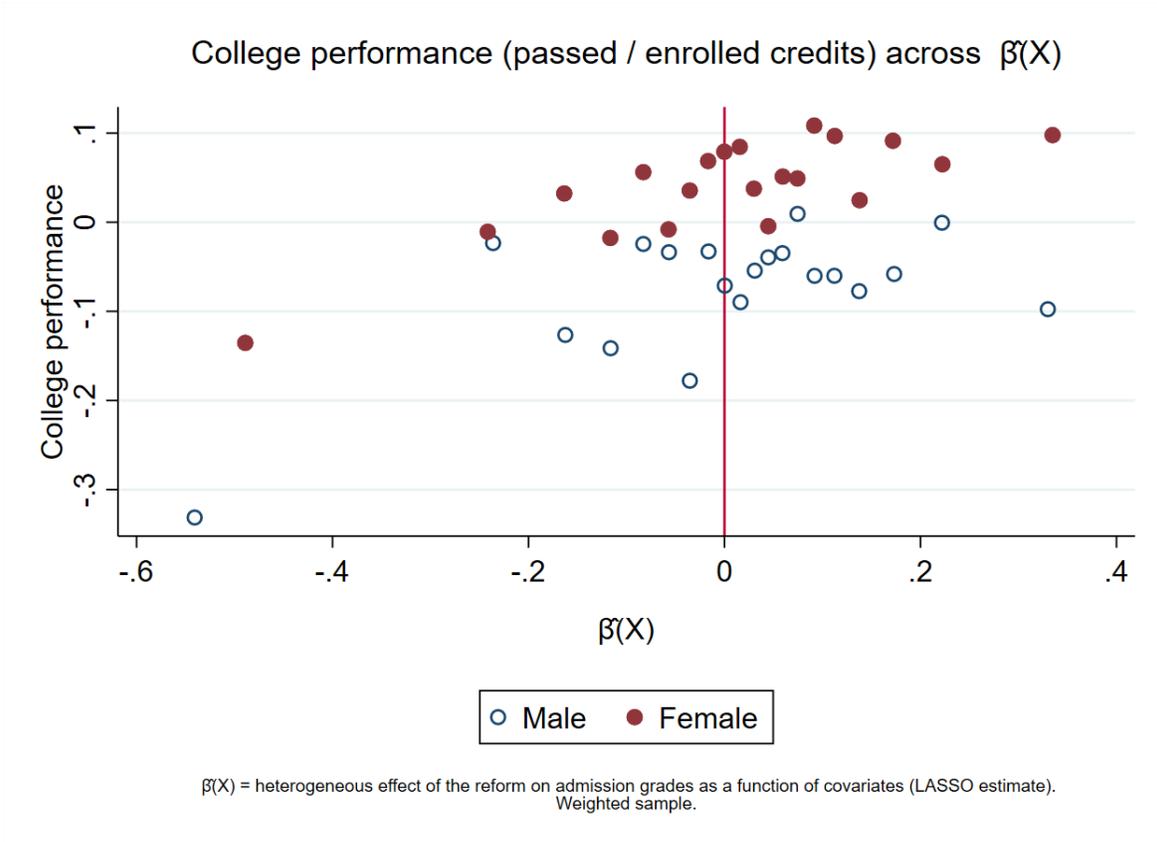
The top panel in figure 8 displays college performance of pre-reform students \widetilde{CP}_i over $\hat{\beta}(X)$. It shows three interesting patterns. First, a positive unconditional and within-gender correlation between college performance and expected gains from the reform, which suggests that high-stakes performance skills correlate positively with college performance skills.

Second, that female students perform better in college than male students with the same expected gains from the reform. This can be seen more precisely in the bottom panel of the figure, which splits college performance by gender and by expected winners and losers from the reform. The college performance of females with $\hat{\beta}(X) < 0$ is larger than the performance of males with $\hat{\beta}(X) < 0$, and the same for expected winners ($\hat{\beta}(X) > 0$).

Third, that females predicted to lose from the reform perform better in college than males predicted to benefit. Again, this is shown more precisely in the bottom panel. The performance of females expected to lose (i.e., with $\hat{\beta}(X) < 0$) is larger than that of males expected to win from the reform (i.e. with $\hat{\beta}(X) > 0$). Table A15 in the Appendix displays point estimates corresponding to both figures.

¹⁸This means that the students from the 2006 and 2007 cohorts are slightly positively selected because we observe them conditional on enrolment in their second or third year. However, dropping those cohorts does not change the results.

Figure 8: College performance and expected gains from the reform



The results show that although within gender, high-stakes performance skills positively correlate with college performance skills, the gender difference in high-stakes performance is negatively related to college performance skills. This means that the gender differences in admission scores induced by the reform may go against policy-makers objective functions aiming at selecting students based on their college performance potential.

Figure A3 in the Appendix displays the same figure for GPA rather than the fraction of credits passed, and figure A4 shows unweighted results, with a very similar pattern.

We also report results disaggregated by field of study in figure A5 the Appendix.¹⁹ The figure indicates that the results are largely driven by social science students. Investigating the mechanisms driving these heterogeneous effects is left for future research.

5.1 Post-reform college outcomes

Finally, we study how the reform affects college performance. We focus on two outcomes. First, passed subjects as a fraction of enrolled subjects. This is an indicator that is often used by universities to assess student performance, labeled as the “efficiency rate”. It relates to time to graduate and also it relates to a non-wasteful use of public resources, since college is subsidized. Second, we also study graduation rates.

Table 4 reports estimates of the effect of the reform on the *efficiency rate*. We find that it led to a decrease in the college performance of female students (column 1). To understand what is driving this result, we study whether females are now enrolling in academic programs where average performance tends to be lower. To this aim, in column (2), the dependent variable is the pre-reform average college performance in the students’ program of enrolment. The results show that due to the reform, females enroll in programs where performance tends to be lower. In column (3), the dependent variable is the difference between the individual

¹⁹The field composition in our sample vs. the population is the following: Arts-Humanities (14% in sample vs. 10% population), Science (17% vs. 9%), Social Sciences (51% vs. 43%), Health Science (12% vs. 15%), and Engineering (5% vs. 23%).

level performance (the outcome of column (1)) and the average pre-reform performance in the program (the outcome of column (2)). The effect on this measure is close to zero, which suggests that females do worse because they are enrolling in programs where performance tends to be lower. Finally, in column (4) we estimate the same specification of column (1), but controlling for $\widehat{\beta}(X)$ (i.e., the expected gains from the reform). The results show that this fully accounts for the gender differences in performance after the reform. This is reassuring, as it suggests that the performance effect is driven by those students enrolling into different programs as a result of the reform.

It may seem surprising that if female students are moving to less selective degrees, their performance declines. In table 5, we examine the relationship between performance and students' and program characteristics. It turns out that more selective programs (i.e., with lower threshold grades) feature a higher college performance on average (even beyond individuals' admission grades). Hence, moving to less selective programs means moving to programs where performance tends to be lower, which can explain why the effect of the reform on females' college performance is negative.

We also study how the reform affects the probability of graduation in table 6.²⁰ Here, we observe a similar pattern. Females are less likely to graduate after the reform, in part because they move to programs with lower graduation rates (which also tend to be less selective programs, as shown in table 7). Again, this is largely driven by female students expected to lose from the reform, since the effect shrinks and becomes non-significant once we control for $\widehat{\beta}(X)$.

²⁰Although we do not observe it directly for all the sample, we compute the maximum number of credits passed by a student in each program and cohort, and define graduation as obtaining more than 90% of that figure. That figure is not exactly the same for all students that graduate since a few credits can be validated with non-academic activities, and do not enter the count.

Table 4: Effects of the reform on college performance

	CP	CP ^{Pre-treat}	Δ CP	CP
	(1)	(2)	(3)	(4)
Female	0.0955*** (0.00430)	0.0576*** (0.00182)	0.0379*** (0.00388)	0.0976*** (0.00440)
Female \times Post 2009	-0.0129* (0.00679)	-0.0145*** (0.00273)	0.00159 (0.00621)	-0.00618 (0.00689)
Year FE	✓	✓	✓	✓
$\widehat{\beta}(X)$ & $\widehat{\beta}(X) \times$ Post 2009				✓
Mean Dep. Var	0.745	0.742	0.00248	0.745
N	43886	43886	43886	43886

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample. CP: ratio passed/enrolled credits. CP^{Pre-treat}: pre-reform college performance in the program of enrolment. Δ CP = CP - CP^{Pre-treat}.

$\widehat{\beta}(X)$: controls for expected gains from the reform.

Table 5: Performance across programs and students

	CP (credits passed/enrolled)		
	(1)	(2)	(3)
Admission Grade	0.125*** (0.00148)		0.105*** (0.00769)
Threshold Grade		0.117*** (0.00713)	0.0364*** (0.00792)
Year FE	✓	✓	✓
Mean Dep. Var	0.745	0.745	0.745
N	43886	43886	43886

Standard errors clustered by program of enrolment in parentheses.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample. CP: ratio passed/enrolled credits.

Table 6: Effects of the reform on graduation

	Grad.	Grad. ^{Pre-treat}	Δ Grad.	Grad.
	(1)	(2)	(3)	(4)
Female	0.105*** (0.00643)	0.0684*** (0.00327)	0.0363*** (0.00553)	0.107*** (0.00656)
Female \times Post 2009	-0.0190* (0.0102)	-0.0477*** (0.00453)	0.0287*** (0.00980)	-0.00977 (0.0104)
Year FE	✓	✓	✓	✓
$\hat{\beta}(X)$ & $\hat{\beta}(X) \times$ Post 2009				✓
Mean Dep. Var	0.605	0.587	0.0178	0.605
N	43886	43886	43886	43886

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample. Grad.: indicator for graduating. Grad^{Pre-treat}: pre-reform graduation rate in the program of enrolment. Δ Grad. = Grad. – Grad.^{Pre-treat}.

$\hat{\beta}(X)$: controls for expected gains from the reform.

Table 7: Graduation across programs and students

	Graduation		
	(1)	(2)	(3)
Admission Grade	0.113*** (0.00239)		0.0961*** (0.0122)
Threshold Grade		0.106*** (0.0189)	0.0320* (0.0174)
Year FE	✓	✓	✓
Mean Dep. Var	0.605	0.605	0.605
N	43886	43886	43886

Standard errors clustered by enrolment program in parentheses.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample.

Grad.: indicator for graduating.

6 Conclusions

The Gale-Shapley algorithm is one of the most popular college allocation algorithms around the world. A crucial policy question in its setting is designing admission priorities for students, understanding how they disadvantage certain demographic groups, and whether these differences are related to differences in college performance potential (i.e., whether these differences are fair). Our results, evaluating a policy change with administrative data, show that giving more weight to high-stakes exams for admissions has important cross-gender effects. In general, female students tend to outperform male students in high school, but gender differences in high-stakes performance are much smaller. We find a significant negative effect on female college admission scores of a reform that increased the weight of the comprehensive high-stakes exam at the end of high school for college admissions. A very substantial part of this effect is due to a re-weighting of the baseline high school vs. high-stakes performance differences, but the overall effect is slightly larger, suggesting that the effect of the reform is amplified by behavioral responses.

We further document that these effects have important consequences for the allocation of students to college. Most gender differences in admission scores induced by the reform happen at the top of the ability distribution, and as a result, the reform does not affect college enrolment. Nevertheless, the percentage of female students in the most selective degrees decreases significantly, and this comes along with a decline in their career prospects, widening expected gender gaps in the labour market.

Finally, we study whether the reform entails a trade-off between gender inequality and match quality. We find that within gender, good college performers tend to benefit from the reform. However, the results show that female students expected to lose from the reform are better college performers than male students expected to gain from the reform. Hence, the results show that gender differences in high-stake exam performance are not positively

related to determinants of college performance (if anything, these are negatively related). This is an important result for policy-makers designing college admission policies aiming at maximizing college performance potential in admissions while also taking into account gender differences in performance in different settings.

Finally, in future research, we hope to understand the implications of the changes in the gender and skill (high vs. low stakes performance) composition of graduating individuals (and therefore, of workers) for the economy and for inequality. Indeed, recent research by Buser *et al.* (2021) shows that individual competitiveness is a good predictor for education and labor market outcomes and that gender differences in competitiveness can explain 5-10 percent of the observed gender differences in education and labor market outcomes.

References

- ABADIE, A., ATHEY, S., IMBENS, G. W. and WOOLDRIDGE, J. M. (2023). When should you adjust standard errors for clustering? *The Quarterly Journal of Economics*, **138** (1), 1–35.
- ABDULKADIROGLU, A. and ANDERSSON, T. (2022). School choice. *NBER Working Paper*, (w29822).
- APARICIO-FENOLL, A. (2016). Returns to education and educational outcomes: The case of the Spanish housing boom. *Journal of Human Capital*, **10** (2), 235–265.
- ARENAS, A., CALSAMIGLIA, C. and LOVIGLIO, A. (2021). What is at stake without high-stakes exams? students’ evaluation and admission to college at the time of covid-19. *Economics of Education Review*, **83**, 102143.
- and MALGOUYRES, C. (2018). Countercyclical school attainment and intergenerational mobility. *Labour Economics*, **53**, 97–111.

- ATHEY, S. and IMBENS, G. W. (2019). Machine learning methods that economists should know about. *Annual Review of Economics*, **11**, 685–725.
- AZMAT, G., CALSAMIGLIA, C. and IRIBERRI, N. (2016). Gender differences in response to big stakes. *Journal of the European Economic Association*, **14** (6), 1372–1400.
- BALAFOUTAS, L. and SUTTER, M. (2012). Affirmative action policies promote women and do not harm efficiency in the laboratory. *Science*, **335** (6068), 579–582.
- BUSER, T., NIEDERLE, M. and OOSTERBEEK, H. (2021). *Can competitiveness predict education and labor market outcomes? Evidence from incentivized choice and survey measures*. Tech. rep., National Bureau of Economic Research.
- CAI, X., LU, Y., PAN, J. and ZHONG, S. (2018). Gender Gap under Pressure: Evidence from China’s National College Entrance Examination. *The Review of Economics and Statistics*.
- CALSAMIGLIA, C. and GÜELL, M. (2018). Priorities in school choice: The case of the boston mechanism in barcelona. *Journal of Public Economics*, **163**, 20–36.
- DUCKWORTH, A. L. and SELIGMAN, M. E. (2006). Self-discipline gives girls the edge: Gender in self-discipline, grades, and achievement test scores. *Journal of educational psychology*, **98** (1), 198.
- EBENSTEIN, A., LAVY, V. and ROTH, S. (2016). The long-run economic consequences of high-stakes examinations: Evidence from transitory variation in pollution. *American Economic Journal: Applied Economics*, **8** (4), 36–65.
- GALE, D. and SHAPLEY, L. S. (1962). College admissions and the stability of marriage. *The American Mathematical Monthly*, **69** (1), 9–15.
- GNEEZY, U., NIEDERLE, M. and RUSTICHINI, A. (2003). Performance in competitive environments: Gender differences. *Quarterly Journal of Economics*, **118** (3), 1049–1074.

- and RUSTICHINI, A. (2004). Gender and competition at a young age. *American Economic Review*, **94** (2), 377–381.
- GORTÁZAR, L., MAYOR, D. and MONTALBÁN, J. (2023). Residence-based priorities and school choice. *Economics of Education Review*, **95**, 102384.
- GRAETZ, G. and KARIMI, A. (2022). Gender gap variation across assessment types: Explanations and implications. *Economics of Education Review*, **91**, 102313.
- HAINMUELLER, J. (2012). Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies. *Political analysis*, pp. 25–46.
- IRIBERRI, N. and REY-BIEL, P. (2019). Competitive pressure widens the gender gap in performance: Evidence from a two-stage competition in mathematics. *The Economic Journal*, **129** (620), 1863–1893.
- JURAJDA, Š. and MÜNICH, D. (2011). Gender gap in performance under competitive pressure: Admissions to Czech universities. *American Economic Review*, **101** (3), 514–18.
- KAUFMANN, K. M., MESSNER, M. and SOLIS, A. (2021). *Elite Higher Education, the Marriage Market and the Intergenerational Transmission of Human Capital*. Crc tr 224 discussion paper series, University of Bonn and University of Mannheim, Germany.
- KIRKEBØEN, L., LEUVEN, E. and MOGSTAD, M. (2021). *College as a marriage market*. Tech. rep., National Bureau of Economic Research.
- KIRKEBØEN, L. J., LEUVEN, E. and MOGSTAD, M. (2016). Field of study, earnings, and self-selection. *The Quarterly Journal of Economics*, **131** (3), 1057–1111.
- KOHAVI, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th international joint conference on Artificial intelligence-Volume 2*, pp. 1137–1143.
- KUHN, M. and JOHNSON, K. (2013). *Applied predictive modeling*, vol. 26. Springer.

- MACLEOD, W. B., RIEHL, E., SAAVEDRA, J. E. and URQUIOLA, M. (2017). The big sort: College reputation and labor market outcomes. *American Economic Journal: Applied Economics*, **9** (3), 223–61.
- MONTOLIO, D. and TABERNER, P. A. (2018). Gender differences under test pressure and their impact on academic performance: a quasi-experimental design. *IEB WP 2018/21*.
- MORIN, L.-P. (2015). Do men and women respond differently to competition? Evidence from a major education reform. *Journal of Labor Economics*, **33** (2), 443–491.
- MULLAINATHAN, S. and SPIESS, J. (2017). Machine learning: an applied econometric approach. *Journal of Economic Perspectives*, **31** (2), 87–106.
- NIEDERLE, M. (2015). Gender. In *Handbook of Experimental Economics*, Princeton University Press.
- , SEGAL, C. and VESTERLUND, L. (2013). How costly is diversity? Affirmative action in light of gender differences in competitiveness. *Management Science*, **59** (1), 1–16.
- and VESTERLUND, L. (2007). Do women shy away from competition? do men compete too much? *The Quarterly Journal of Economics*, **122** (3), 1067–1101.
- and — (2010). Explaining the gender gap in math test scores: The role of competition. *Journal of Economic Perspectives*, **24** (2), 129–44.
- and — (2011). Gender and competition. *Annu. Rev. Econ.*, **3** (1), 601–630.
- ORS, E., PALOMINO, F. and PEYRACHE, E. (2013). Performance gender gap: does competition matter? *Journal of Labor Economics*, **31** (3), 443–499.
- ROTH, A. E. and SOTOMAYOR, M. (1992). Two-sided matching. *Handbook of game theory with economic applications*, **1**, 485–541.
- SAYGIN, P. O. (2018). Gender bias in standardized tests: evidence from a centralized college admissions system. *Empirical Economics*, pp. 1–29.

- SCHLOSSER, A., NEEMAN, Z. and ATTALI, Y. (2019). Differential Performance in High vs. Low Stakes Tests: Evidence from the Gre. *Economic Journal*, **129** (10), 2916–2948.
- SÖDERSTRÖM, M. and UUSITALO, R. (2010). School choice and segregation: Evidence from an admission reform. *Scandinavian Journal of Economics*, **112** (1), 55–76.
- UNESCO (2017). Global education monitoring report.

Appendix

Table A1: Distribution of weights by gender

Male students			Female students		
Weight	Percent	Cum.	Weight	Percent	Cum.
40/140	15.60	15.60	40/140	14.15	14.15
50/140	3.55	19.15	50/140	3.70	17.85
60/140	23.60	42.74	60/140	25.02	42.87
70/140	6.68	49.43	70/140	7.98	50.85
80/140	50.57	100.00	80/140	49.15	100.00
Total	100.00		Total	100.00	

Table A2: Heterogeneity in field subjects' weights

	(1)	(2)
	1(Taking all exams)	Average weight
Female	0.00631* (0.00343)	0.0663 (0.0514)
Year FE	✓	✓
Mean Dep. Var	0.585	13.69
<i>N</i>	84677	84677

Sample: post-reform. Robust standard errors in parentheses.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Avg. weight = $\frac{W_B + W_C}{2}$, with $or = 0$ if the exam is not taken.

Note: column (1) - OLS regression of an indicator for taking all exams on a female indicator and year FE.

column (2) - OLS regression of the average weight at admission on a female indicator and year FE.

Table A3: Date of birth effect

Dependent variable: admission grade		
	(1)	(2)
Born in January	0.0727*** (0.0115)	0.0724*** (0.0115)
Year FE		✓
Mean Dep. Var	-0.00120	-0.00120
<i>N</i>	30255	30255

Sample: born in January or December.

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A4: Parental education gradient in admission grades

Dependent variable: admission grade				
	(1)	(2)	(3)	(4)
Both college educated	0.461*** (0.00557)			
At least one college educated		0.388*** (0.00464)		
Mother college educated			0.407*** (0.00491)	
Father college educated				0.384*** (0.00494)
Year FE	✓	✓	✓	✓
Mean Dep. Var	7.08e-08	7.08e-08	7.08e-08	7.08e-08
Mean Indep. Var	0.233	0.455	0.347	0.341
<i>N</i>	183451	183451	183451	183451

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: OLS regressions of admission grades on one control at a time and year FE in all cases.

Table A5: Dependent Variable: Admission Grade Rank

	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.0205*** (0.00275)	-0.0222*** (0.00547)	-0.0229*** (0.00251)	-0.0192*** (0.00497)
Female	0.0430*** (0.00187)		0.0328*** (0.00173)	
Female	✓	✓	✓	✓
Year FE	✓	✓	✓	✓
Gender-specific trends		✓		✓
Controls			✓	✓
Mean Dep. Var	0.499	0.499	0.499	0.499
<i>N</i>	183451	183451	182259	182259

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Rank is equal to one for the highest score within a cohort and zero for the lowest.

Controls: Mother and father education and occupation, year and month of birth, nationality, high school, postal code.

Table A6: Robustness Admission Grades

	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.0506*** (0.00952)	-0.0608*** (0.00872)	-0.0661*** (0.00761)	-0.0533*** (0.00712)
Female	0.142*** (0.00645)	0.109*** (0.00600)	0.142*** (0.00645)	0.108*** (0.00596)
Female	✓	✓	✓	✓
Year FE	✓	✓	✓	✓
Post \times Subject weights			✓	✓
Controls		✓		✓
Mean Dep. Var	-1.57e-08	-0.00244	7.08e-08	-0.00240
Subject weights	19	19	Baseline (enrolment)	Baseline (enrolment)
N	183451	182259	183451	182259

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Controls: Mother and father education and occupation, year and month of birth, nationality, high school, postal code.

Table A7: Re-weighting effects of the reform

	Admission Grade		Δ Admission Grade Admission Grade, based on pre-treatment formula	
	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.0641*** (0.0165)	-0.0621*** (0.0152)	-0.0475*** (0.00256)	-0.0452*** (0.00254)
Female	0.141*** (0.0106)	0.114*** (0.00985)		
Year FE	✓	✓	✓	✓
Controls		✓		✓
Mean Dep. Var	-0.000151	-0.0000259	-0.000338	-0.000303
N	70228	70067	70228	70067

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Sample: public schools, weighted to match admission scores by gender/year in the population via entropy balancing. Pre-treatment formula: pre-treatment weights for high school vs. high stakes and for core vs. field subjects within the high stakes GPA.

Controls: Mother and father education and occupation, year and month of birth, nationality, high school, postal code.

Table A8: Re-weighting effects of the reform, public schools, unweighted.

	Admission Grade		Δ Admission Grade Admission Grade, based on pre-treatment formula	
	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.0354** (0.0149)	-0.0389*** (0.0140)	-0.0482*** (0.00264)	-0.0458*** (0.00261)
Female	0.103*** (0.00978)	0.0839*** (0.00916)		
Year FE	✓	✓	✓	✓
Controls		✓		✓
Mean Dep. Var	-0.139	-0.138	-0.000646	-0.000609
N	70228	70067	70228	70067

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Sample: public schools. Pre-treatment formula: pre-treatment weights for high school vs. high stakes and for core vs. field subjects within the high stakes GPA.

Controls: Mother and father education and occupation, year and month of birth, nationality, high school, postal code.

Table A9: Effect of the reform along the performance distribution

	Dependent variable: admission grades			
	(1)	(2)	(3)	(4)
Female \times Post 2009	-0.00511 (0.0171)	-0.0294 (0.0179)	-0.112*** (0.0186)	-0.109*** (0.0187)
Female	✓	✓	✓	✓
Year FE	✓	✓	✓	✓
Sample	Predicted 1-25%	Predicted 25-50%	Predicted 50-75%	Predicted 75-100%
Mean Dep. Var	-0.396	-0.164	0.0646	0.496
N	45863	45863	45863	45862

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A10: Within faculty-field-municipality autocorrelation in outcomes

	(1)	(2)	(3)
	Thresh. Grade	#Enrolled Students	%Female Enrolled
Lagged T.Grade	0.949*** (0.0167)		
Lagged #Enrolled Students		0.987*** (0.00750)	
Lagged %Female Enrolled			0.951*** (0.00840)
Mean Dep. Var	-0.873	341.6	0.584
N	874	874	874

Standard errors clustered by the panel unit faculty-field-municipality.

* $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Estimates weighted by the number of enrolled students.

Note: OLS regressions of the dependent variable on its lagged value.

Table A11: Threshold grades, program of enrolment

	Thresh G. (pre-treat).		Thresh G. (actual).	
	(1)	(2)	(3)	(4)
Female	0.168*** (0.00498)		0.168*** (0.00498)	
Female \times Post 2009	-0.0262*** (0.00769)	-0.0537*** (0.0155)	-0.0720*** (0.00795)	-0.107*** (0.0158)
Year FE	✓	✓	✓	✓
Gender-specific trends		✓		✓
Mean Dep. Var	-0.833	-0.833	-0.872	-0.872
N	166372	166372	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Thresh G. (pre-treat): based on pre-reform avg. values by faculty-field-municipality.

Thresh G. (actual): based on pre and post-reform averages by faculty-field-municipality.

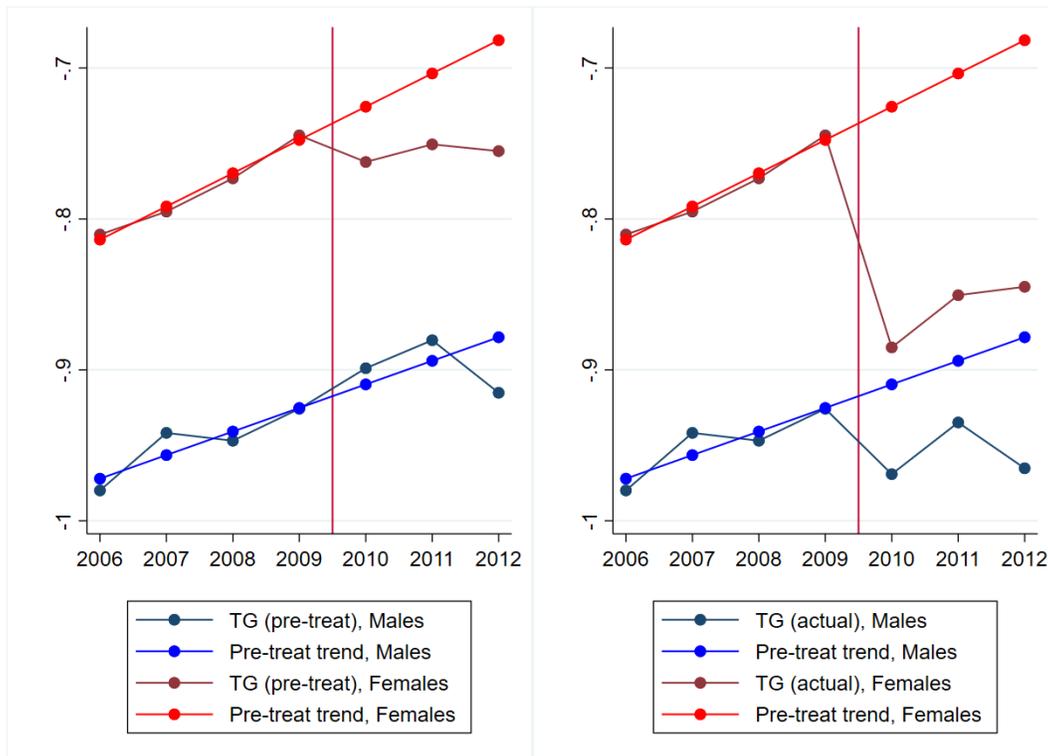


Table A12: Allocation change across vs. within-field

	TG(pre-treat)	TG(pre-treat field avg.)	TG(pre-treat) - TG(pre-treat field avg.)
	(1)	(2)	(3)
Female	0.168*** (0.00498)	0.0535*** (0.00173)	0.115*** (0.00467)
Female \times Post 2009	-0.0262*** (0.00769)	-0.00223 (0.00270)	-0.0220*** (0.00715)
Year FE	✓	✓	✓
Mean Dep. Var	-0.833	-0.841	0.00763
N	166372	169955	166372

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

TG(pre-treat): threshold grade of enrolment, based on pre-reform avg. values by faculty-field-municipality.

TG(pre-treat field avg.): average by field of study.

Table A13: Date of birth effect

Dependent variable: threshold grade		
	(1)	(2)
Born in January	0.0585*** (0.0112)	0.0591*** (0.0112)
Year FE		✓
Mean Dep. Var	-0.870	-0.870
N	28063	28063

Sample: born in January or December.

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A14: Parental education gradient in threshold grades

Dependent variable: threshold grade				
	(1)	(2)	(3)	(4)
Both college educated	0.390*** (0.00584)			
At least one college educated		0.318*** (0.00457)		
Mother college educated			0.330*** (0.00496)	
Father college educated				0.328*** (0.00499)
Year FE	✓	✓	✓	✓
Mean Dep. Var	-0.872	-0.872	-0.872	-0.872
Mean Indep. Var	0.230	0.452	0.344	0.337
N	170082	170082	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Note: OLS regressions of threshold grades on one control at a time and year FE in all cases.

Table A15: College performance and expected gains from the reform

College performance (passed/enrolled credits)	
(1)	
Female, $\widehat{\beta}(X) > 0$	0.0597*** (0.0122)
Male, $\widehat{\beta}(X) < 0$	-0.113*** (0.0190)
Male, $\widehat{\beta}(X) > 0$	-0.0540*** (0.0127)
Intercept (Female, $\widehat{\beta}(X) < 0$)	0.00575 (0.00839)
Mean Dep. Var	7.77e-11
N	39159
Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample.	
College performance (passed/enrolled credits)	
(1)	
$\widehat{\beta}(X)$	0.242*** (0.0588)
Female	0.116*** (0.0109)
Female \times $\widehat{\beta}(X)$	0.0731 (0.0691)
Intercept	-0.0753*** (0.00900)
Mean Dep. Var	7.77e-11
N	39159
Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample.	

Table A16: Predicted effect of the reform and pre-reform admission grades

	$\widehat{\beta}(X)$
	(1)
Admission Grade	-0.0120*** (0.000952)
Female	-0.0736*** (0.00121)
Female \times Admission Grade	-0.0194*** (0.00132)
Intercept	0.0445*** (0.000868)
Mean Dep. Var	-3.35e-08
N	98774

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Pre-reform sample.

Table A17: College performance potential and expected gains from the reform, OLS prediction of $\widehat{\beta}(X)$

	College performance (passed/enrolled credits)
	(1)
Female, $\widehat{\beta}(X) > 0$	0.0696*** (0.0121)
Male, $\widehat{\beta}(X) < 0$	-0.0832*** (0.0162)
Male, $\widehat{\beta}(X) > 0$	-0.0488*** (0.0136)
Intercept (Female, $\widehat{\beta}(X) < 0$)	-0.000615 (0.00862)
Mean Dep. Var	7.77e-11
N	39159

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Weighted sample.

Figure A1: Effect of the reform on admission grades by gender

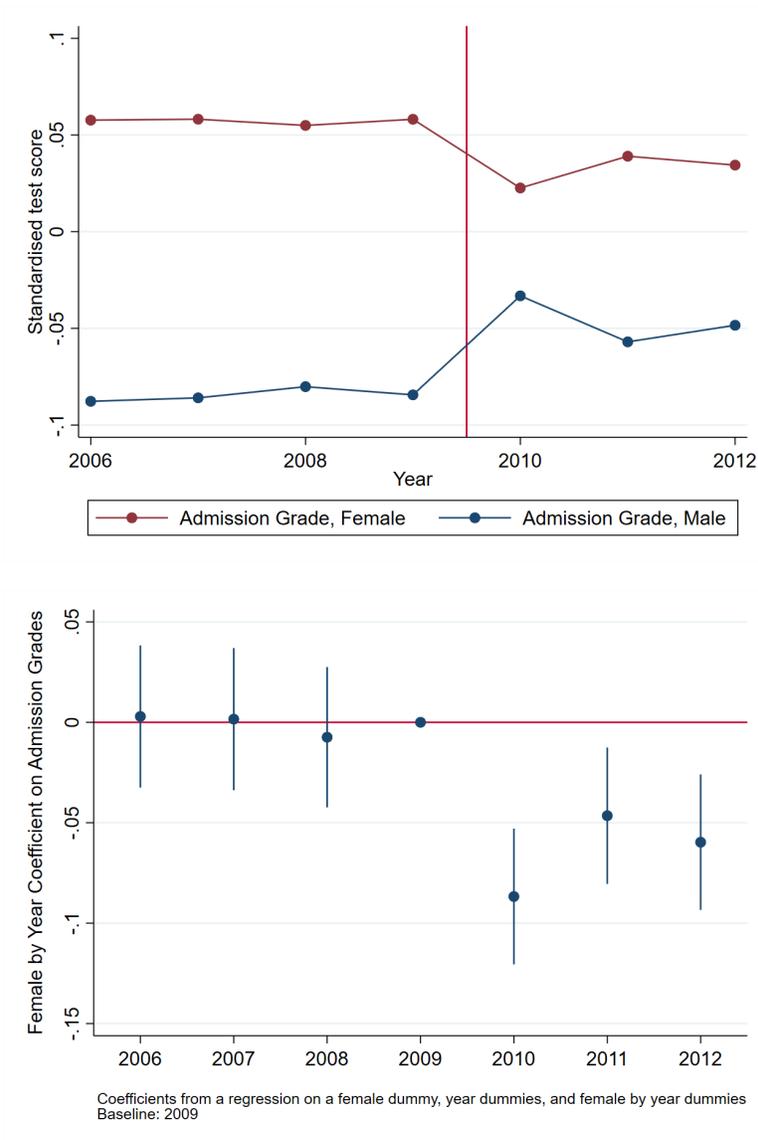


Figure A2

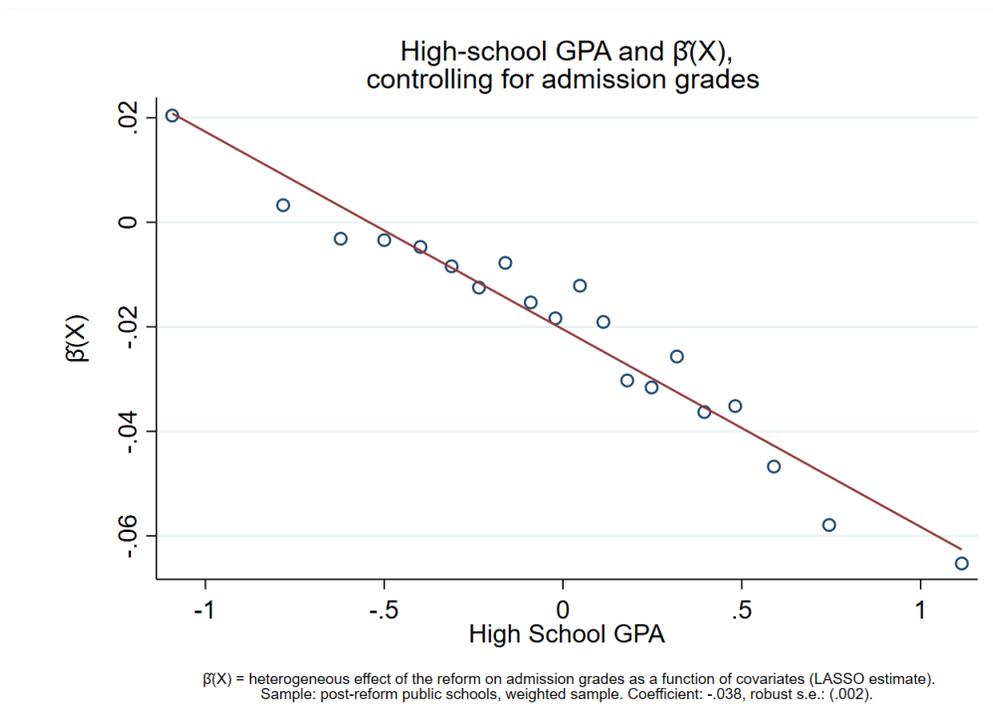


Figure A3: College performance and expected gains from the reform

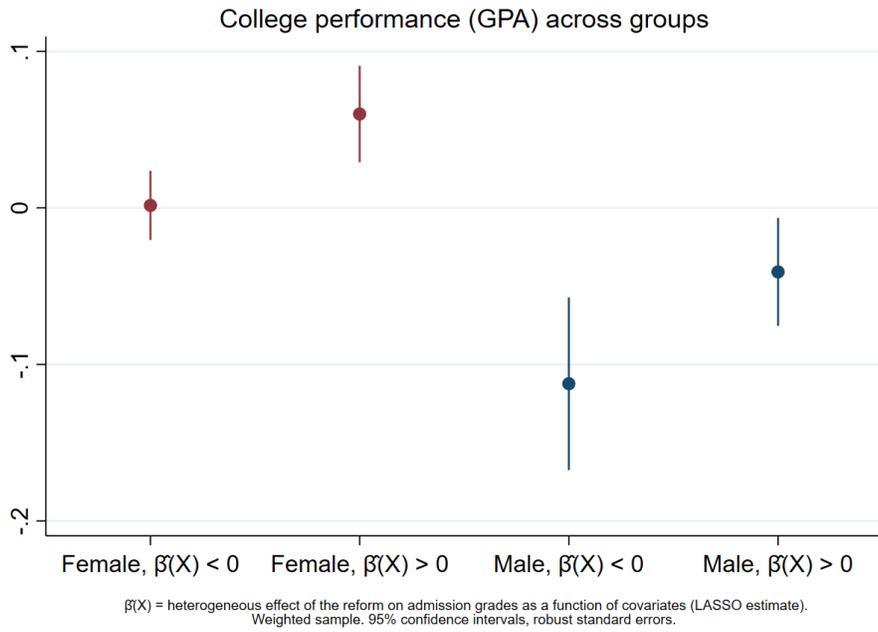
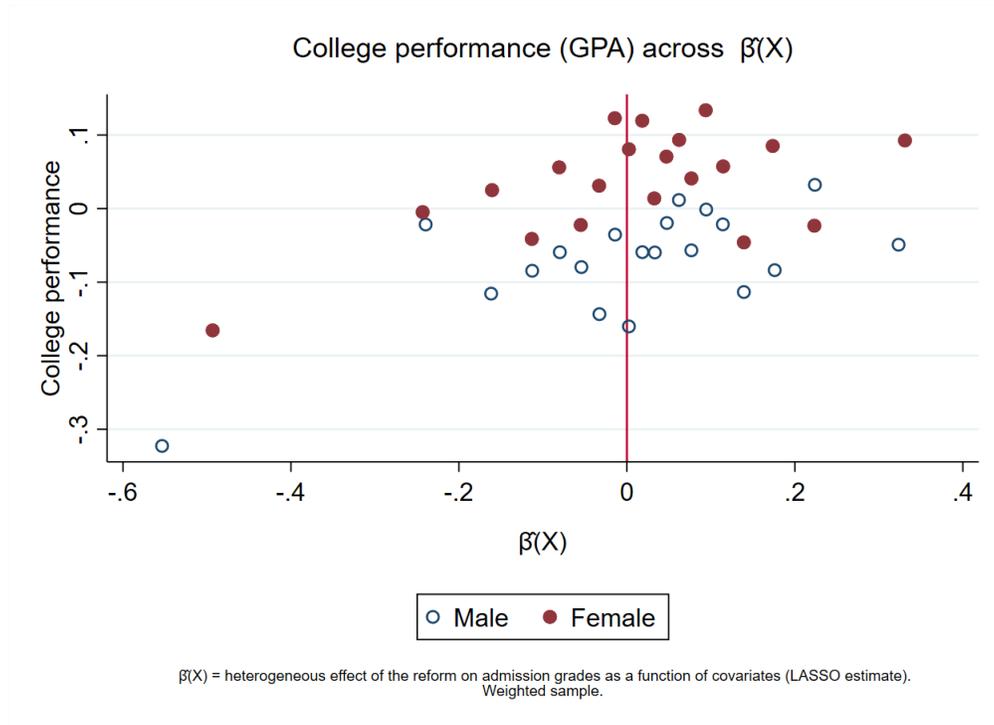


Figure A4: College performance and expected gains from the reform (unweighted)

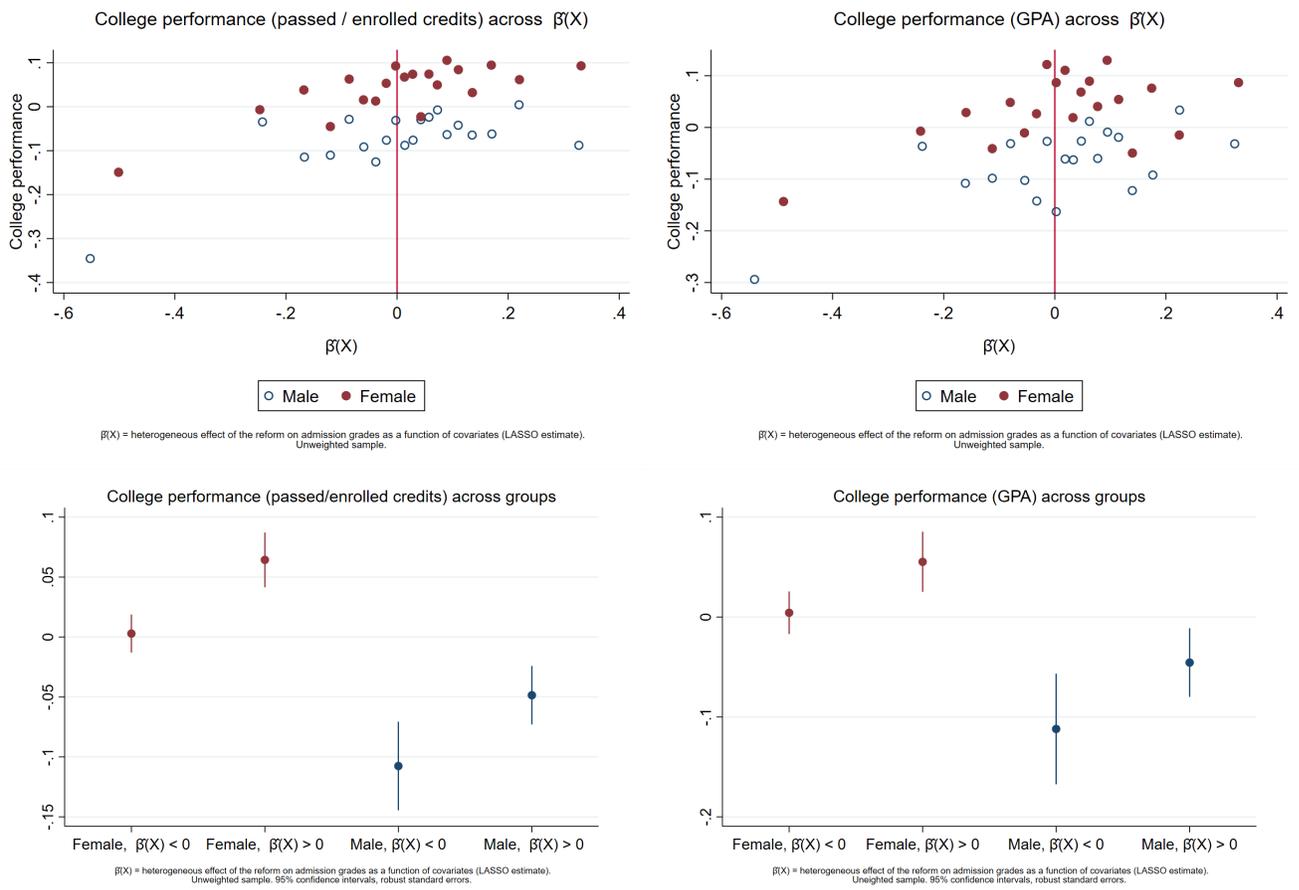
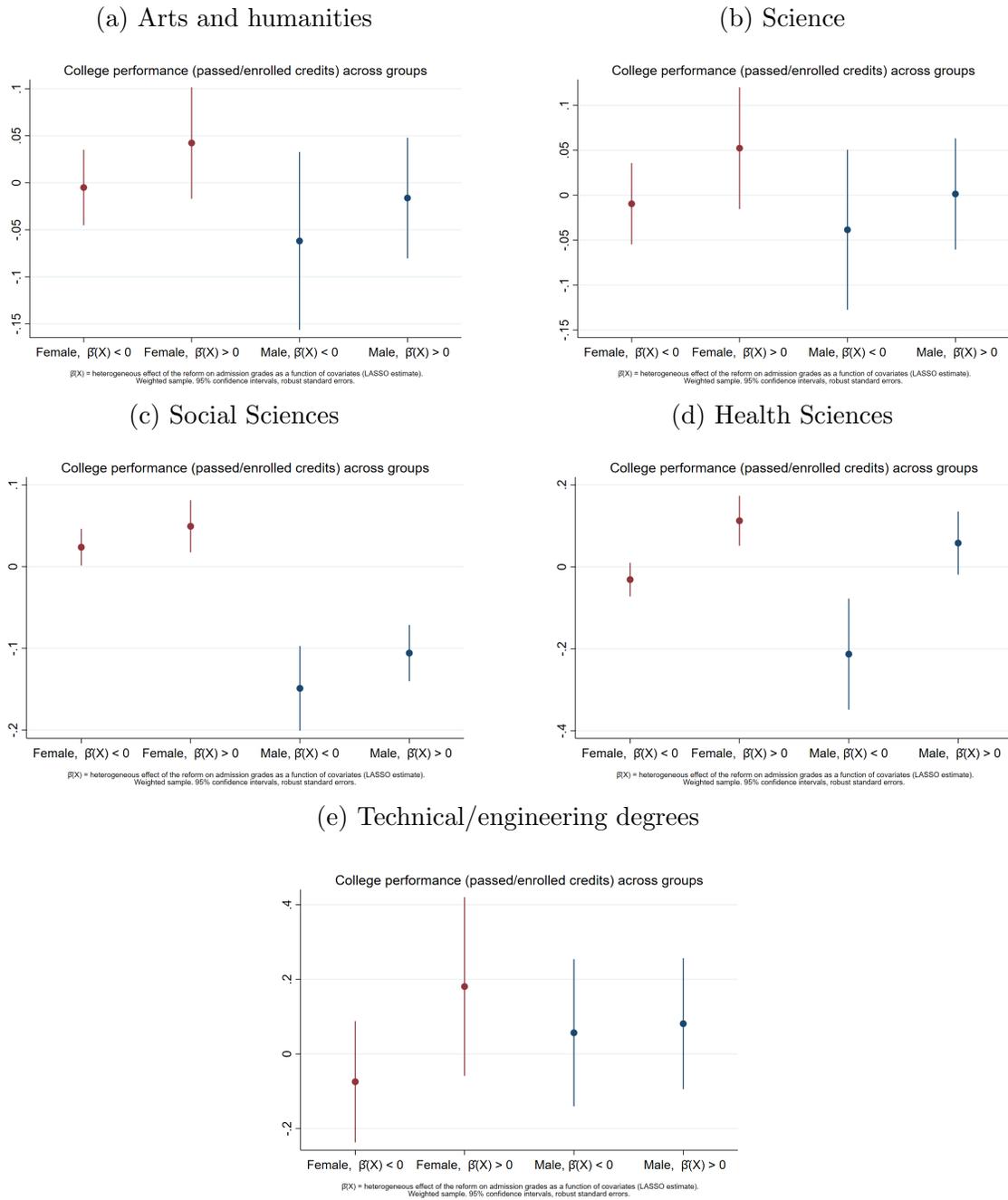


Figure A5: College performance and expected gains from the reform across fields



Appendix B: Career Prospects

In this Appendix, we study how the change in students' allocation changes students' career prospects. To this aim, we use survey data on a sample of pre-treatment college graduates from Catalan universities, with information on labour market outcomes four years after graduation. For every student, we do not observe the exact academic programme, but area (i.e., field of study) indicators, and the university. There are enough area indicators (more than 50) which combined with the university of enrolment make it a meaningful measure, despite some measurement error. Figure A6 displays the social science classification to illustrate the level of detail of the field of study that we observe.²¹

Figure A6: Degree classification example: social science

Catàleg de titulacions

CODI	ENSENYAMENT	SUBÀMBIT DETALLAT (1r NIVELL)		SUBÀMBIT AMPLIAT (2n NIVELL)		ÀMBIT	
2010101	Economia	20101	Economia	201	Economia, Empresa i Turisme	2	Ciències socials i jurídiques
2010102	Comptabilitat i finances	20101	Economia				
2010201	Administració i direcció d'empreses	20102	Administració d'Empreses				
2010202	Màrqueting i investigació de mercats						
2010203	Ciències empresarials						
2010204	Estudis internacionals d'economia i empresa						
2010301	Turisme	20103	Turisme	202	Dret, laboral i polítiques		
2020101	Dret	20201	Dret				
2020201	Criminologia	20202	Laboral				
2020202	Relacions laborals						
2020203	Ciències del treball						
2020204	Prevençió i seguretat integral						
2020301	Gestió i administració pública	20203	Polítiques	203	Comunicació i Documentació		
2020302	Ciències polítiques i de l'administració	20204	Sociologia, Geografia				
2020401	Sociologia						
2020402	Antropologia social i cultural	20301	Comunicació	204	Educació		
2020403	Geografia						
2030101	Comunicació audiovisual	20302	Documentació				
2030102	Periodisme						
2030103	Publicitat i relacions públiques						
2030201	Informació i documentació	20401	Mestres	205	Intervenció Social		
2040101	Educació infantil						
2040102	Educació primària						
2040103	Mestre. Especialitat d'Educació Especial						
2040104	Mestre. Especialitat d'Educació Física						
2040105	Mestre. Especialitat d'Educació Musical						
2040106	Mestre. Especialitat de Llengua Estrangera						
2040201	Pedagogia	20402	Pedagogia i Psicopedagogia	210	Titulacions Mixtes		
2040202	Psicopedagogia						
2040203	Formació de professorat	20501	Treball i educació social				
2050101	Treball social						
2050102	Educació social						
2050201	Psicologia social i organitzacional	20502	Psicologia	210	Titulacions Mixtes		
2100101	Titulacions Mixtes	21001	Titulacions Mixtes				

²¹We do not observe academic programmes (columns 1 and 2), but sub-sub-area indicators (columns 3 and 4).

In this representative survey, although girls outperform boys in educational attainment, females earn 23% less than males on average (9.3% less when accounting for field of study), as reported by table A18.²²

Table A18: Gender wage gap of college graduates

Dependent variable: ln(wage)			
	(1)	(2)	(3)
Female	-0.232*** (0.00915)	-0.0923*** (0.00930)	-0.0939*** (0.00930)
Cohort FE	✓	✓	✓
Field of study FE		✓	✓
University-by-field of study FE			✓
Mean Dep. Var	9.662	9.662	9.662
<i>N</i>	11729	11729	11724

Field of study: sub-sub-area.

Sample of 2006-2009 cohorts, 4 years after graduation.

All regressions control for year of survey FE

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

The left panel of figure A7 shows that there is almost no unconditional correlation between college selectivity and earnings. However, the right panel shows that a positive correlation exists within field of study. This is because selectivity is determined by capacity constraints and demand, and some high-paying technical degrees have good career prospects but low capacity constraints and low demand; while some degrees in the humanities have worse career prospects but high capacity constraints and demand. However, within field of study, where demand and capacity constraints are more homogeneous, the correlation is positive, as one would expect.

Figure A8 displays wages against the gender composition of academic programmes. First, it shows that within academic programmes, females earn lower wages. Second, it also shows that programmes with a higher percentage of female students tend to pay less (for both males and females). This is important because the right panel of figure A8 shows that due to the reform, females enrol less in programmes with a higher pre-reform percentage of female

²²We use the 2014 and 2017 waves of the survey, conducted by the Catalan Agency for the Quality of Universities (AQU), which include students from the 2006-2009 enrolling cohorts.

students.

Figure A7: Threshold grades and wages

Unconditionally

Within field

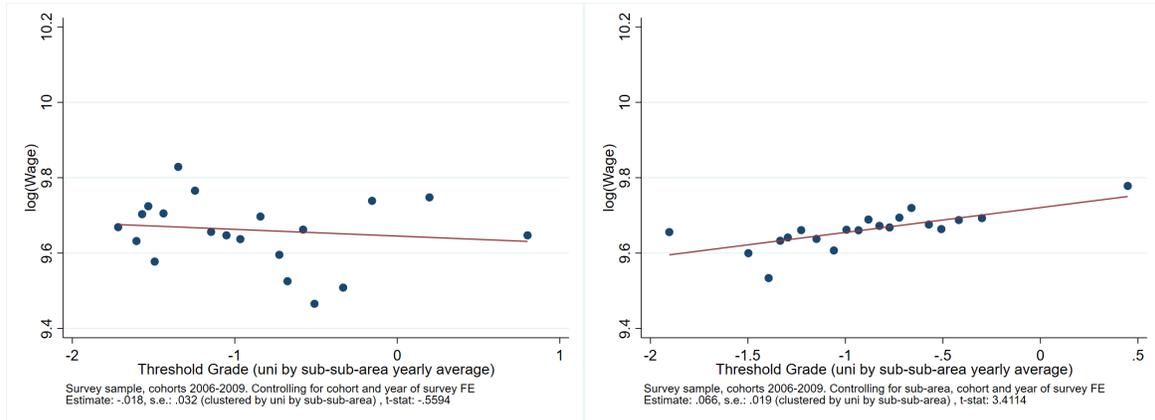
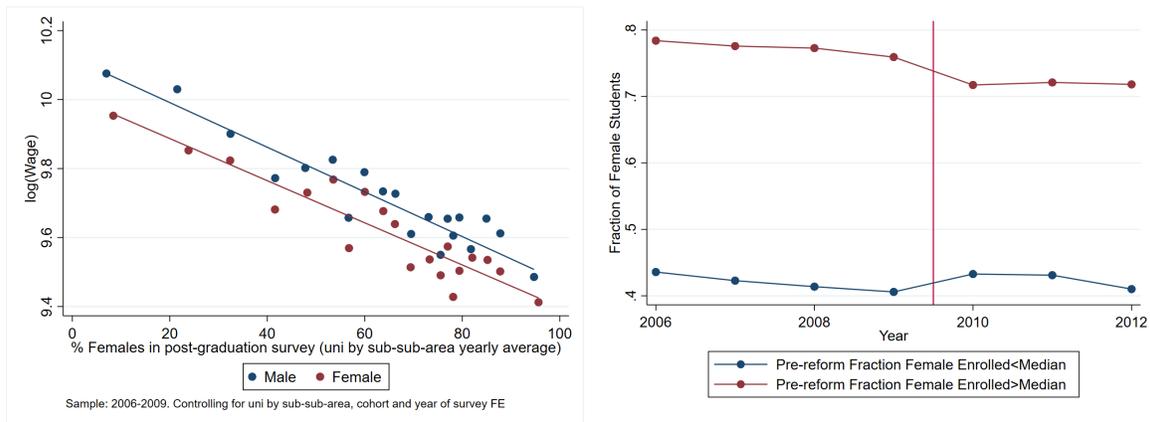


Figure A8

% of female students
and wages

% female enrolment
by pre-treatment female enrolment



Hence, given these patterns of earnings across academic programmes, the effect of the reform on career prospects by gender is not straightforward. To estimate it, we first estimate expected labour market outcomes by academic programme using the survey data, which includes cohorts enrolling into college between 2006 and 2009 (i.e., pre-reform cohorts):

$$Outcome_{it} = \delta Female_i + \alpha(Area \times Uni)_i + \beta Trend_t + \gamma Trend_t \times (Area \times Uni)_i + \epsilon_{it}$$

Where labour market outcomes of student i in enrolling cohort t are measured for the 2006-2009 enrolling cohorts (and survey FE have been partialled out), and where $Area \times Uni$ are dummies for study subarea (or sub-sub-area) by university.

In a 2nd step, we combine the predicted labour market outcomes from the previous regression with the college enrolment data from the Selectivitat dataset, and we estimate:

$$\widehat{Outcome}_{it} = \delta_t + \gamma Female_i + \beta Female_i \times Post_t + \epsilon_{it}$$

Table A19 reports point estimates, indicating an increase of around 2.5pp in the gap, on top of a 9.3pp pre-reform gap within field of study (and an unconditional 23pp pre-reform gap). Table A20 reports point estimates on the expected employment rate. Given the high employment rate among Catalan university graduates (around 87% according to the survey), the magnitude of the effect is smaller, but still significant. To benchmark the magnitude of these effects, it is interesting to compare them with the findings in Ebenstein *et al.* (2016) that pollution in matriculation exam days leads to lower test scores, resulting in a decline in post-secondary education and earnings. It turns out that the effect on female test scores and career prospects is similar in magnitude to the effect of one standard deviation in pollution exposure on the day of the exam. Tables A21 and A22 report placebo tests showing that the change in the post-treatment period is large and significant compared to any changes within the pre-treatment period.

Table A19: Dependent Variable: Predicted log(wage)

	(1)	(2)	(3)	(4)
Female	-0.229*** (0.00167)	-0.224*** (0.00172)	-0.244*** (0.00195)	-0.229*** (0.00203)
Female × Post 2009	-0.0212*** (0.00276)	-0.0140*** (0.00282)	-0.0361*** (0.00341)	-0.0274*** (0.00354)
Year FE	✓	✓	✓	✓
Main Predictor	Sub-field × Uni	Sub-field × Uni × Female	Sub-sub-field × Uni	Sub-sub-field × Uni × Female
Mean Dep. Var	9.684	9.688	9.661	9.665
<i>N</i>	170082	170082	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A20: Dependent Variable: Predicted Employment Rate

	(1)	(2)	(3)	(4)
Female	-0.0109*** (0.000515)	-0.00983*** (0.000600)	-0.0148*** (0.000695)	-0.0134*** (0.000787)
Female × Post 2009	-0.00340*** (0.00104)	-0.00110 (0.00110)	-0.00868*** (0.00169)	-0.00551*** (0.00174)
Year FE	✓	✓	✓	✓
Main Predictor	Sub-field × Uni	Sub-field × Uni × Female	Sub-sub-field × Uni	Sub-sub-field × Uni × Female
Mean Dep. Var	0.873	0.872	0.874	0.872
<i>N</i>	170082	170082	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A21: Dependent Variable: Predicted log(wage)

	(1)	(2)	(3)	(4)
Female	-0.228*** (0.00241)	-0.224*** (0.00249)	-0.243*** (0.00276)	-0.229*** (0.00288)
Female × Post 2009	-0.0203*** (0.00319)	-0.0142*** (0.00326)	-0.0353*** (0.00392)	-0.0276*** (0.00407)
Female × Post 2007	-0.00181 (0.00334)	0.000519 (0.00344)	-0.00152 (0.00390)	0.000392 (0.00406)
Year FE	✓	✓	✓	✓
Main Predictor	Sub-field × Uni	Sub-field × Uni × Female	Sub-sub-field × Uni	Sub-sub-field × Uni × Female
Mean Dep. Var	9.684	9.688	9.661	9.665
<i>N</i>	170082	170082	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Table A22: Dependent Variable: Predicted Employment Rate

	(1)	(2)	(3)	(4)
Female	-0.0117*** (0.000853)	-0.0113*** (0.000986)	-0.0144*** (0.000997)	-0.0137*** (0.00116)
Female × Post 2009	-0.00408*** (0.00109)	-0.00242** (0.00117)	-0.00835*** (0.00182)	-0.00583*** (0.00189)
Female × Post 2007	0.00143 (0.00105)	0.00277** (0.00122)	-0.000702 (0.00139)	0.000686 (0.00158)
Year FE	✓	✓	✓	✓
Main Predictor	Sub-field × Uni	Sub-field × Uni × Female	Sub-sub-field × Uni	Sub-sub-field × Uni × Female
Mean Dep. Var	0.873	0.872	0.874	0.872
<i>N</i>	170082	170082	170082	170082

Robust standard errors in parentheses. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$

Appendix C: exam structure

Table A23: Exam Structure Comparison, sample of subjects

Subject	2009 (pre-reform)	2010 (post-reform)
Catalan	Can choose between two exam models with same structure Part 1: reading. 4 OE questions Part 2: writing. 1 OE question Part 3: reflection on language. 4 OE, 1 MCQ	Can choose between two exam models with same structure Part 1: reading. 4 OE, 2 MCQ Part 2: writing. 3 OE questions Part 3: reflection on language. 2 OE, 2 MCQ
Economics	Can choose between two options with same structure 6 OE questions	Can choose between two options with same structure 6 OE questions
English	Part 1: reading. 1 page text + 8 MCQs Part 2: writing. Choose between 2 topics (100 words) Part 3: listening. 8 MCQs	Part 1: reading. 1 page text + 8 MCQs Part 2: writing. Choose between 2 topics (100 words) Part 3: listening. 8 MCQs
History	Can choose between two options with same structure For each option: text + 5 OE and infographic + 5 OE	Can choose between two options with same structure For each option: text + 5 OE and infographic + 5 OE
Industrial Tech.	Can choose between two options with same structure 5 MCQs + 3 OE exercises	Can choose between two options with same structure 5 MCQs + 3 OE exercises
Latin	Can choose between two options with same structure 3 OE, 1 MCQ	Can choose between two options with same structure 3 OE, 1 MCQ
Mathematics	6 OE questions Must answer 3 out of 4 questions (2 points each), 1 out of 2 questions (4 points each))	6 OE questions (Must answer 5 out of 6 questions (2 points each))
Philosophy	Can choose between two options with same structure For each option: text + 5 OE	Can choose between two options with same structure or each option: text + 5 OE
Physics	Can choose between two options with same structure For each option: 6 OE questions	Can choose between two options with same structure For each option: 5 OE questions
Spanish	Can choose between two options with same structure Part 1: reading. 5 OE, 1 MCQ Part 2: writing. 1 OE question Part 3: reflection on language. 3 OE, 2 MCQ	Can choose between two options with same structure Part 1: reading. 4 OE, 2 MCQs Part 2: writing. 3 OE questions Part 3: reflection on language. 2 OE, 2 MCQ

Note: OE = open-ended question; MCQ = multiple choice question.

Source: <https://www.selecat.cat/>

2019

- 2019/1, Mediavilla, M.; Mancebón, M. J.; Gómez-Sancho, J. M.; Pires Jiménez, L.:** “Bilingual education and school choice: a case study of public secondary schools in the Spanish region of Madrid”
- 2019/2, Brutti, Z.; Montolio, D.:** “Preventing criminal minds: early education access and adult offending behavior”
- 2019/3, Montalvo, J. G.; Piolatto, A.; Raya, J.:** “Transaction-tax evasion in the housing market”
- 2019/4, Durán-Cabré, J.M.; Esteller-Moré, A.; Mas-Montserrat, M.:** “Behavioural responses to the re)introduction of wealth taxes. Evidence from Spain”
- 2019/5, Garcia-López, M.A.; Jofre-Monseny, J.; Martínez Mazza, R.; Segú, M.:** “Do short-term rental platforms affect housing markets? Evidence from Airbnb in Barcelona”
- 2019/6, Domínguez, M.; Montolio, D.:** “Bolstering community ties as a means of reducing crime”
- 2019/7, García-Quevedo, J.; Massa-Camps, X.:** “Why firms invest (or not) in energy efficiency? A review of the econometric evidence”
- 2019/8, Gómez-Fernández, N.; Mediavilla, M.:** “What are the factors that influence the use of ICT in the classroom by teachers? Evidence from a census survey in Madrid”
- 2019/9, Arribas-Bel, D.; Garcia-López, M.A.; Viladecans-Marsal, E.:** “The long-run redistributive power of the net wealth tax”
- 2019/10, Arribas-Bel, D.; Garcia-López, M.A.; Viladecans-Marsal, E.:** “Building(s and) cities: delineating urban areas with a machine learning algorithm”
- 2019/11, Bordignon, M.; Gamalerio, M.; Slerca, E.; Turati, G.:** “Stop invasion! The electoral tipping point in anti-immigrant voting”

2020

- 2020/01, Daniele, G.; Piolatto, A.; Sas, W.:** “Does the winner take it all? Redistributive policies and political extremism”
- 2020/02, Sanz, C.; Solé-Ollé, A.; Sorribas-Navarro, P.:** “Betrayed by the elites: how corruption amplifies the political effects of recessions”
- 2020/03, Farré, L.; Jofre-Monseny, J.; Torrecillas, J.:** “Commuting time and the gender gap in labor market participation”
- 2020/04, Romarri, A.:** “Does the internet change attitudes towards immigrants? Evidence from Spain”
- 2020/05, Magontier, P.:** “Does media coverage affect governments’ preparation for natural disasters?”
- 2020/06, McDougal, T.L.; Montolio, D.; Brauer, J.:** “Modeling the U.S. firearms market: the effects of civilian stocks, crime, legislation, and armed conflict”
- 2020/07, Veneri, P.; Comandon, A.; Garcia-López, M.A.; Daams, M.N.:** “What do divided cities have in common? An international comparison of income segregation”
- 2020/08, Piolatto, A.:** “Information doesn’t want to be free’: informational shocks with anonymous online platforms”
- 2020/09, Marie, O.; Vall Castello, J.:** “If sick-leave becomes more costly, will I go back to work? Could it be too soon?”
- 2020/10, Montolio, D.; Oliveira, C.:** “Law incentives for juvenile recruiting by drug trafficking gangs: empirical evidence from Rio de Janeiro”
- 2020/11, Garcia-López, M.A.; Pasidis, I.; Viladecans-Marsal, E.:** “Congestion in highways when tolls and railroads matter: evidence from European cities”
- 2020/12, Ferraresi, M.; Mazzanti, M.; Mazzarano, M.; Rizzo, L.; Secomandi, R.:** “Political cycles and yardstick competition in the recycling of waste. evidence from Italian provinces”
- 2020/13, Beigelman, M.; Vall Castelló, J.:** “COVID-19 and help-seeking behavior for intimate partner violence victims”
- 2020/14, Martínez-Mazza, R.:** “Mom, Dad: I’m staying” initial labor market conditions, housing markets, and welfare”
- 2020/15, Agrawal, D.; Foremny, D.; Martínez-Toledano, C.:** “*Paraísos fiscales*, wealth taxation, and mobility”
- 2020/16, Garcia-Pérez, J.I.; Serrano-Alarcón, M.; Vall Castelló, J.:** “Long-term unemployment subsidies and middle-age disadvantaged workers’ health”

2021

- 2021/01, Rusteholz, G.; Mediavilla, M.; Pires, L.:** “Impact of bullying on academic performance. A case study for the community of Madrid”

- 2021/02, Amuedo-Dorantes, C.; Rivera-Garrido, N.; Vall Castelló, J.:** “Reforming the provision of cross-border medical care evidence from Spain”
- 2021/03, Domínguez, M.:** “Sweeping up gangs: The effects of tough-on-crime policies from a network approach”
- 2021/04, Arenas, A.; Calsamiglia, C.; Loviglio, A.:** “What is at stake without high-stakes exams? Students’ evaluation and admission to college at the time of COVID-19”
- 2021/05, Armijos Bravo, G.; Vall Castelló, J.:** “Terrorist attacks, Islamophobia and newborns’ health”
- 2021/06, Asensio, J.; Matas, A.:** “The impact of ‘competition for the market’ regulatory designs on intercity bus prices”
- 2021/07, Boffa, F.; Cavalcanti, F.; Piolatto, A.:** “Ignorance is bliss: voter education and alignment in distributive politics”

2022

- 2022/01, Montolio, D.; Piolatto, A.; Salvadori, L.:** “Financing public education when altruistic agents have retirement concerns”
- 2022/02, Jofre-Monseny, J.; Martínez-Mazza, R.; Segú, M.:** “Effectiveness and supply effects of high-coverage rent control policies”
- 2022/03, Arenas, A.; Gortazar, L.:** “Learning loss one year after school closures: evidence from the Basque Country”
- 2022/04, Tassinari, F.:** “Low emission zones and traffic congestion: evidence from Madrid Central”
- 2022/05, Cervini-Plá, M.; Tomàs, M.; Vázquez-Grenno, J.:** “Public transportation, fare policies and tax salience”
- 2022/06, Fernández-Baldor Laporta, P.:** “The short-term impact of the minimum wage on employment: Evidence from Spain”
- 2022/07, Foremny, D.; Sorribas-Navarro, P.; Vall Castelló, J.:** “Income insecurity and mental health in pandemic times”
- 2022/08, Garcia-López, M.A.; Viladecans-Marsal, E.:** “The role of historic amenities in shaping cities”
- 2022/09, Cheshire, P. C., Hilber, C. A. L., Montebruno, P., Sanchis-Guarner, R.:** “(IN)convenient stores? What do policies pushing stores to town centres actually do?”
- 2022/10, Sanchis-Guarner, R.:** “Decomposing the impact of immigration on house prices”

2023

- 2023/01, Garrouste, M., Lafourcade, M.:** “Place-based policies: Opportunity for deprived schools or zone-and-shame effect?”
- 2023/02, Durán-Cabré, J.M., Esteller-Moré A., Rizzo L., Secomandi, R.:** “Fiscal Knowledge and its Impact on Revealed MWTP in COVID times: Evidence from Survey Data”
- 2023/03, Esteller-Moré A., Galmarini U.:** “Optimal tax administration responses to fake mobility and underreporting”
- 2023/04, Armijos Bravo, G., Vall Castelló, J.:** “Job competition in civil servant public examinations and sick leave behavior”
- 2023/05, Buitrago-Mora, D., Garcia-López, M.A.:** “Real estate prices and land use regulations: Evidence from the law of heights in Bogotá”
- 2023/06, Rodríguez-Planas, N., Secor, A.:** “College Students’ Social Capital and their Perceptions of Local and National Cohesion”
- 2023/07, Obaco, M., Davi-Arderius D., Pontarollo, N.:** “Spillover Effects and Regional Determinants in the Ecuadorian Clean-Cooking Program: A Spatiotemporal Econometric Analysis”
- 2023/08, Durán-Cabré, J.M., Esteller-Moré, A., Rizzo, L., Secomandi, R.:** “Has Covid Vaccination Success Increased our Marginal Willingness to Pay Taxes?”
- 2023/09, Borrella-Mas, M.A., Millán-Quijano, J., Terskaya, A.:** “How do Labels and Vouchers Shape Unconditional Cash Transfers? Experimental Evidence from Georgia”
- 2023/10, Messina, J., Sanz-de-Galdeano, A., Terskaya, A.:** “Birds of a Feather Earn Together. Gender and Peer Effects at the Workplace”
- 2023/11, Pelegrín, A., Vidal, Ll., González, I.:** “Diversifying Economic Risks: Japan’s Economic Hedging Towards China”
- 2023/12, Rodríguez-Planas, N., Secor, A., De Balanzó Joue, R.:** “Resilience-thinking Training for College Students: Evidence from a Randomized Trial”

