

ACLA WORKING PAPER SERIES

EDUCATIONAL EQUITY AND TEACHER DISCRETION EFFECTS IN
HIGH STAKES EXAMS

ILJA CORNELISZ
MARTIJN MEETER
CHRIS VAN KLAVEREN

Working Paper 20182

<http://www.acla.amsterdam/workingpapers-wp20182>

AMSTERDAM CENTER FOR LEARNING ANALYTICS
van der Boechorststraat 7
1081 BT, Amsterdam
December 2018

Educational Equity and Teacher Discretion Effects in High Stake Exams*

Ilja Cornelisz
Martijn Meeter
Chris van Klaveren

Abstract

This study examines teacher discretion effects in Dutch secondary education for the period 2007-2012. Stark discontinuities are observed in the exam grade distribution for high-stakes retaking students and are located at important graduation thresholds. This phenomenon is systematically related to the level of discretion when grading the exam, with results suggesting that approximately 12% of all graduating retakers did so because of teacher discretion. This yields unequal graduation opportunities that are the result of school- and subject choice patterns, since teacher discretion is structurally and selectively exerted at the school-level with the objective to let students on the margin graduate.

JEL-code: I2

Keywords: Teacher Discretion, Grading, Equity

1 Introduction

Standardized tests serve to provide objective measures on student performance and these can be high stakes for students as they often determine, at least in part, retention and graduation

*The authors are affiliated with the Amsterdam Center for Learning Analytics (*ACLA*, acla.amsterdam), Faculty of Behavioral and Movement Sciences, VU University Amsterdam, The Netherlands. Ilja Cornelisz (i.cornelisz@vu.nl) and Chris van Klaveren (c.p.b.j.van.klaveren@vu.nl) are the corresponding authors. of the Vrije Universiteit Amsterdam. This authors gratefully acknowledge funding by the Dutch Ministry of Education, Culture and Science (OC&W). We would like to thank Hedvig Horvath, Hessel Oosterbeek for helpful comments. Also, we would like to thank Lianne de Vries for excellent research assistance. All remaining errors are our own.

decisions (Dee et al., 2016). These standardized tests also have become increasingly central to accountability policies with the objective to evaluate, for example, school and teacher performance. The main intent of test-based accountability policies is to provide incentives that maximize student learning, but perverse incentives resulting from poorly designed accountability policies can have significant, unintended and undesirable consequences (Jacob, 2005). The existing widespread concerns over test validity and the manipulation of scores are therefore not surprising (Dee et al., 2016), yet until recently there has been surprisingly little empirical evidence related to test-based accountability and how it may induce manipulation of student test scores (Jacob, 2005).

Two recent empirical evaluations performed in the United States (Dee et al., 2016) and Sweden (Diamond and Persson, 2016) provide strong evidence that allowing for teacher discretion in grading standardized exams gives all the more reason for policy makers to be concerned. Dee et al. (2016) examine the causes and consequences of test score manipulation of high-stakes exit exams for New York State secondary-school students and find that teachers purposefully moved students just over predefined performance thresholds when grading their own students. Moreover, results varied systematically across and within schools and this had heterogeneous implications with respect to subsequent student outcomes. Notably, conditional on scoring near a proficiency cutoff, white and Asian students, students with better baseline scores, and those with good behavioral records are more likely to benefit from such teacher discretion. Diamond and Persson (2016) corroborate the existence of test score manipulation for Swedish compulsory schools, and similarly identify ‘a bad test day’-effect, suggesting that teachers exploit their discretion to undo potentially harmful consequences of idiosyncratic student performance. In contrast to the results in Dee et al. (2016), their estimates do not suggest that test score manipulation is related to student background characteristics. Furthermore, they find relative homogeneous positive implications for subsequent educational, labor market and life outcomes, highlighting that potential signaling mechanisms resulting from graduation could enhance a student’s academic motivation and/or teachers’ perception of academic ability.

This study adds to this emerging body of literature on local grading, teacher discretion and test score manipulation (see, also: Lavy, 2008; Hanna and Linden, 2012; Burgess and Greaves, 2013) by evaluating scores on high-stakes standardized exams at the end of secondary education in the Netherlands. It empirically investigates the existence of teacher discretion in grading and potential consequences for inequalities of student graduation opportunities. A specific contribution of this study is the ability to empirically expose the

underlying dynamics of teacher discretion mechanisms, by exploiting variation in information, stakes *and* teacher discretion opportunities, as to validate a per-student utility model of teacher discretion.

For this purpose, a unique feature of the Dutch exam system is exploited in that subject teachers grade two similar standardized exams of some of their students twice over a short span of time. Yet, these two attempts differ vastly in terms of the *stakes* at hand and the *information* available, which can impact the validity of the observed student performance measure (Neal, 2013) and the associated teacher grading practices (McMillan and Nash, 2000). To be specific, students are allowed to retake *one* exam for one subject which takes place within a week after the results of the first attempt have become known. A retake is often observed if the grade point average (GPA) across all subjects after the first term is insufficient for passing the matriculation examination. The stakes at the retake are thus even higher than in the first-term exams when graduation depends solely on the outcome of this single retake exam. Also, the information to students and teachers between the first and second term is distinctively different. The Dutch Testing Agency (CITO) announces the subject-specific conversion formula used to translate achieved exam points into grades only when the results of the first term are made public. Yet, this same conversion formula then also applies for the subject retake exam that is still yet to be administered and graded. As a result, both teachers and students do not precisely know how many points are needed to pass a particular subject in the first period attempt, but know exactly how many points are needed on the second attempt in order to pass the subject and graduate. We show that when students require a retake exam for graduation, it holds that the optimal strategy of students is to perform as well as they can, and that grade manipulation by means of teacher discretion should explicitly reveal itself in the grading of exams taken in the second term.

For the empirical evaluations, administrative data for the Netherlands is used for the period 2007-2012, covering 99 percent of the Dutch secondary school exam population. This data is augmented with information on the proportion of open questions on the (retake) exam.¹ Since observed improvement gains on the retake exam can potentially be related with teacher discretion, student ability boosting and mean reversion, we use the proportion of open questions as an instrument to identify the effect of teacher discretion. The identifying assumption is that that this measure affects potential teacher discretion (see also Schuurs et al., 2017), as it determines the degree of freedom a teacher has to manipulate grades, but

¹This information is manually obtained from the exam booklets, as uploaded to the ministerial website of the Commission for Tests and Exams (<https://www.examenblad.nl/>)

not students' ability boosting and mean reversion mechanisms.

Next, the grading effect of teacher discretion is estimated by focusing on first attempts observed in the second term. A first attempt in this retake period is only allowed for if students were sick at the time of the first attempt, or when another unanticipated event occurred. Even though these specific first-attempts are considered to be incidental and unforeseen in nature, the information that teachers have at their disposal is markedly different since the formula to convert points scored to grades is known in this second term. Comparing the sample characteristics and grade distribution of students with a first attempt in the first term to those observed with a first attempt in the second term provides an additional test for the existence of teacher discretion effects in a context when potential mean reversion is absent. Finally, this paper examines whether teacher discretion raises concerns of educational inequity with respect to unequal graduation opportunities; both between and within schools. For this purpose, variation in the school-location student proportion that graduated as a result of the retake (i.e. transferal rate) is evaluated. In addition, the aforementioned subject-specific instrument of proportion open-ended questions is exploited to estimate if within-school variation in teacher discretion effects is related with gender and ethnicity.

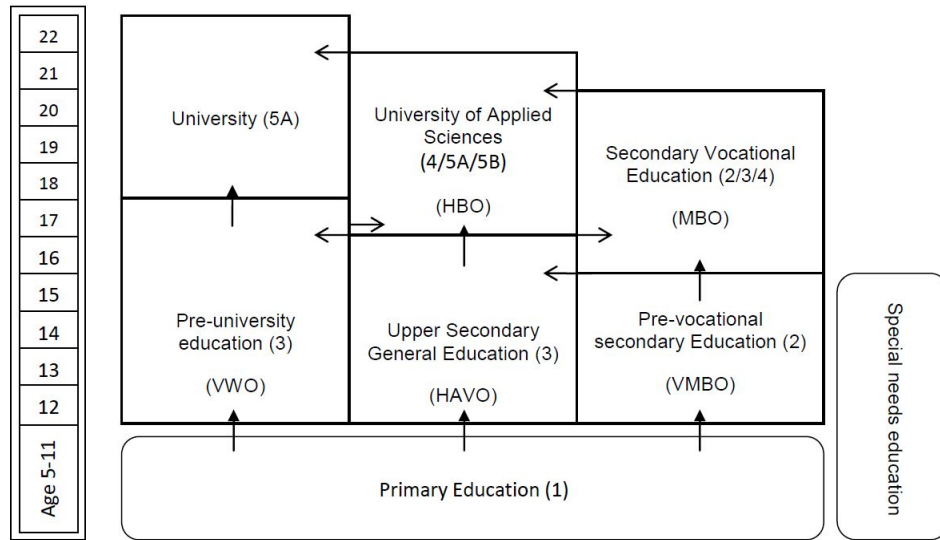
This paper proceeds as follows. Section 2, outlines the Dutch institutional background and explains the exam and grading system in detail. Section 3 introduces a so-called graduation game that emerges in this context based on information and stakes and integrates these insights in a theoretical model for teacher discretion in (Dutch exam) grading. Section 4 reports on the data and descriptive statistics. Section 5 shows the empirical findings, and Section 6 summarizes and provides a discussion of the results and their potential policy implications.

2 Dutch Institutional Background

2.1 Secondary Education in the Netherlands

Upon finishing primary education, children in the Netherlands are tracked into different secondary education levels, with their final track determined after the first or second year of secondary education (Figure 1). The decision to assign students to a particular track at the start of secondary education is based on both a standardized assessment in taken grade 8 and the advice of the primary school teacher. Three distinct tracks in secondary education can be distinguished. Pre-vocational education (4 years) prepares students for vocational education and comprises 4 separate sub-tracks, secondary general education (5

Figure 1. Dutch education system



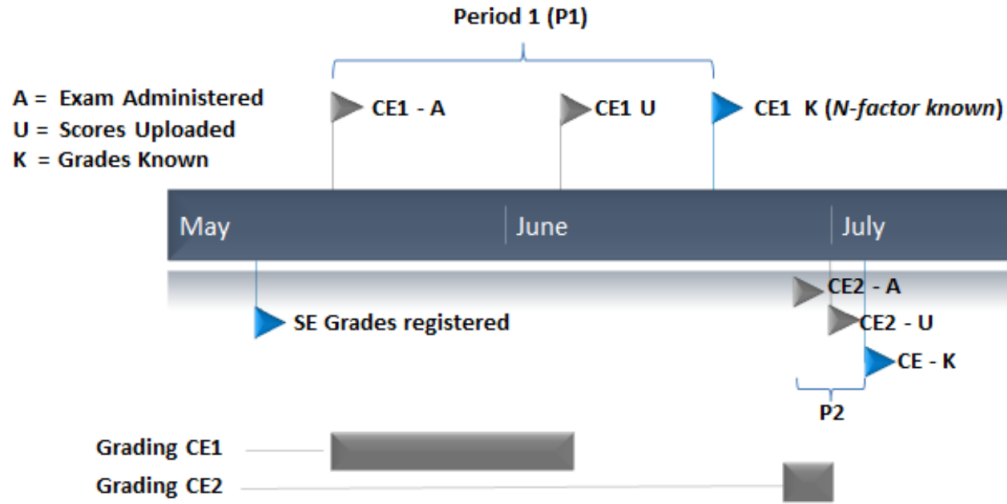
Note: ISCED levels are shown in parenthesis. The Figure is taken from Cornelisz and Van Klaveren(2018).

years) prepares for universities of applied sciences, and pre-university education (6 years) for academic universities. Each track has a matriculation examination in place, in which students take exams in a variety of subjects. This study focuses on students from all three tracks enrolled in the final grade after which they -are expected to- matriculate.

2.2 Secondary School Exams and Grading System

The school-leaving (matriculation) examination for secondary education in the Netherlands consists, for each subject, of a school examination (*SE*) and a national written examination (hereafter referred to as *Central Exam, CE*) at the end of the final school year. Depending on the level of education, students take a CE in roughly 6 to 8 different subjects. The Ministry of Education, Culture and Science prescribes the topics that must be covered in the *SEs* of each subject, but schools have discretion in constructing their own school exams. These school exams usually comprise two or more tests per subject, and can be oral, practical or written. The *CE* for each subject is one test, constructed by the Ministry of Education, Culture and Science, and takes place at a fixed date and time at the end of the final year. The grading scale of each subject is from 1 (lowest) to 10 (highest) and the final grade (grade point average, GPA) for each subject is the arithmetic average of the grades achieved on the

Figure 2. Timeline of Exam Grading Activities



school and the central examination (i.e. $GPA_{s(ubject)} = \frac{1}{2} \cdot SE_s + \frac{1}{2} \cdot CE_s$). A GPA of 5.5 is required to pass a particular subject, but since the school leaving examinations consist of 6 or more subjects, there are explicit rules determining whether a student graduates. The specific graduation rules are outlined in Appendix A, but the main determinant for graduation is whether a student has passed (nearly) all individual subjects.

The Dutch examination system give students the opportunity to retake the *CE*, but for *one* subject only. This retake takes place within a week after the first-term results have become known and the highest score on both attempts is used towards determining whether a student has met the requirements for graduation. The formula determining grade point average per subject is then represented by $GPA_{s(ubject)} = \frac{1}{2} \cdot SE_s + \frac{1}{2} \cdot \max(CE_{s1}, CE_{s2})$. Students who fail to graduate will not be eligible to enroll in the post-secondary education sector their track was preparing them for.

Figure 2 outlines the timing of exam grading activities that are relevant for identifying the possible existence of teacher discretion effects in grading high stakes exams. The figure covers the period from the moment teachers have registered the *SE* grades in a secured digital environment until the moment the grades of the retake exams are publicly announced. Early in May, teachers upload the *SE* grades of their students in *WOLF*, a (web-based) program to exchange exam-related files.² Upon successful uploading, teachers can no longer change

²WOLF is developed by the Dutch National Institute for Educational Measurement (CITO) charged with all logistics surrounding the national examination in secondary education.

the *SE* grades registered. At some point mid-May, the CE exams are administered for all subjects and teachers are provided with explicit and strict guidelines regarding the grading procedure. A student's own subject teacher has two weeks to assign a *score* to each answer based on these guidelines and the assigned scores *per question* are uploaded in *WOLF*. Once the scores are uploaded, teachers can no longer change the assigned scores. The Dutch National Institute for Educational Measurement (CITO) has assigned a teacher from a different school (but same subject) to check and re-mark the work (the so-called second corrector). The second corrector also has two weeks to review the answers and registers any deviating scores, after which (s)he is redirected to a negotiation page. Also for the second corrector it holds that registered deviations cannot be altered once landed on the negotiation page. Deviations are shown on the negotiation page and both correctors then contact each other to reach agreement. Once agreement has been reached, the *second corrector* alters the assigned scores and the first corrector must approve that these revised scores are correctly changed in *WOLF*. Upon approval, the assigned scores are final and stored in *WOLF*. In practice, these final scores are very close to the initial scores (i.e. 0.3% lower on average) as uploaded by the first corrector Kuhlemeier and Kremers (2012).³

Relevant for the possibility to exert teacher discretion is that CITO announces the conversion formula required to translate points to an actual grade only after the final scores of all exam students are stored in *WOLF*. This conversion formula contains a subject-specific factor which varies from year to year as to control for erratic differences in the difficulty of a CE exam. In practice, depending on the level of this factor, this can mean a difference in *CE* grade of (over) 2 points on a scale from 1 to 10. Not knowing this factor thus makes that students and teachers cannot have an accurate prior expectation about the number of points required to (just) pass a particular subject. In the conversion formula, this factor is denoted by N . This N -factor is announced in mid-June, after which students can determine whether they need a retake exam. If a retake exam is required, students will take this exam within a week and the grading procedure will then be similar to the procedure described above. One pivotal difference, given that the N -factor is known at the time of grading, is that teachers now know exactly how many points are needed at the retake exam for the student 'at risk' to graduate.

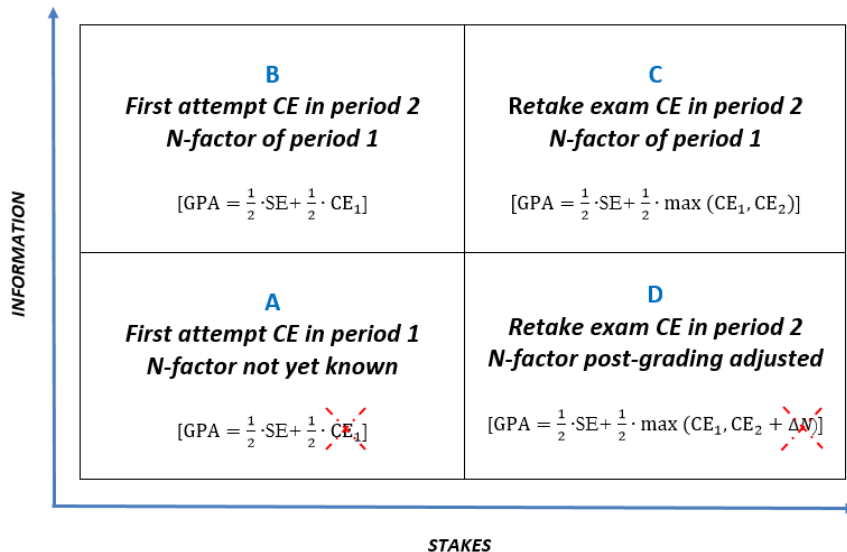
³If both correctors cannot reach agreement, school boards from both schools are asked to mediate in the process. If still no consensus is reached, schools can inform the Educational Inspectorate, who can decide to appoint a third independent corrector. In practice, school boards are asked to mediate in only 1% of the cases and the Educational Inspectorate is asked to intervene in only 0.3% of the cases Kuhlemeier and Kremers (2013)

3 The Graduation Game and a Theoretical Model of Teacher Discretion in Grading

Exams are graded in vastly different contexts in terms of both the stakes at hand and the information available. These different grading environments create a setting we refer to as the graduation game, illustrated in Figure 3. Distinguishing between these different grading contexts gives insights into when teachers are will exert their discretion in grading an exam as to enable student to -just- graduate. Figure 3 depicts four possible situations (*A*, *B*, *C* and *D*) that occur throughout the exam period and the corresponding *GPA* formula in each situation. This formula represents not only how the *GPA* of the student can be calculated, but also what part of that formula is (un)observed by the teacher *when grading*. In situation *A*, teachers do not observe the *N*-factor and, therefore, do not observe CE_1 and cannot precisely determine *GPA* when grading. Moreover, teachers are also unaware of the *GPA* for any of the other subjects part of a student's school-leaving examination, making it relatively difficult to pinpoint exactly which students are at risk of not graduating. Yet, in situation *B* (i.e. a student has a first attempt for a subject in the second term), the situation is markedly different. In this case, teachers and students are aware of the *N*-factor and of the grades achieved through first attempts on other subjects in the first term. Relative to situation *A*, this additional piece of information enable the subject-teacher in charge of grading the exam to figure out whether a score does or does not yield graduation. This situation somewhat resembles the situation of the retake exam (i.e. situation *C*) in terms of the amount of information available, but the stakes in situation *B* are lower, because students can still make use of the retake opportunity⁴. Situation *C* represents the standard retake exam opportunity observed in the second term and, if a student failed to matriculate on the basis of the first term results, both the student and the teacher know that graduation depends solely on the result of this retake exam (i.e. CE_2). One could say that this situation is similar than situation *B*, but that the urgency (stakes) for exerting teacher discretion and/or manipulating scores as to benefit students is higher. Situation *D* refers to a specific situation in which the *N*-factor of the retake exam is adjusted afterwards. Importantly, the *N*-factor can only be adjusted upwards, such that the initial *N*-factor is always a correct lower bound of the final grade.

⁴Given the incidental nature of situation *B*, the schedule for this "period 3" examinations is determined on a post-hoc basis. These additional retake exams usually take place early August.

Figure 3. The Graduation Game



A Model of Teacher Discretion in Grading

To gain more insight in the workings of teacher discretion in grading, we integrate the aforementioned graduation game in the model of test score manipulation by Diamond and Persson (2016). The formula used to determine the CE -grade can be represented by:⁵

$$CE_i(S_i, N) = (10 - N) \cdot \frac{S_i}{S_{total}} + N, \quad (1)$$

where S_{total} represents the maximum number of points that can be achieved, S_i the actual number of points achieved by student i , and N the N -factor that can take on any value between 0 and 2. The formula highlights that CE is determined by both the N -factor and the proportion of points obtained on the exam. In the absence of teacher discretion in

⁵We note that in reality the formula is $CE_i = 9/(S_{total}/S_i) + N$. This is a rather inconvenient formula because the N -factor can take any value between 0 and 2 such that CE may be (lower) higher than the (minimum) maximum CE of (1) 10 points when the student answered all exam question (incorrectly) correctly. The examination board, which is a ministerial but independent organization that has the responsibility that the quality and the logistics are guaranteed (see <https://www.cvte.nl/>) therefore formulated the so-called border-relationships. These border-relationships are grade corrections for when the N -factor is unequal to 1. Information on the exact standardization of the exams can be found at <http://wetten.overheid.nl/BWBR0010538/1999-07-07>. For the developed model in this section it is important that CE is a continuous monotonic increasing function of S_i and N and the alternative presented formula for determining CE is therefore convenient and without loss of generality.

grading, equation (1) can be rewritten as:

$$CE_i(S_i(SE_i, \varepsilon_i), N) = (10 - N) \cdot \frac{S_i(a_i, \varepsilon_i)}{S_{total}} + N. \quad (2)$$

The achieved number of points on the CE then depend on student ability, a_i , and an error term, ε_i , which captures the fact that the performance on the central exam may deviate from the true ability of the student. As such, students can have a good test day ($S_i > a_i$) or a bad test day ($S_i < a_i$), reflecting idiosyncratic performance differences with respect to a student's (unobserved) "true" ability (cf. Diamond and Persson (2016)).

Students will only pass a specific subject if $CE \geq 11 - SE$ and teachers can use this information to determine whether they exert discretion in grading (e.g. to ensure a student graduates), thereby awarding additional points to student i , indicated by Δ_i .⁶ When we combine this passing rule with Equation 2 -and rearrange terms-, the threshold of total points (i.e. $S_i + \Delta_i$) required to pass a subject is given by:

$$S_i(a_i, \varepsilon_i) + \Delta_i = \frac{S_{total}(11 - SE - N)}{10 - N} \quad (3)$$

The left-hand side of Equation 3 indicates that the total points obtained is the sum of points achieved by the student ($S_i(a_i, \varepsilon_i)$) and any additional points awarded by the teacher by means of exerting discretion (Δ_i). The right-hand side of Equation 3 highlights that the required total points to pass the exam is conditional on the N -factor and the already registered SE -grade. We refer to this right-hand side as κ , or $\kappa(N)$, to illustrate that the teacher knows how many points is required only if the N -factor is known.

Based on Equation 3 and threshold κ , we can then define the following indicator function:

$$t_i = t(a_i, \varepsilon_i, \Delta_i, N) = \left\{ \begin{array}{ll} 1 & \text{if } S_i(a_i, \varepsilon_i) + \Delta_i \geq \kappa(N) \\ 0 & \text{Otherwise} \end{array} \right\}. \quad (4)$$

Equation 4 reflects that student i passes the exam for a certain subject if the obtained points $S_i(a_i, \varepsilon_i) + \Delta_i$ surpasses threshold $\kappa(N)$. A distinct difference between the model of Diamond and Persson (2016) is that our model is not a *full information model* in the sense that teachers do not have full information regarding the threshold value if the N -factor has not yet been announced (i.e. situation A). This means that teachers can only effectively target an artificial increase in CE grade - exerting their discretion in grading- if the N -

⁶This can be easily seen, because the subject-specific $GPA = \frac{1}{2} \cdot SE + \frac{1}{2} \cdot CE$ and students pass the subject if $GPA \geq 5.50$. When we substitute $GPA = 5.50$ and rearrange terms we obtain the rule that student pass if $CE \geq 11 - SE$.

factor is known. Equation 4 also captures that teachers will only increase the test score of a student with Δ if this will result in graduation.

Assume that student i is taught by teacher j . When teacher j is manipulating student i 's test score effectively, Δ_i is chosen such that the utility function of each student is maximized:

$$\begin{aligned} u_{ij}(\Delta_i) &= \beta_{ij}t_i(SE_i, \varepsilon_i, \Delta_i, N) - c_{ij}(\Delta_i), \\ c'_{ij}(\Delta_i) &> 0, \quad c''_{ij}(\Delta_i) > 0 \end{aligned} \tag{5}$$

Parameter β_{ij} reflects teacher j 's student-specific desire to raise student i 's grade from a fail to a pass, or as Diamond and Persson (2016) remark β_{ij} "... permits the teacher to use her discretion both in a "corrective" and "discriminatory" fashion" (p.12). This distinction is important, as the term corrective refers to teachers who may have a preference for compensating a bad test day, while the term discriminatory refers to teachers who may have a preference for increasing test scores of students with certain (background) characteristics. Increasing the points obtained by a student with Δ_i comes at a cost (i.e. $c_{ij}(\Delta_i)$) and these costs are assumed to be strictly increasing and convex. This implies that it becomes increasingly difficult for the teacher to award additional points by means of exerting discretion, given that (1) exams only have a limited set of points that are subjective to teacher discretion, (2) additional points would require rewarding answers that are clearly wrong, and (3) teachers (schools) have to justify their grading results to the second corrector (educational inspectorate). Teachers thus optimally exert their discretion up to the point where the marginal benefits of doing so just offset the marginal costs:

$$\frac{\partial u_{ij}}{\partial \Delta_i} = 0 \implies \beta_{ij} \frac{\partial t_i}{\partial \Delta_i} = \frac{\partial c_{ij}}{\partial \Delta_i} > 0 \tag{6}$$

From the model above, a number of empirical hypotheses are derived. Equation 6 illustrates that teachers will only engage in the costly exertion of teacher discretion if (s)he has a positive student-specific desire to do so ($\beta_{ij} > 0$) and when doing so alters the student's subject grade from fail to pass. Furthermore, a teacher would not add points beyond the passing threshold, given that it's costly to do and does not further alter the grade in terms of pass or fail status. As such, if teacher discretion is exerted while grading, this will be targeted effectively and should emerge as a discontinuity in the test score distribution centered around the subject-specific pass-fail threshold. Next, given that knowing exactly the threshold level of points required is contingent on observing the N-factor (i.e. $\kappa(N)$), and given the positive, increasing, convex cost function of adding points by means of teacher discretion ($c_{ij}(\Delta_i) > 0$, $c'_{ij}(\Delta_i) > 0$, $c''_{ij}(\Delta_i) > 0$), test score manipulation will be partic-

ularly observed when information available and stakes at hand are both high (i.e. retake attempts of high-risk students). Lastly, the magnitude of points added by the teacher (Δ_i) is conjectured to be increasing in a teacher’s student-specific desire (β_{ij}) to engage in this behavior and decreasing in the costs associated with it ($c_{ij}(\Delta_i)$); the latter suggesting that test score manipulation will be more prevalent when the potential to exert teacher discretion is relatively high.

4 Data and Descriptive Statistics

This study uses student-level administrative data on 1.12 million students who are in their final secondary school year in the period 2007-2012. The data contains information on students enrolled in publicly-funded schools, covering 99% of the exam student population. For each student, a list of background characteristics is known, together with the results on school examinations and central exams (for all subjects and both terms). Information about the N -factor was derived from the ministerial website of the Commission for Tests and Exams (<https://www.examenblad.nl/>). The average N -factor in period 1 was 0.95 (SD = 0.45) and the N -term adjustment in the retake was 0.29 (SD = 0.19).

Table 1 compares student characteristics between the full student population and the population of students who make use of their retake opportunity. Among the population of retakers, students who required a retake in order to graduate are labeled ‘high-stakes’. The average student is around 16 years old and has achieved a school exam grade of 6.52, on a 1 (lowest) to 10 (highest) scale. The proportion of students with a migrant background or living in an impoverished neighborhood is, respectively, 0.20 and 0.13. The proportions related to education-level show that most students are enrolled in pre-vocational education (55 percent) and the least students are enrolled in pre-university education (17 percent). Students who used their retake opportunity have, on average, lower school exam grades, are somewhat older and more frequently have a migrant background or live in an impoverished neighborhood. Also they are more frequently enrolled in upper general or pre-university education. Whether students required a retake in order to graduate (i.e. the retake is high stakes) is reflected in the lower achieved SE -grade of 5.74. Even though these students have scored a lower grade on their school examination, they have otherwise rather similar characteristics to other retakers, except that upper-general education is relatively over-represented.

Table 1. Comparing Characteristics of Full Population with Retakers

	Full Population		Retake Population			
	Mean	SD	All		High Stakes	
			Mean	SD	Mean	SD
Male	0.50	0.50	0.47	0.50	0.48	0.50
Age on October 1 st	16.08	0.96	16.25	1.00	16.37	1.02
Non-Dutch Background	0.20	0.40	0.30	0.46	0.32	0.46
Impoverished Neighborhood	0.13	0.34	0.18	0.39	0.19	0.40
<i>SE</i>	6.52	0.84	6.04	0.86	5.74	0.69
Pre-vocational education	0.55	0.50	0.48	0.52	0.43	0.50
Upper general education	0.26	0.44	0.29	0.45	0.34	0.47
Pre-university education	0.17	0.38	0.23	0.42	0.23	0.43
<i>N</i>	1,118,650		253,796		136,638	

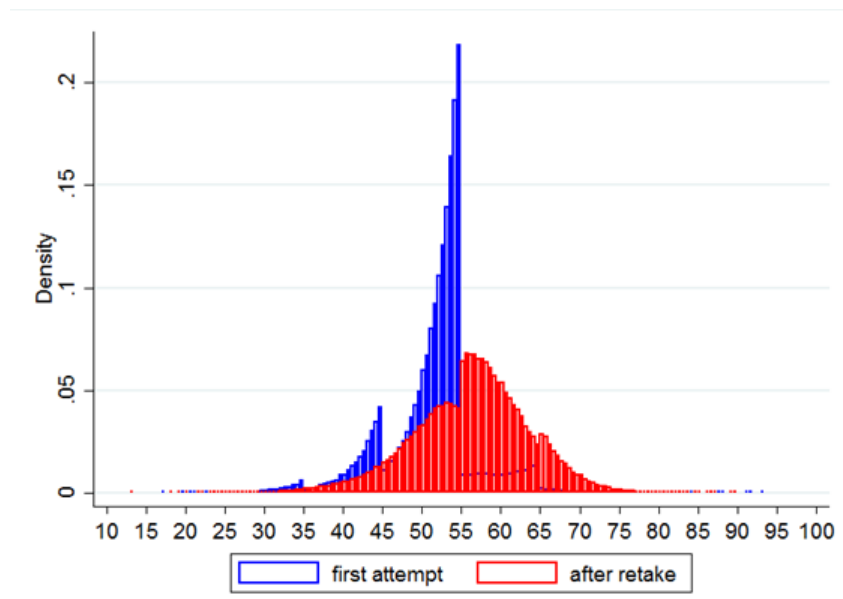
5 Findings

5.1 Teacher discretion effects for marginal students in retake exams

The theoretical model predicts that teachers will (primarily) exert their discretion when the information available and stakes at hand are both high and *only* if doing so would lead to graduation. In total, 253,796 students make use of their retake opportunity, but for 136,638 (53.8%) of them required a retake in order to graduate (i.e. high-stakes retakers). Figure 6 shows for high-stakes retakers the final grades distributions based on only the first attempt (i.e. CE_1) and based on the highest achieved grade achieved in the first attempt and the retake (i.e. $\max(CE_1, CE_2)$). Of these students, 108,972 students (79.8%) experienced a GPA gain by means of the retake exam and 69,280 students graduate as a result of this GPA gain (50.7%). The figure shows a large and significant discontinuity at the passing threshold of 5.5, indicating that a substantial fraction of students is transferred from the left to the right of the passing threshold. There are three reasons that could cause such a transfer:

1. ability boosting, in that students put in a lot of effort in preparing for the retake as to improve their performance,
2. mean reversion, in that students performed relatively low on the initial attempt ($S_i < a_i$),
3. teacher discretion, in that teachers exert their discretion to award additional points on the retake exam ($\Delta_i > 0$),

Figure 4. Grade distributions based on first and retake attempts



N= 136698

Discontinuity= 0.439***

SE = 0.017

To distinguish teacher discretion effects from the effects of ability boosting and mean reversion, we exploit the variation in open questions (i.e. non-multiple choice questions). Given that both the first attempt and the retake exam have the same proportion of open questions and are equally difficult⁷, the identifying assumption is that the proportion open questions is unrelated with ability boosting by students and mean reversion, but can be positively related to teacher discretion. Information on the nature of (retake) exams was obtained from the ministerial website of the Commission for Tests and Exams (<https://www.examenblad.nl/>), as to determine the proportion of open questions.

Mean reversion departs from the recognition that the grade of a student is drawn from his/her own (normal) grade distribution, and when the first attempt produced a grade low in this distribution ($S_i < a_i$), the probability that the retake produces a higher grade than the grade achieved in the first attempt is relatively high. Mean reversion can be considered a ‘bad test day’-effect and cannot be structurally related to the proportion open questions

⁷We note that both the first attempt and the retake exam are constructed and validated at the same moment.

on the exam.

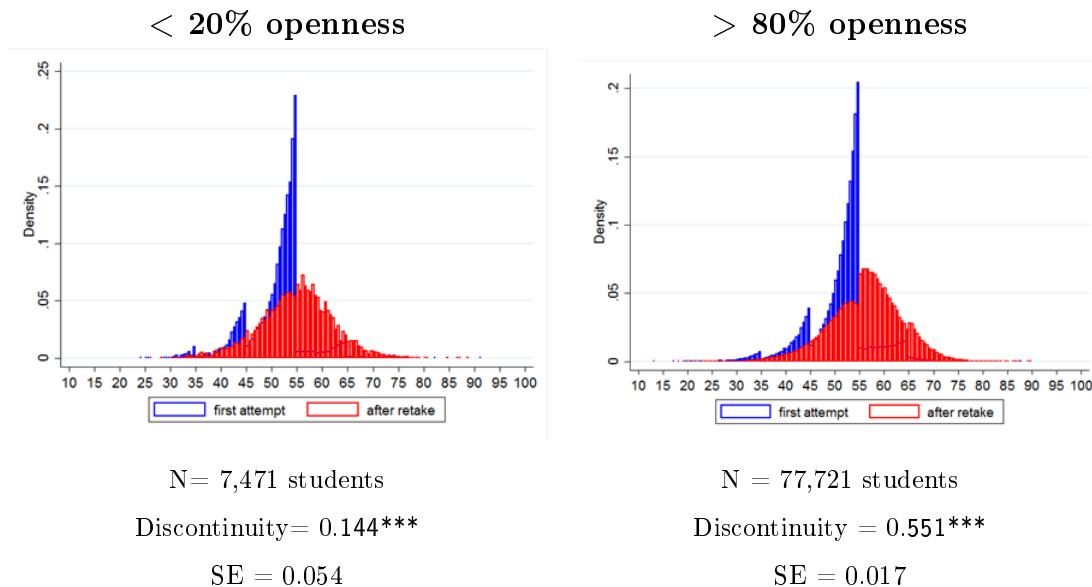
Ability boosting is the process in which students (temporarily) raise their performance level with the objective to graduate. The optimal strategy for high-stakes students is to boost their ability and score as many points as possible on the retake exam, thereby maximizing the probability of graduating and minimizing the uncertainty with respect to whether the scored points are indeed enough for graduation. Since students are able to convert scored points on the retake exam to grades -as the N -factor is available to them- it can be argued that it is optimal to achieve only exactly the required number of points needed for graduation, since this maximizes graduation at minimal effort. However, students do not know how much effort is required and cannot exactly convert points on the retake exam to grades when the proportion of open questions is high, which further exemplifies that it is optimal for students to boost their ability as much as possible and to score as many points as possible on the retake exam.

Finally, the theoretical model of teacher discretion in grading (Section 3) shows that teacher manipulation always results in an improved grade (i.e. $\Delta_i > 0$) and is applied only when when it leads to graduation (i.e. $t_i > 0$). Moreover, the model shows that positive grade manipulation comes at a cost, because exams have a restricted set of open questions that are subjective to teacher discretion. When all exam questions consist of multiple choice questions, then there is no set of points at all that is subjective to teacher discretion. Instead, when there are many open questions, teachers have a larger set of questions to exploit their discretion and increase the grade. Thus, when teachers structurally manipulate exam grades in the retake period, the proportion of open questions of the retake exam should correlate positively with both grade gains and the propensity to graduate. As such, this proportion open questions (or POQ) is used as an instrument for identifying teacher discretion effects. In Appendix B, it is shown that for the 386 subject-year clusters considered in the sample of graduating retakers, the proportion of open questions (Figure B1) is negatively skewed ($M=0.72$ and $SD=0.28$) and -in itself- does not predict (Table B1) first-attempt performance on the central exam (CE_1), supporting its validity as an instrument to identify teacher discretion effects on the retake exam (i.e. CE_2).

Figure 5 illustrates the relationship between the proportion of open questions and GPA gains that students achieve by doing a retake exam. The left panel shows the GPA gains when the proportion of open questions is less than 20 percent, while the right panel shows the GPA gains when the proportion of open questions is more than 80 percent. Both panels show that students are transferred to the right of the passing threshold, but McCrary density

tests (McCrary, 2008) show that the observed discontinuity at the passing threshold (i.e. $GPA = 55$) is much larger when the proportion of open questions is larger (i.e. 0.551 versus 0.144).

Figure 5. Grade distributions based on first and second attempt - by exam question openness



Since the 20 and 80 percent thresholds are chosen arbitrarily, a regression model is estimated in which central exam achievement gains are related to the proportion open questions. The theoretical model shows that in the presence of teacher discretion we must observe central exam improvement gains ($\Delta_i > 0$) and graduation ($t_i > 0$), which is why the Table 2 estimates include only the 69,280 students who graduated because of the retake exam. Baseline model 1 includes not only the proportion open question, but also CE -grade in term 1 and the interaction between both. This is a necessary addition, because the teacher will never use his or her discretion if graduation has already been achieved by means of ability boosting and/or mean reversion. The higher the CE -grade in the first term, the more likely it becomes that the student graduates as a result of ability boosting and mean reversion, and not because of teacher discretion. It follows that the teacher discretion effect measured by the instrument proportion open questions (POQ) can be correlated with CE -grade and by including CE -grade and the interaction term between CE -grade and the proportion of open questions, we ensure that the parameter that is associated with POQ captures only

the effect of teacher discretion.

Model 1 shows that the estimated coefficient for *POQ* when controlling only for *CE*-grade in term 1 is 4.8 and highly significant. This parameter estimate indicates that, on average, the central exam gains are 0.48 points higher (on the final 1-10 scale) in the presence of a teacher discretion effect. Models 2, 3, and 4 show that this estimate remains unchanged when *SE*-grade, student level controls, and education level, year and school location fixed effects are included in the model. This confirms the notion that the teacher discretion effect estimate is independent of mean reversion and ability boosting and provides evidence that teachers exploit their discretion to let students graduate who would otherwise not have graduated.

To get an idea of the number of students this concerns, the observed central exam improvement gains for individual graduating retaking students are "corrected" for teacher discretion, using the proportion of open questions on the retake exam and the -95% confidence interval of the *POQ* parameter estimate obtained in Table 2. Using these simulated gains, subject-specific final GPA scores are re-calculated and for each student it is determined whether or not (s)he would have still passed the pivotal graduating threshold. The result of translating this result from Table 2 to graduation effects returns a range of 6.5-16.5 percent (with a mean of 11.8 percent). This suggests that approximately 12% of all students who graduated by means of a retake did so because of teacher discretion. On a yearly basis, this translates to roughly 1,400 students who are transferred to graduation as a result of a teacher exerting discretion while grading.

5.2 *Teacher discretion effects and selective participation in first exam attempts*

For the vast majority of students ($N=1,111,971$) the first attempts are observed in the first term (i.e. situation A), when the *N*-factor required for points-to-grade conversion is unknown. First attempts in the second term (i.e. situation B) are arguably incidental and random (e.g. sickness), such that average student characteristics related to exam performance should not correlate with the incidence of observing a first attempt in the second term. Although this situation indeed occurs for a relatively small group of students ($N=6,679$), it is an interesting situation because in term 2 the *N*-factor is known to both teachers and students such that exam scores can be precisely converted to a grade. Furthermore, mean reversion cannot confound the interpretation when comparing first-attempt subject-exam grade distributions across situations A and B. Lastly, we can compare first attempt results for the same group of students, thereby exploiting variation in whether the *N*-factor is known

Table 2. *CE* improvement gains and proportion of open-ended exam questions for graduating marginal retakers

	<i>Dependent variable: Central Exam Improvement Gains</i>							
	1		2		3		4	
	Coeff	SE	Coeff	SE	Coeff	SE	Coeff	SE
Prop. Open Questions (POQ)	4.827***	(1.100)	4.553**	(1.023)	4.562***	(0.962)	4.511***	(0.972)
CE grade - attempt 1 · POQ	-0.048***	(0.016)	-0.023	(0.015)	-0.008	(0.014)	-0.007	(0.014)
CE grade - attempt 1	-0.482***	(0.014)	-0.508***	(0.012)	-0.514***	(0.011)	-0.514***	(0.012)
SE grade			0.060***	(0.016)	0.074***	(0.015)	0.080***	(0.014)
Constant	35.53***	(0.943)	46.50***	(2.433)	29.49***	(1.497)	29.63***	(1.918)
Student-level controls	<i>No</i>		<i>Yes</i>		<i>Yes</i>		<i>Yes</i>	
Education Level dummies	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
Year dummies	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
School Location Fixed Effects	<i>No</i>		<i>No</i>		<i>No</i>		<i>Yes</i>	
N	69280		69280		69280		69280	
R^2	0.214		0.236		0.245		0.275	
# subject-year clusters	386		386		386		386	

Note: Robust standard errors in parentheses. */**/** denote significance at a 10/5/1 percent confidence level. Outcome variable in Models 1-4 is the observed improvement in CE-grade for a student on the retake exam. Standard errors are clustered at the subject-year level. Student-level controls are: boy, age, non-Dutch background, impoverished neighborhood, SE-grade and a dummies for whether a covariate is missing. Education level dummies are: pre-vocational education, upper general education.

(i.e. situation A versus B), variation in stakes at hand (i.e. all students versus high-stakes retakers), and variation in teacher discretion (i.e. all subjects versus subjects with at least 80% open questions).

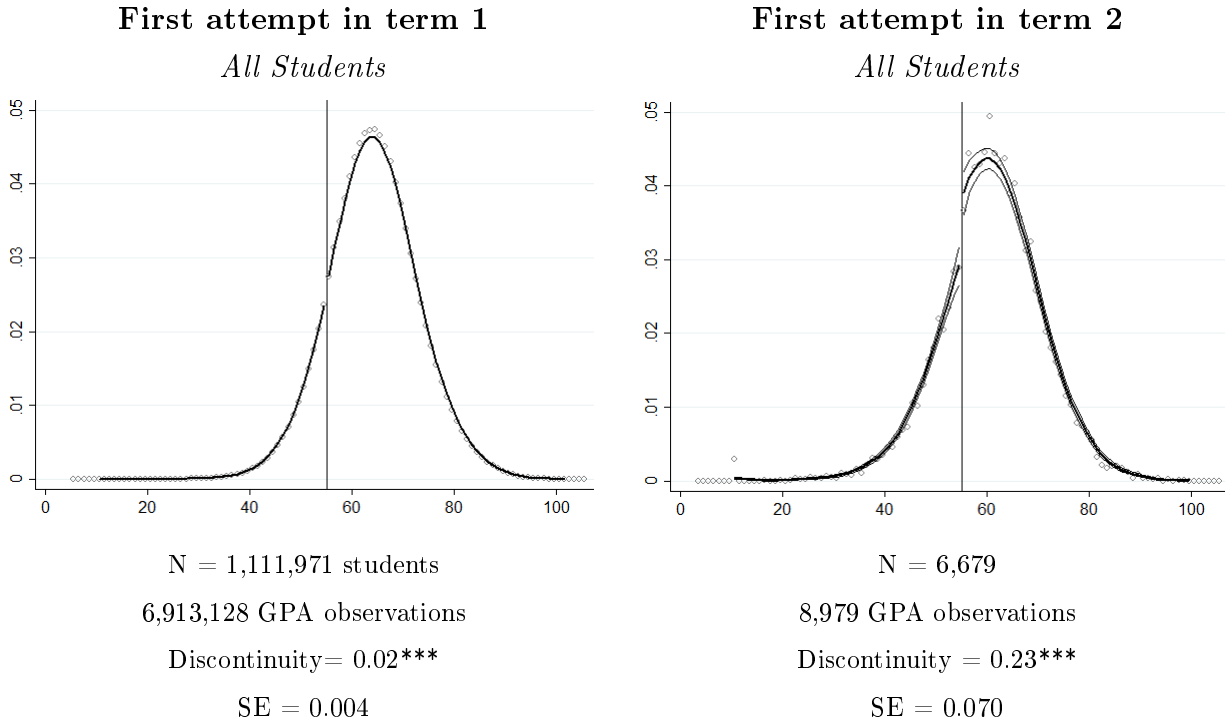
Table 3 compares characteristics of students who had their first exam attempt in term 1 with those who had their first attempt in term 2. The table indicates that sample of students with first attempt in term 2 (situation B) are not a random sample of the total student population. Students with a first attempt in term 2 are, on average, somewhat older, have achieved lower school exam grades, more frequently have a migrant background and live in an impoverished neighborhood. These differences indicate that a selective group of relatively lower performing students had their first attempt in term 2.

Table 3. Selective participation in first exam attempt

	First Attempt in			
	Term 1		Term 2	
	Mean	SD	Mean	SD
Male	0.50	0.50	0.50	0.50
Age on October 1 st	16.07	0.96	16.18	1.05
Migrant Background	0.20	0.40	0.28	0.45
Impoverished Neighborhood	0.13	0.34	0.18	0.39
<i>SE</i>	6.52	0.84	6.36	0.91
Pre-vocational education	0.55	0.50	0.50	0.50
Upper general education	0.26	0.44	0.29	0.46
Pre-university education	0.19	0.39	0.21	0.41
<i>N</i>	1,111,971		6,679	

Figure 6 then compares the subject-specific GPA distributions of first exam attempts for term 1 and 2 separately. The left panel shows no sizable discontinuity (i.e. 0.02) at the passing threshold of 55 (the statistical significance is primarily the result of the large number of observations in situation A). This is as expected, since for a first attempt in term 1 it is not possible to convert points to grades and concurrently to target grade manipulation effectively. Yet, the right panel (situation B) shows a substantial and significant discontinuity at the passing threshold for first attempts observed in term 2, indicating that a significant proportion of students is transferred from the left to right of the test score distribution. The McCrary density discontinuity result of 0.23 thus indicates that either teachers use their discretion to enable students to pass the subject (and graduate) in the first attempt ($\Delta_i > 0$) and/or that students have (temporarily) boosted their ability to pass the subject.

Figure 6. Final test score distribution first attempts: term 1 vs term 2



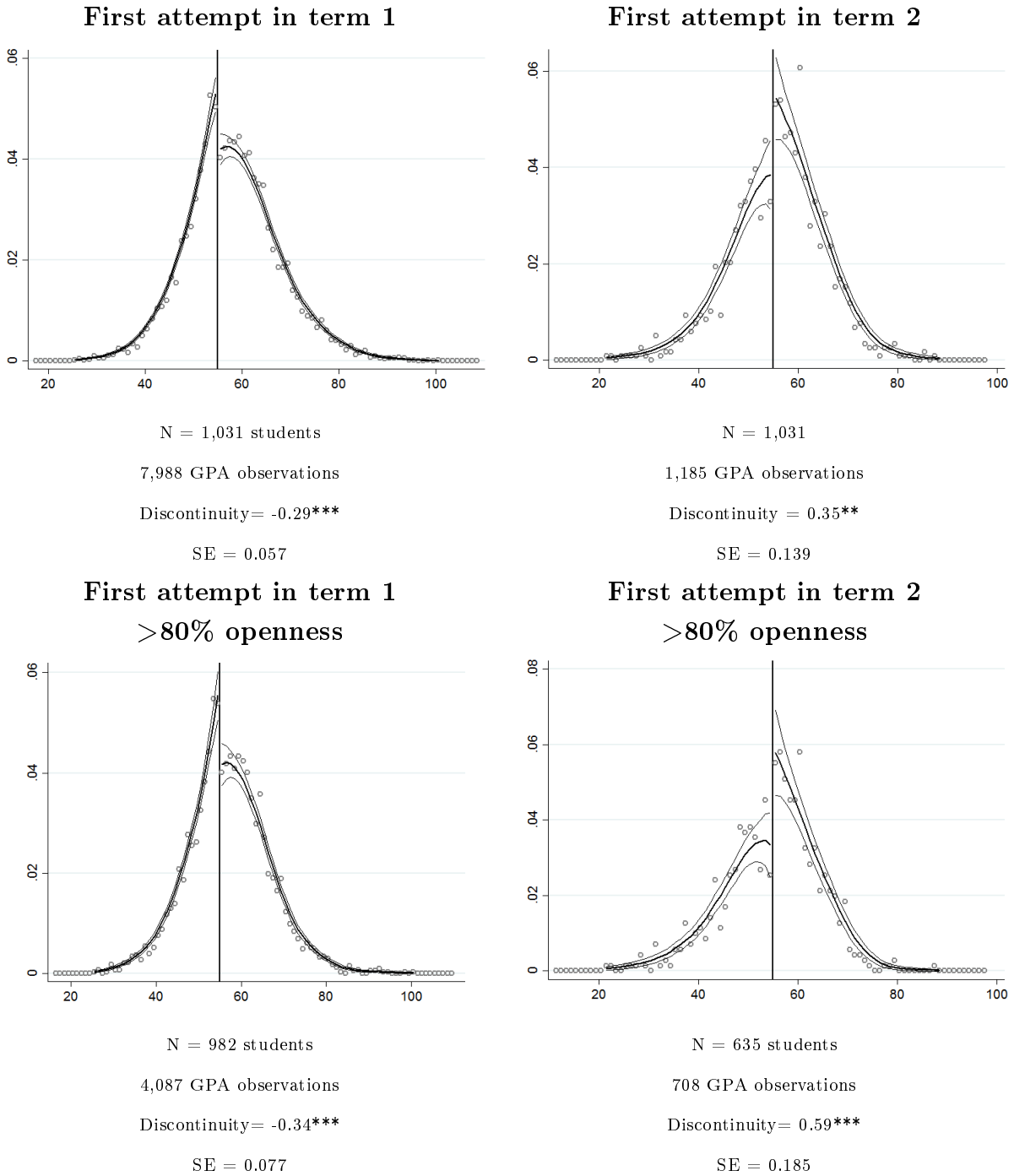
To distinguish whether it is the teacher and/or student creating this difference between situation A and B, the sample is restricted to contain only the high-stakes retaking students ($N=136,638$) introduced in the previous section. Whereas high-stakes retaking students represent roughly 12% (136,638 out of 1,118,650) of the total sample, they make up 15% (1,031 out of 6,679) of the students observed to have at least one first attempt in term 2. This corroborates the aforementioned selective nature of this group. The reason they are selected for the results displayed in Figure 7 is that if a subject exam's first attempt is observed in term 2 for this subsample of students, it by definition is a high stakes event (i.e. given that they need a retake for graduation, failing this subject for which the first attempt is observed in term 2 will be detrimental for their propensity to graduate).

The upper left McCrary density test results relate to the subject-specific first-attempt GPA distribution observed in term 1 and the observed negative discontinuity indicates that high-stakes retaker are more frequently (just) failing a particular subject. Yet, when for these students the GPA distribution is analyzed for the subject(s) for which a first attempt is observed in term 2, a markedly different picture emerges in that now a positive discontinuity

is observed.

This discontinuity is larger than for the overall population of students for which first attempts are observed in term 2 (i.e. 0.35 versus 0.23 overall), which is in line with the high-stakes nature of this subpopulation of students. Furthermore, the bottom panels in Figure 7 indicate that results for exams with more than 80% open questions are similar to other exams in terms of their first attempt term 1 distributions, but even larger positive discontinuities are observed when it concerns a first attempt exam observed in term 2 (i.e. 0.59 versus 0.35 overall). While these results only concern a small subsample of the overall student population, it reaffirms that teachers exploit their discretion to artificially improve performance as to let students graduate who would otherwise not have graduated. Furthermore, this phenomenon is observed in a context when potential mean reversion does not come into play, reveals itself when the N -factor is known, is larger in magnitude when students are at risk of not graduating, and larger when teachers have more discretion when grading the exam.

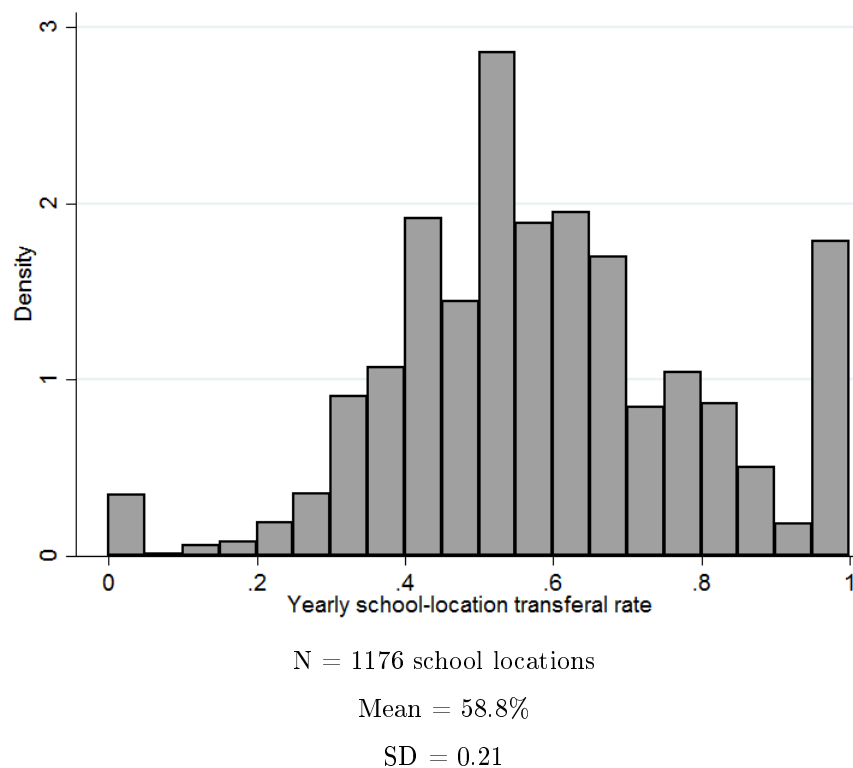
Figure 7. Final test score distribution of high-stakes retakers' first attempt: term 1 vs term 2



5.3 Teacher Discretion and Unequal Graduation Opportunities

Notwithstanding that teacher discretion effects can originate from a genuine desire to help students graduate, it can have undesirable effects in that it can cause inequitable between- and within-school variation in graduation opportunities for high-stakes retakers. The proportion of high-stakes retakers that graduated due to the retake exam result is referred to here as the transferal rate. The average transferal rate for the sample of high-risk retakers is 54.8%. Figure 8 displays the transferal rate *by school location and year* and the distribution shows a substantial amount of variance, indicating that in some schools no high-stakes retakers are transferred towards graduation in a given year, while in other schools all high-stakes retakers are transferred towards graduation in that year.

Figure 8. School-level Transferal Rate Distribution



Observing variation around the mean is not necessarily problematic with respect to between-school inequity, as long as schools are not structurally located high (or low) in this distribution over time. To examine whether this is the case, we estimate several random

effects models shown in Table 4. Baseline model 1 shows that the (weighted) average school-level transferal rate is 58.1% and the intra-class correlation coefficient (ρ) indicates that 65.2 percent of the observed variation in transferal rate is due to between school variation. When education level dummies are included in model 2, the residual intra-class correlation coefficient becomes lower, but is still 58.6 percent. It can be argued there may also be sorting effects of students into schools, in that high-quality schools attract better students. However, the random effect analysis is performed only for high-stakes retakers, thus who all required a retake exam for graduation after taking a nationwide standardized exam. Yet, to control for potential sorting effects, student controls are added in model 3, including the school examination grade and the central exam grade of the first period, and the results show that still 56.4 percent of the variation in school-level transferal rate is between school-level variation. The final model shows that between-school variation is not driven by structural differences between schools in the proportion of open-ended questions of retake exams observed across school locations. These results show that high-stakes retakers are transferred towards graduation structurally more often in some schools than in others. As such, this offers a substantial source of between-school inequity in graduation opportunities for these high-stakes retakers.

With respect to potential concerns of within-school inequity, Table 5 shows the regression results for student-level transferal status. Similar to Table 2, the inclusion of term 1 *CE*-grades ensures that the regression coefficient for *POQ* isolates teacher discretion effects from potential mean reversion and ability boosting effects. The estimation results of Model 1 indicate that the proportion of high-stakes retakers transferred to graduation increases significantly as a result of teacher discretion. Student-level controls and interactions between *POQ* and gender and ethnicity are included in Model 2 to examine if the teacher discretion effects vary with these characteristics *within schools*. The estimation results reveal that teacher discretion effects are statistically significantly smaller for non-Dutch students, but this difference remains very small when education-level and year dummies (Model 3), and school-location dummies (Model 4) are included. It can thus be concluded that teacher discretion effects cause substantial inequitable between-school variation, and that there is marginal within-school variation due to the heterogeneity of teacher discretion effects with respect to a student's immigrant background. These results indicate that teacher discretion effects result in unequal graduation opportunities (by school choice) and that these effects arise because (teachers in) some schools exploit discretion in grading retake exams more than (teachers in) other schools.

Table 4. Between-school inequity: random-effects model

<i>Dependent variable: school-level year transferal rate</i>								
	1		2		3		4	
	Coeff	SE	Coeff	SE	Coeff	SE	Coeff	SE
Constant	0.581***	(0.004)	0.573***	(0.004)	0.534***	(0.009)	0.526***	(0.009)
Education level dummies	<i>No</i>		<i>Yes</i>		<i>Yes</i>		<i>Yes</i>	
Student controls	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
Exam openness	<i>No</i>		<i>No</i>		<i>No</i>		<i>Yes</i>	
# students	136,698		136,698		136,698		136,698	
# school locations	1,176		1,176		1,176		1,176	
sigma_u	0.143		0.125		0.119		0.119	
sigma_e	0.105		0.105		0.104		0.104	
R ²	0.000		0.098		0.070		0.069	
rho	0.652		0.586		0.564		0.564	

Note: */**/** denote significance at 10/5/1% level (two-sided). Outcome variable in Models 1-4 is the observed yearly transferal rate for the school location a student attends. Standard errors are clustered at the subject-year level. Student-level controls are: boy, age, non-Dutch background, impoverished neighborhood, SE-grade and a dummies for whether a covariate is missing. Education level dummies are: pre-vocational education, upper general education. Exam openness is the percentage of open-ended questions of the retake exam. School-location random effects are based on 1176 school locations.

Table 5. Within-school inequity: heterogeneous graduation transferal

<i>Dependent variable: student-level transferal status</i>								
	1		2		3		4	
	Coeff	SE	Coeff	SE	Coeff	SE	Coeff	SE
Proportion Open Questions (POQ)	0.064***	(0.036)	0.104***	(0.033)	0.091***	(0.025)	0.086***	(0.024)
CE grade - attempt 1	0.005***	(0.0005)	0.004***	(0.0004)	0.003***	(0.0003)	0.003***	(0.0003)
CE grade - attempt 1 · POQ	-0.002***	(0.0005)	-0.002***	(0.0004)	0.0003	(0.0003)	0.0003	(0.0003)
boy			0.010*	(0.06)	-0.005	(0.005)	0.002	(0.004)
non-Dutch background			-0.030***	(0.003)	-0.053***	(0.003)	-0.050***	(0.004)
POQ · boy			0.001	(0.007)	0.001	(0.006)	0.001	(0.006)
POQ · non-Dutch background			-0.022***	(0.004)	-0.009***	(0.003)	-0.006**	(0.003)
Constant	0.348***	(0.028)	-0.349***	(0.032)	-0.546***	(0.043)	0.189***	(0.089)
Student-level controls	<i>No</i>		<i>Yes</i>		<i>Yes</i>		<i>Yes</i>	
Level & Year dummies	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
School-location dummies	<i>No</i>		<i>No</i>		<i>No</i>		<i>Yes</i>	
N	136,698		136,698		136,698		136,698	
R^2	0.006		0.038		0.073		0.098	
# subject-year clusters	389		389		389		389	

Note: */**/** denote significance at 10/5/1% level (two-sided). Outcome variable in Models 1-4 is the observed transferal status for a retaking student. Standard errors are clustered at the subject-year level. Student-level controls are: boy, age, non-Dutch background, impoverished neighborhood, SE-grade and a dummies for whether a covariate is missing. Education level dummies are: pre-vocational education, upper general education. Year dummies are: 2008, 2009, 2010, 2011, 2012. School-location dummies are based on 1176 school locations.

6 Discussion

This study shows that teachers structurally use their discretion to increase the performance of their students with the objective to let them graduate. This discretion is targeted at student who find themselves just below the passing threshold and effectuated when the stakes are highest and there is full information on how to convert assigned points to grades. To distinguish teacher discretion effects from the effects of ability exploitation and mean reversion we use the proportion of open-ended questions (summary, essay) as an instrument to identify teacher discretion effects. The identifying assumption is that the proportion open-ended questions is unrelated with ability boosting of students (i.e. students -temporarily- raise their performance level with the objective to graduate) and mean reversion, but can be positively related to teacher discretion. We find that teacher discretion is revealed when the teacher has full information regarding the conversion of points obtained to grades, that the effect is larger in magnitude when students are at risk of not graduating, and larger when teachers have more discretion when grading the exam. The results suggest that approximately 12% of all students who graduated by means of a retake exam did so because of teacher discretion. This roughly translates to 1,400 students in any given exam year. This result is derived only from teacher discretion at retake exams and -given that teachers locally grade many more tests (e.g. school exams)- could therefore be considered to be a lower-bound estimate of the overall implications of teacher discretion effects in Dutch secondary education.

Notwithstanding the good intentions teachers arguably have to artificially improve the grades of students who are on the margin of graduation, it may have undesirable inequity effects. First of all, it results in between-subject variation (i.e. the retake subject choice matters for graduation). Secondly, it can cause both inequitable between-school variation (i.e. school-location differences in transferal rates) and within-school variation due to heterogeneity of teacher discretion effects with respect to student-level characteristics. The structural differences observed between schools indicate that teacher discretion effects result in unequal graduation opportunities (by school choice) and these effects arise because teachers in some schools exploit their discretion more than teachers in other schools.

When objective skills-assessment is a priority of the school-leaving exams, the results presented here show that teacher discretion issues can -at least in theory- be easily resolved by either avoiding teacher discretion when grading, or by imposing that the nation-wide central exams are graded anonymously. Obviously, this can potentially start a different public debate about whether it is desirable that students who are just below the passing threshold have to redo the examination year (either partially or entirely). However, these

valid questions are directly targeted at the functioning of the exam system, and stand alone in the fundamental argument that students should have equal educational opportunities.

References

- Burgess, Simon and Ellen Greaves (2013), ‘Test scores, subjective assessment, and stereotyping of ethnic minorities’, *Journal of Labor Economics* **31**(3), 535–576.
- Cornelisz, Ilja and Chris Van Klaveren (2018), ‘Student engagement with computerized adaptive practicing: Ability, task value and difficulty perceptions’, *Journal of Computer Assisted Learning* pp. 1–15.
- Dee, Thomas S, Will Dobbie, Brian A Jacob and Jonah Rockoff (2016), The causes and consequences of test score manipulation: Evidence from the new york regents examinations, Technical report, National Bureau of Economic Research.
- Diamond, Rebecca and Petra Persson (2016), The long-term consequences of teacher discretion in grading of high-stakes tests, Technical report, National Bureau of Economic Research.
- Hanna, Rema N and Leigh L Linden (2012), ‘Discrimination in grading’, *American Economic Journal: Economic Policy* **4**(4), 146–168.
- Jacob, Brian A (2005), ‘Accountability, incentives and behavior: The impact of high-stakes testing in the chicago public schools’, *Journal of public Economics* **89**(5), 761–796.
- Kuhlemeier, Hans and Ed Kremers (2012), De praktijk van de eerste en tweede correctie van het cse. verslag van een landelijke enquête, Technical report, Arnhem: Cito.
- Kuhlemeier, Hans and Ed Kremers (2013), De praktijk van eerte en tweede correctie. samenvatting van onderzoek naar het functioneren van het cse, Technical report, Arnhem: Cito.
- Lavy, Victor (2008), ‘Do gender stereotypes reduce girls’ or boys’ human capital outcomes? evidence from a natural experiment’, *Journal of public Economics* **92**(10), 2083–2105.
- McCrary, Justin (2008), ‘Manipulation of the running variable in the regression discontinuity design: A density test’, *Journal of Econometrics* **142**(2), 698 – 714.
URL: <http://www.sciencedirect.com/science/article/pii/S0304407607001133>

McMillan, James H and Suzanne Nash (2000), 'Teacher classroom assessment and grading practices decision making.'

Neal, Derek (2013), 'The consequences of using one assessment system to pursue two objectives.', *Journal of Economic Education* **44**(4), 339–352.

Schuurs, Uriël, Hans Kuhlemeier and Hugo Gitsels (2017), 'De invloed van het lt-examenverslag op de scores', *Levende Talen Tijdschrift* **18**(4), 25–35.

Appendix A Graduation Rules

For the evaluation window considered in this study, the rules state (REF) that students in pre-vocational education pass the school-leaving examinations if one of the following situations hold:

1. GPA for all subjects is at least 5.5
2. GPA for one subject is between 4.5 and 5.45, for all other subjects at least 5.5
3. GPA for one subject is between 3.5 and 4.45, for one subject at least 6.5, and for all other subjects at least 5.5
4. GPA for two subjects is between 4.5 and 5.45, for one subject at least 6.5, and for all other subjects at least 5.5

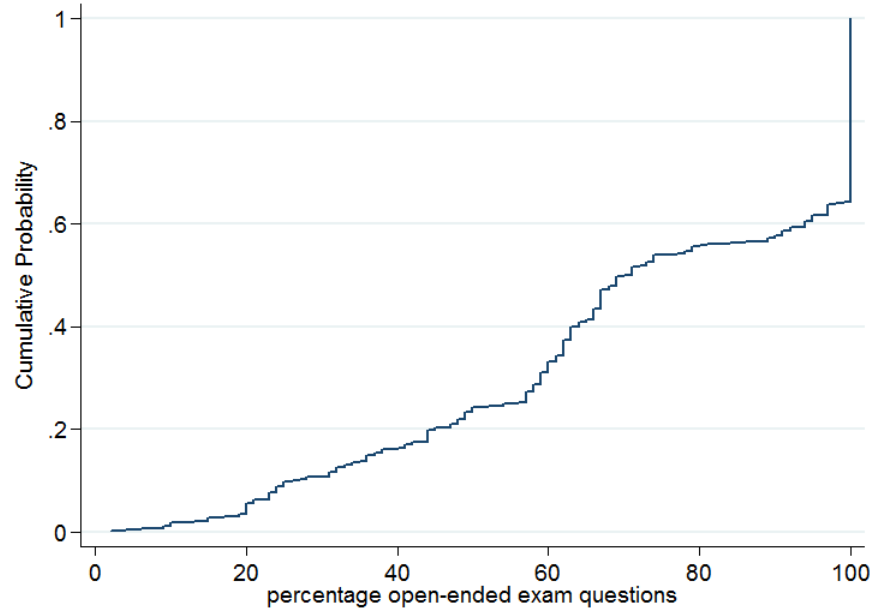
For students who are in secondary general education, and pre-university education these rules are:

1. GPA for all subjects is at least 5.5
2. GPA for one subject is between 4.5 and 5.45, for all other subjects at least 5.5
3. GPA for one subject is between 3.5 and 4.45 and for all other subjects at least 5.5 and overall GPA is at least 6.0 based on subject grades rounded to whole integers.
4. GPA for two subjects is between 4.5 and 5.45, and for all other subjects at least 5.5 and overall GPA is at least 6.0 based on subject grades rounded to whole integers.
5. GPA for one subject is between 3.5 and 4.45, for one subject between 4.5 and 5.45, and for all other subjects at least 5.5 and overall GPA is at least 6.0 based on subject grades rounded to whole integers.

For 2012, an additional graduation requirement for students in secondary upper general and pre-university education is that the average CE-grade across all subjects is at least 5.5.

Appendix B Proportion Open Questions and CE_1 performance

Figure B.1: Cumulative Distribution Function of Proportion Open Questions



N = 386 subject-year clusters

Mean = 71.7%

SD = 0.28

Table B.1: CE_1 scores and Proportion Open Questions for graduating retakers

	1		2		3		4	
	Coeff	SE	Coeff	SE	Coeff	SE	Coeff	SE
Prop. Open Questions (POQ)	0.312	(0.736)	0.609	(0.876)	0.632	(0.994)	0.743	(0.945)
Constant	44.72***	(0.509)	42.97***	(5.582)	40.90***	(3.660)	41.88***	(3.092)
Student-level controls	<i>No</i>		<i>Yes</i>		<i>Yes</i>		<i>Yes</i>	
Education Level dummies	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
Year dummies	<i>No</i>		<i>No</i>		<i>Yes</i>		<i>Yes</i>	
School Location Fixed Effects	<i>No</i>		<i>No</i>		<i>No</i>		<i>Yes</i>	
N	69280		69280		69280		69280	
R^2	0.029		0.051		0.061		0.092	
# subject-year clusters	386		386		386		386	

Note: Robust standard errors in parentheses. */**/** denote significance at a 10/5/1 percent confidence level. Outcome variable in Models 1-4 is the observed first-attempt CE-grade for a graduating retaker. Standard errors are clustered at the subject-year level. Student-level controls are: boy, age, non-Dutch background, impoverished neighborhood, SE-grade and a dummies for whether a covariate is missing. Education level dummies are: pre-vocational education, upper general education.