

A quantitative analysis of the uncertainty in grading of written exams in mathematics and physics

Type:

Research paper

Abstract:**Background**

The most common way to grade students in courses at university and university college level is to use final written exams. The aim of final exams is generally to provide a reliable and a valid measurement of the extent to which a student has achieved the learning outcomes for the course. A source of uncertainty in grading students based on an exam is that such exams only consist of a limited number of exercises.

Material and methods

We investigate the extent of this uncertainty by means of a statistical analysis of the results of 23 different examinations taken by 2788 students.

Results

The amount of uncertainty is substantial and typically ranges over three grades.

Conclusions

Increasing the duration of the examination decreases the uncertainty, however.

Keywords:

Grading, uncertainty, quantitative research, examination duration, written exam

1 **Running head: Uncertainty in grading written exams**

2 **Title: A quantitative analysis of uncertainty in the grading of written exams**
3 **in mathematics and physics**

4 **State of the literature**

- 5
- 6 • There is a long tradition of using statistical approaches to analyze the reliability and
7 validity of written exams, see, e.g., Lord (1952), Lord (1953), Lord & Novick (1968).
 - 8 • More recent advances: Item Response models (Hambleton, Swaminathan, and
9 Rogers, 1991; Lord and Novick, 1968; Lord 1980) and the Generalized Partial Credit
10 model (Muraki and Bock, 2002; Muraki 1997)
 - 11 • References focusing on the number of exercises that should be included in a test or
12 exam to alleviate the challenges of uncertainty in grading: Bird and Yucel (2013) and
Burton (2006).

13 **Contribution of this paper to the literature**

- 14
- 15 • We have not found any paper analyzing real exam correction data to measure
16 uncertainty in the grading of written exams in mathematics and physics.
 - 17 • The most likely reason is that none of the traditional models, such as the Item
18 Response model or the Generalized Partial Credit model, are applicable to such data.
19 Instead, we construct a suitable model based on the less common Beta Regression
20 framework.
 - 21 • We demonstrate that exam correction data are a very valuable source of information
22 for measuring uncertainty in the grading of written exams in mathematics and physics.
23 Our results show that the uncertainty is substantial and typically ranges over three
24 grades.

A quantitative analysis of uncertainty in the grading of written exams in mathematics and physics

Abstract

The most common way to grade students in courses at university and university college level is to use final written exams. The aim of final exams is generally to provide a reliable and a valid measurement of the extent to which a student has achieved the learning outcomes for the course. A source of uncertainty in grading students based on an exam is that such exams only consist of a limited number of exercises. We investigate the extent of this uncertainty by means of a statistical analysis of the results of 23 different examinations taken by 2788 students. The amount of uncertainty is substantial and typically ranges over three grades. Increasing the duration of the examination decreases the uncertainty, however.

Keywords: examination duration, grading, quantitative research, uncertainty, written exam

Introduction

Background to the study

Quality in higher education has been the focus of much policy development worldwide (Blanco-Ramírez and Berger, 2014), often in response to external incentives or to comply with norms that are considered legitimate (Vukasovic, 2013). This situation has resulted in increased focus on devising and implementing quality assurance systems in institutions of higher education (Westerheijden et al., 2014). The notion of academic quality is multifaceted and has been described as “an inherently vague concept” (Wittek and Kvernbekk, 2011). It may also be noted that, when measuring the quality of education, the emphasis has shifted from educational inputs (what the teacher conveys and how) to learning outputs (what the students have achieved in terms of learning outcomes), as reported in Hughes (2013). In that respect, it is natural to consider the quality of assessment methods as an inherent part of the quality of an educational program (Boyas, Bryan, and Lee, 2012).

Examination results have been referred to as a form of “currency” (William, 1996, Simpson and Baird, 2013) that is dependent on trust in order to retain its status and value. In tertiary education, the primary purpose of final examination grades is to communicate a student’s achievement to future employers and to other institutions to which the student might apply for further degrees. It is generally accepted that grades can be of critical importance to a student’s future educational path and career, and that it is therefore crucial to ensure that they accurately reflect the student’s proficiency level.

Issues relating to the reliability and validity of assessment have been the focus of attention in the literature on assessment. Reliability has been defined as “the repeatability of

62 an assessment and its results” (Irwin and Hepplestone, 2012, 777). An assessment form can be
63 said to be reliable if it is not affected by factors that lay outside the student’s control, such as
64 the background or the views of the examiner (Harlen, 2005). Changes in difficulty levels from
65 one year to another, or from one semester to another, may also endanger the reliability of an
66 assessment form (DeVellis, 2012).

67 Validity refers to the extent to which an assessment form “measures what it is
68 designed to measure” (Russell et al., 2006, 466). In that respect, an assessment form is only
69 valid if it allows the students to demonstrate effectively whether and to what degree they have
70 achieved the learning goals that were set for the course. A valid assessment is therefore one
71 that prevents students from over-performing or under-performing compared with their actual
72 level of mastery of the curriculum. Altogether, validity necessitates reliability, but reliability
73 is not in itself sufficient to ensure validity.

74 Institutions of higher education throughout the world are currently under pressure to
75 increase productivity due to budget cuts (Agasisti and Bonomi, 2014) and growing student
76 numbers (Allais, 2014). In an educational climate where cost-efficiency is emphasized,
77 institutions may feel pressure to reduce the duration of examinations in order to reduce the
78 costs associated with remunerating invigilators, renting examination rooms, and compensating
79 faculty members, teaching assistants or external examiners for marking the examination
80 papers. In addition, institutions may face examination timetabling problems due to the limited
81 number of weeks that can be allocated to examinations and because of growing student
82 numbers (as suggested in, e.g., Mumford (2010) and Abdul-Rahman et al. (2014)).

83 Compared with other types of summative assessment, such as oral examinations,
84 closed-book written examinations at the end of the module, semester, or academic year are
85 relatively inexpensive. Computer-based approaches (Delen, 2015) (as described in, e.g.,

87 Delen (2015) and Kuo et al. (2015)) are another inexpensive approach, although they are
88 much less used than written exams. Written examinations are regarded as particularly suitable
89 for testing students' learning outcomes in mathematics and other science subjects (Davis et
90 al., 2005). Relatively little attention has been devoted, however, to the reliability and validity
91 of written examinations. Interestingly, the literature on assessment seems to be more
92 concerned with the validity of other assessment forms, such as practical examinations (Vu et
93 al., 2006), modified essay questions (Palmer et al., 2010), or portfolio assessment (Admiraal
94 et al., 2011).

95 There are a few notable exceptions, however, for example the work of Bird and Yucel
96 (2013) and Burton (2006). The latter suggests that the optimal length of an academic test
97 consisting of short-answer and multiple-choice questions with dichotomous scoring (either 0
98 or 1) might be around 300 questions or test items. This is in order to allow for different levels
99 of difficulty in the questions, unevenness of knowledge among the students taking the test,
100 and the possibility that some questions may be badly phrased, while other questions may be
101 so similar to the textbook material that students can answer them correctly more from
102 memory than by reasoning. He also points out that testing more than two separable facts in
103 one dichotomously scored test item provides an additional level of uncertainty and
104 recommends avoiding the use of such "double questions" (p. 576).

105 There is a long tradition of using statistical approaches to analyze the reliability and
106 validity of written exams. Foundational works such as Lord (1952) and Lord (1953) have
107 highlighted the need to differentiate between ability scores, which are test-independent, and
108 observed scores and true scores, which are test-dependent. Other works, such as Lord &
109 Novick (1968), have described classical test theory as relying on the assumption that test
110 scores are the result of a combination of true scores and measurement error.

112 The study described in this article aims to provide insights into the reliability and
113 validity of written exams in mathematics and physics. To that end, we present an extensive
114 analysis of uncertainty in the grading of written exams. Such exams typically consist of 10 to
115 20 exercises from different parts of the curriculum, and the reliability is affected, among other
116 things, by the number of exercises included. We analyze the reliability of such exams using a
117 quantitative approach based on an extensive dataset consisting of the marking of 34 800
118 examination answers from 2788 students based on exams from two universities and one
119 university college in Norway. We analyze the data using a Generalized Linear model (Dobson
120 and Barnett 2008a). Generalized Linear models have been applied extensively in educational
121 measurement or educational assessment through models such as Item Response models
122 (Hambleton, Swaminathan, and Rogers, 1991; Lord and Novick, 1968; Lord, 1980) and the
123 Generalized Partial Credit model (Muraki and Bock, 2002; Muraki, 1997). It is worth noting,
124 however, that all these models assume that the test scores are discrete (e.g., right/wrong). This
125 suggests that traditional assessment models cannot be used to shed light on data material
126 where the scores are continuous, as is the case in our study. The analysis in this article is thus
127 based on a less common Generalized Linear model called Beta Regression.

128 To the best of our knowledge, there is little published research that takes a quantitative
129 approach to analyzing the reliability and validity of written exams. We assume that the reason
130 for this is that it is not possible to analyze continuous data using the traditional statistical
131 assessment models described above. Our decision to use beta regression may therefore
132 represent a significant contribution to the field of assessment, since it provides new insights
133 into the reliability and validity of written exams.

Examples

In order to ensure both the reliability and validity of exams in mathematics or physics, such exams must include a sufficient amount of exercises to test the students' actual level of mastery. The following example could be used to illustrate this claim. Let us consider two students with very different levels of mastery of the course curriculum, which consists of 10 main parts. If student A masters only one of the ten parts and student B masters nine of the ten, an examination with only one exercise aimed at testing just one part of the curriculum might give a totally erroneous picture of the students' actual level of mastery. If the exercise happens to be on the one part of the curriculum that student A masters, he or she will get a good grade, which does not reflect his or her actual level of mastery of the curriculum. Conversely, if the exercise happens to be on the one part of the curriculum that student B does not master, he or she will be awarded a poor grade that does not reflect his or her level of mastery either. In order to reduce the random effect of luck (or lack thereof) on examination scores and thereby increase their validity and reliability, it is necessary to ensure that each examination consists of a sufficient amount of exercises.

A second example may further illustrate the problem. We assume that an exam consists of an equal amount of very difficult, difficult, easy, and very easy exercises from different parts of the curriculum. For an average student, we assume that the probabilities of the student managing exercises on the different levels of difficulty are 0.2, 0.4, 0.6, and 0.8, respectively. We assume that the time allotted per exercise is 15 minutes, which is typical for traditional exams in mathematics and physics. The exam score will be the mean of the scores for each exercise. Figure 1 shows the probability of different exam scores for the student in this simple binomial model. The exam scores are rescaled to a 0 to 100 scale. We see that, for a one-hour exam (four exercises), there is a probability of approximately 4% that the student

160 will get all the answers wrong and an equal probability that the student will get all the answers
161 correct, which means that the reliability of such a short exam is very poor. In the case of four-
162 hour and eight-hour exams, the possible exam results are spread over several grades as well,
163 which means quite poor reliability. In the rest of this paper, a similar analysis will be
164 performed where the uncertainty (lack of reliability) is estimated based on the data from the
165 marking of the 34 800 examination answers.

166 **Figure 1:** Probability of different exam scores in a binomial exam model.[About here]

167 **Material and methods**

168 **Exam correction data**

169 Our analysis was based on the marking of 23 exams for introductory courses in
170 mathematics, statistics, and physics from the University of Oslo, the Norwegian University of
171 Science and Technology, and Oslo and Akershus University College of Applied Sciences.
172 The material consists of marks awarded to 2 788 different students for 301 different exercises.
173 The marking of each exercise for all the students in all the exams will be used in the analysis,
174 ending up with a total of 34 800 observations. For each exam, the marks (scores) for the
175 exercise answers were normalized to the $[0, 1]$ interval, where completely wrong and
176 completely correct answers were awarded zero and one point, respectively.

177 The characteristics of the dataset were as follows. Each exercise in the data material
178 required the student to perform some kind of calculations, i.e., no multiple choice exercises

180 where the student could guess the answer. All exams were traditional written exams using pen
181 and paper.

182 The duration of the exams varied between three and five hours. The time allotted to
183 solving each exercise varied between the different exams, ranging from 12 minutes to 18
184 minutes. For exercises where the students were given a long time per exercise, the exercises
185 typically consisted of many subtasks or longer computations.

186 *Methodological issues*

187 The markings (scores) of exercises from earlier written exams are an extremely useful
188 source of information for measuring the reliability and validity of written exams, as will be
189 seen in the results section. Quite surprisingly, we have not found any research papers that take
190 advantage of this valuable source of information. The most likely explanation is that the data
191 are available in a format that does not easily lend itself to analysis. In this article, we show
192 that Beta Regression, a type of Generalized Linear model (Dobson and Barnett, 2008b), is a
193 suitable choice. A motivation for and detailed description of the statistical model is provided
194 below.

195 **Statistical model**

196 **Figure 2:** Beta distribution for a variety of values of the shape parameters (a, b). The black
197 curves show typical distributions of scores on exam exercises. [About here]

198 In this section, we describe a statistical model that quantifies the amount of uncertainty
199 in the grading of written exams in mathematics and physics. The statistical model has much in
200 common with Item Response models and Generalized Partial Credit Interval, but our model
201 differs from these models in that the responses (test scores) are not discrete, but continuous on

203 a limited interval. Naturally, if the uncertainty in grading is high, the reliability and validity of
204 the exam will be low. Let M quantify the level of mastery for a student taking an exam. A
205 student with a high level of proficiency in the subject will have a large value of M , while a
206 student with a low level of proficiency in the subject will have a small value. Further, let D
207 quantify the level of difficulty of an arbitrary exercise in an exam. An easy exercise will have
208 a low value of D , while a difficult exercise will have a large value of D . Let
209 S_1, S_2, \dots, S_n denote the scores for a student for the different exercises in an exam. We assume
210 that each score is given on the $[0, 1]$ interval, where a completely wrong answer results in the
211 score zero, a completely correct answer in the score one and a partly correct answer
212 somewhere in between. Let S_E denote the resulting exam score (grade) for this student based
213 on the exercise scores S_1, S_2, \dots, S_n . The most common way to compute the exam score, S_E , is
214 to take the average of each exercise score S_i and multiply by 100

$$215 \quad S_E = \frac{100}{n} \sum_{i=0}^n S_i \quad (1)$$

216 For an exam to have high reliability and validity, the uncertainty of the exam score, S_E , must
217 be low. The main source of uncertainty in the exam score is that, if a student is given many
218 exercises of the same level of difficulty, D , the student will by chance alone get some
219 exercises correct, some wrong, and some partly correct. If the student is lucky, an exam will
220 consist of many exercises that the student, by chance, is able to solve. If the student is
221 unlucky, the exam will consist of many exercises that the student, by chance, is not able to
222 solve. Recall the two examples at the end of the introduction. The best way to reduce this
223 source of randomness in exam score, S_E , is to include many exercises that are well-suited to
224 testing different parts of the curriculum. Since the exam score is typically the average of the

226 exam scores (Equation (1)), by the law of large numbers, the exam score will approach the
227 student's true level of mastery, m , when the number of exercises increases.

228 As described above, a student who is given many exercises of the same level of
229 difficulty, D , will simply by chance get some exercises correct, some wrong, and some partly
230 correct. Let $p(s; m, d)$ denote a probability distribution that summarizes this property. More
231 specifically, $p(s; m, d)$ is the probability distribution of exercise scores, S , a student with a
232 level of mastery M will be awarded for an exercise of difficulty level D . If this probability
233 distribution is narrow (small variance), the uncertainty in exercise scores S_1, S_2, \dots, S_n is small
234 and the resulting uncertainty in exam score, S_E , will be small (recall equation (1) and the law
235 of large numbers). If the distribution $p(s; m, d)$ is wide, the uncertainty in exam score, S_E ,
236 will be large. We also expect that, if a student has a high level of mastery M or the exercise is
237 easy (low value of D), the distribution will shift toward high values of exercise scores, and
238 shift toward low values if the student has a low level of mastery or the exercise is difficult.

239 We estimate the probability distribution $p(s; m, d)$ using a regression model where the
240 exercise score S is the dependent variable and M and D are the independent variables,
241 modelled as random effects. The distributions for level of mastery (M) of each student, level
242 of difficulty (D) of exercises, and the relation between M , D , and S are estimated using the
243 marking (scores) of the 34 800 exercises in the data material. Traditional statistical
244 assessment models assume that the response is binomial. Since the exercise
245 scores S_1, S_2, \dots, S_n represent continuous variables on the $[0, 1]$ interval in our data, the
246 binomial regression models cannot be used. We instead used the less common alternative of
247 assuming that the exercise scores are outcomes from a beta distribution. The beta distribution
248 is a highly flexible continuous distribution on the $[0, 1]$ interval depending on the choices of
249 the model parameters, as shown in Figure 2. It is thus an ideal distribution for modelling the

251 exercise scores S_1, S_2, \dots, S_n . The exercise scores are typically distributed with “U”-shapes
252 like the black curves in Figure 2, since it is most common to score answers to exercises as
253 either completely wrong (zero points) or completely correct (one point) (see Figure 3). A
254 more detailed description of the beta regression model is provided in the Appendix.

255 Results

256 As described above, the data material consists of the marking (score) of each exercise
257 for all the students in the 23 exams, resulting in a total of 34 800 scores. Figure 3 shows a
258 histogram of all these scores.

259 **Figure 3:** Distribution of scores for all the answers in the 23 exams. [About here]

260 We see that the most common scores are, as expected, zero and one, but also the scores $1/6$,
261 $1/4$, $1/3$, $1/2$, $2/3$, $3/4$, and $5/6$ are used to some extent.

262 We now fit 23 beta regression models, one for each exam.

263 We start by showing results from one out of the 23 exams, as representative of the
264 results of all the 23 exams. The estimated values for the parameters in the regression model
265 are shown in the Appendix (Table 2).

267 **Figure 4:** Distribution of exercise scores, $p(s; m, d)$, for an average student for exercises of
268 varying levels of difficulty. [About here]

269 Figure 4 shows the distribution of exercise scores, $p(s; m, d)$, for different levels of
270 difficulty for a student who on average scores 50 out of 100 points in the exam (average level
271 of mastery M). Easy and very easy exercises are represented by exercises being one and two
272 standard deviations easier than average exercises. Similarly, difficult and very difficult
273 exercises are represented by exercises being one and two standard deviations more difficult
274 than average exercises. We see that the distributions acquire the characteristic “U” shape,
275 which is as expected since the most common exercise scores in the data material are zero and
276 one (Figure 3).

277 **Figure 5:** Distribution of exam scores for an average student. [About here]

278 Now, suppose that the average student faces an exam with an equal amount of very
279 difficult, difficult, easy, and very easy exercises. For the exam that we will now study, the
280 time per exercise was 15 minutes, which means that a four-hour exam consists of 16
281 exercises, four exercises on each level of difficulty. The exam score is computed from the
282 exercise scores using equation (1). The distribution of possible exam scores for the student is
283 shown in Figure 5. As expected, the uncertainty in the exam score is reduced as the number of
284 exercises in the exam is increased (recall Equation (1) and the law of large numbers). Thus,
285 the reliability and validity of the exam increases when the number of exercises increases.

287 From Figure 5, we see that almost all possible exam scores that the average student can get
288 for a four-hour exam fall between 30 and 70 points.

289 **Figure 6:** Distribution of exam scores for a strong student. [About here]

290 Figure 6 shows the same as Figure 5, but for a student who is two standard deviations
291 stronger than an average student. By comparing Figures 5 and 6, we make the interesting
292 observation that the uncertainty in the exam score is smaller for strong students than for
293 average students.

294 **Figure 7:** Uncertainty in exam scores as a function of duration of exam. [About here]

295 We now present results for all the 23 exams. As a measure of uncertainty in exam
296 scores, we use the difference between the 95% and 5% quantile in the distribution of exam
297 scores. For example, for the four-hour exam in Figure 5, the 95% and 5% quantiles are 63.2
298 and 36.9, respectively, resulting in a difference of 26.3 points. As described in the methods
299 section, the time allotted to solving an exercise varies between exams (12 to 18 minutes). A
300 comparison based on the number of exercises is therefore not correct. An exam with a few
301 exercises but with little variability in the exercise scores can be better than an exam with
302 many exercises and high variability. Since we know the time allotted per exercise for the
303 different exams, we can re-compute from the number of exercises given to the duration of

305 examination, and compare this to the uncertainty in exam scores. The results are presented in
306 Figure 7 as described below.

307 Figure 7 shows the relationship between uncertainty in the exam score (the difference
308 between the 95% and 5% quantile as described above) and the duration of examinations for
309 all the 23 exams. We have two curves (dashed and solid) for each exam, representing cases
310 with little and much uncertainty in the exam score. The main contribution to the varying
311 uncertainty (difference between the 95% and 5% quantile) is the level of mastery of the
312 students. There is less variability in exam scores for weak and strong students compared to
313 average students (recall Figures 5 and 6). The uncertainty from the estimation of the true
314 regression parameters is also included. As expected, in Figure 7, we see that the uncertainty is
315 reduced when the duration of examinations increases (recall Equation (1) and the law of large
316 numbers). For example, in a two-hour exam, the uncertainty for an average student potentially
317 reaches above 50 points (out of 100), while in a four-hour exam, the uncertainty is rarely
318 above 35 points and, for a six-hour exam, rarely above 25 points. For strong and weak
319 students, the uncertainty is rarely above 40, 25, and 20 points, for two-hour, four-hour, and
320 six-hour exams, respectively. We see some differences between the 23 exams, but overall the
321 different exams have more or less the same amount of uncertainty in exam scores.

322 Discussion and conclusion

323 The analysis shows that there is substantial uncertainty in grading written exams due
324 to the limited duration of examinations. This means that the reliability and validity of written
325 exams in mathematics and physics are critically low. By increasing the duration of
326 examinations, the uncertainty will decrease and the reliability and validity improve. From
327 Figure 7, however, we see that the reduction in uncertainty is less when we go from a two-
328 hour to a four-hour exam compared to going from a four-hour to a six-hour exam.

330 The conversion of an exam score on the interval $[0, 100]$ to specific grades varies a lot
331 around the world, but the ECTS system (A-F) with conversions as shown in Table 2 is very
332 common (Radboud University Nijmegen, 2011). For all international grading systems, the
333 interval for each grade is typically between 5 and 20 points wide (except for the interval for
334 ‘fail’). This means that the uncertainties documented in Figure 7 span several grades. For
335 example, for the grading systems in Table 1, for a four-hour exam, an average student can be
336 awarded all grades between F and C, while a strong student can be awarded all grades
337 between C and A, on a purely chance basis. This means that the reliability and validity of such
338 written exams is low.

339 **Table 1:** Typical conversions to ECTS grading system (A-F) [about here].

340 The analysis in this paper confirms that increasing the length of examinations has a
341 significant effect on reducing the amount of uncertainty in marking. Such results suggest that
342 institutions should strive to use as long a duration as practically possible for written exams.
343 Fatigue as a result of the long duration of an examination may be an issue, but previous
344 research on examinations in other subject areas documented that performance increased with
345 examination length (Jensen et al. (2013), Ackerman & Kanfer (2009)).

346 It can be noted that the results from our research were obtained by studying
347 examinations where the various exercises covered as much of the curriculum as possible
348 (typically, each exercise would be used to test the student’s mastery of a different area of the
349 curriculum). If, for any reason, an examination is designed in such a way that it only aims to
350 test parts of the curriculum (for example, if it includes several exercises that are related to the
351 same part of the curriculum, and no exercises that are related to other parts of the curriculum),
352 then increasing the length of the examination might not result in a decrease in marking

354 uncertainty. In such cases, a longer examination might neither increase its reliability nor its
355 validity, as it would be based on a biased sample of curriculum parts.

356 It can also be noted that increasing the length of an examination will not contribute to
357 reducing marking uncertainty if the examination is not designed to test mastery levels in a
358 time-efficient way. In other words, increasing the length of an examination solely by
359 including lengthy and tedious calculations in the exercises will not increase its reliability or its
360 validity. In order to reduce uncertainty, it is necessary to design examinations in such a way
361 that the time that students spend on answering questions is used as effectively as possible. For
362 example, when an examination question only aims to test the students' mastery of recalling
363 facts, a multiple-choice form may be a better alternative than a lengthy exercise.

364 It can be inferred from the data and from our analysis that there is generally a large
365 degree of uncertainty associated with using a summative assessment of one subject as an
366 indicator of a student's level of mastery of the curriculum in that subject. It is therefore
367 unlikely that one grade will provide an accurate picture of a student's abilities. In order to
368 reduce this uncertainty, it is necessary to have access to a larger number of examination
369 grades. Typically, a student takes between 25 and 75 exams within the framework of an
370 educational program, and, because of the law of large numbers, the average grade based on all
371 the individual course grades will normally reflect the student's ability with much less
372 uncertainty than an individual grade.

373 Of course, there are other possible challenges to the validity and reliability of an
374 average grade based on several exams, and they could be the subject of further research. Such
375 challenges might include differences in strictness levels from one examiner to another, from
376 one subject area to another, and from one institution to another. Another challenge may be
377 that some examiners and institutions use norm-referenced grades (i.e., grades that to a greater

379 extent reflect where the examination paper stands in comparison with the level of the other
380 examination papers), rather than criterion-referenced grades (i.e., grades that reflect the
381 intrinsic quality of the paper, independently of the rest of the group). Although this practice
382 has been pinpointed as unethical (Sadler 2009), it is commonly used in various educational
383 settings, for example in order to prevent “grade inflation” (Cliffordson, 2008). It is therefore
384 important that further research encompassing a broad variety of examinations and
385 examination results ascertains the degree of integrity of the grading systems, i.e., the extent to
386 which they are criterion-based rather than norm-based.

387 Appendix

388 Beta Regression model

389 In this section, we provide a more detailed description of the beta regression model
390 used in this paper. Let S denote a stochastic variable for the score on an arbitrary exercise for
391 an arbitrary student, normalized to the $[0, 1]$ interval. We assume that S is beta-distributed

$$392 \quad \pi(s) = \frac{1}{B(a, b)} s^{a-1} (1-s)^{b-1}, \quad a > 0, b > 0$$

393 where

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

where $\Gamma(x)$ is the Gamma function.

We now link the expectation of the beta distribution to a linear predictor of some covariates using a link function. We use the reparameterization

$$\mu = \frac{a}{a+b}, \quad 0 < \mu < 1$$

$$\phi = a+b, \quad \phi > 0$$

to arrive at

$$E(S) = \mu$$

Further, we get

$$Var(S) = \frac{\mu(1-\mu)}{1+\phi}$$

where ϕ is known as the precision parameter, since for a fixed μ , the larger ϕ , the smaller the variance in S . We link μ to a the linear predictor η using the logit-link function

$$\mu = \frac{e^\eta}{1+e^\eta}$$

Such a model is called beta regression. We use the following linear predictor in the beta regression model

$$\eta = k + M + D$$

412 where k is the fixed effect interception and M and D random effects representing the
413 variability in the level of mastery (M) of students taking the exam and the difficulty level (D)
414 of the exercises in the exam, respectively. We assume that M and D are normally distributed
415 with zero expectations and variances $1/\tau_M$ and $1/\tau_D$, respectively. The parameters τ_M and
416 τ_D are the inverse of the variance, which is referred to as precision. We assume that, given the
417 random effects, the observations (score for a particular exercise for a particular student) are
418 independent. We use a Bayesian approach and add prior distributions to the unknown
419 parameters ϕ , k , τ_M and τ_D . For details about the prior distributions, we refer to the INLA
420 web page (Rue, 2014).

421 *Estimated values of parameters in a regression model*

422 Table 2 shows properties of the posterior distributions of the variables ϕ , k , τ_M and τ_D
423 for one of the 23 exams.

424 **Table 2:** Properties of the posterior distributions for the variables ϕ , k , τ_M and τ_D
425 [about here].

426 We see that k is less than zero, showing that, on average, the students scored below 0.5 on the
427 exercises in this particular exam. We also observe that the largest estimation uncertainty is in
428 the estimation of τ_D , variability in levels of difficulty on the exercises.

Grade	Score intervals 1	Score intervals 2
A	90 – 100	90 – 100
B	80 – 89	80 – 89
C	70 – 79	60 – 79
D	60 – 69	50 – 59
E	50 – 59	40 – 49
F	0 – 49	0 – 39

Table 1: Typical conversions to ECTS grading system (A-F).

Variable	Mean	Stdev	5% quantile	50% quantile	95% quantile
k	-0.461	0.165	-0.732	-0.461	-0.191
ϕ	1.168	0.024	1.128	1.168	1.208
τ_M	1.357	0.135	1.148	1.349	1.592
τ_D	2.478	0.821	1.339	2.374	3.984

Table 2: Properties of the posterior distributions for the variables ϕ , k , τ_M and τ_D .

References

- 447 Abdul-Rahman, Syariza, Edmund Burke, Andrzej Bargiela, Barry McCollum, and Ender Özcan. 2014. "A
 448 constructive approach to examination timetabling based on adaptive decomposition and
 449 ordering." *Annals of Operations Research* 218 (1):3-21. doi: 10.1007/s10479-011-0999-8.
- 450 Ackerman, Phillip L., and Ruth Kanfer. 2009. "Test Length and Cognitive Fatigue: An Empirical
 451 Examination of Effects on Performance and Test-Taker Reactions." *Journal of Experimental*
 452 *Psychology: Applied* 15 (2):163-181.
- 453 Admiraal, Wilfried, Mark Hoeksma, Marie-Therese van de Kamp, and Gee van Duin. 2011.
 454 "Assessment of Teacher Competence Using Video Portfolios: Reliability, Construct Validity,
 455 and Consequential Validity." *Teaching and Teacher Education: An International Journal of*
 456 *Research and Studies* 27 (6):1019-1028.
- 457 Agasisti, Tommaso, and Francesca Bonomi. 2014. "Benchmarking universities' efficiency indicators in
 458 the presence of internal heterogeneity." *Studies in Higher Education* 39 (7):1237-1255. doi:
 459 10.1080/03075079.2013.801423.
- 460 Allais, Stephanie. 2014. "A critical perspective on large class teaching: the political economy of
 461 massification and the sociology of knowledge." *Higher Education* 67 (6):721-734. doi:
 462 10.1007/s10734-013-9672-2.
- 463 Bird, Fiona L., and Robyn Yucel. 2013. "Improving marking reliability of scientific writing with the
 464 Developing Understanding of Assessment for Learning programme." *Assessment &*
 465 *Evaluation in Higher Education* 38 (5):536-553. doi: 10.1080/02602938.2012.658155.
- 466 Blanco-Ramírez, Gerardo, and Joseph B. Berger. 2014. "Rankings, accreditation, and the international
 467 quest for quality Organizing an approach to value in higher education." *Quality Assurance in*
 468 *Education: An International Perspective* 22 (1):88-104. doi: 10.1108/QAE-07-2013-0031.
- 469 Boyas, Elise, Lois D. Bryan, and Tanya Lee. 2012. "Conditions affecting the usefulness of pre- and
 470 post-tests for assessment purposes." *Assessment & Evaluation in Higher Education* 37
 471 (4):427-437. doi: 10.1080/02602938.2010.538665.
- 472 Burton, Richard F. 2006. "Sampling Knowledge and Understanding: How Long Should a Test Be?"
 473 *Assessment & Evaluation in Higher Education* 31 (5):569-582.
- 474 Cliffordson, Christina. 2008. "Differential Prediction of Study Success across Academic Programs in
 475 the Swedish Context: The Validity of Grades and Tests as Selection Instruments for Higher
 476 Education." *Educational Assessment* 13 (1):56-75.
- 477 Davis, L. E., M. C. Harrison, A. S. Palipana, and J. P. Ward. 2005. "Assessment-driven learning of
 478 mathematics for engineering students." *International Journal of Electrical Engineering*
 479 *Education* 42 (1):63-72.
- 480 Delen, Erhan. 2015. "Enhancing a Computer-Based Environment with Optimum Item Response
 481 Time." *Eurasia Journal of Mathematics, Science and Technology Education* 11 (6):1457-1472.
- 482 DeVellis, Robert F. . 2012. *Scale Development: Theory and Applications*. 3rd ed. London: Sage.
- 483 Dobson, Annette J. , and Adrian G. Barnett. 2008a. *An Introduction to Generalized Linear Models,*
 484 *Texts in Statistical Science*. Boca Raton, FL: Chapman & Hall/CRC Press.
- 485 Dobson, Annette J., and Adrian G. Barnett. 2008b. *Introduction to Generalized Linear Models*. 3rd ed.
 486 London: Chapman and Hall/CRC.
- 487 Hambleton, Ronald K, Hariharan H Swaminathan, and Jane Rogers. 1991. *Fundamentals of item*
 488 *response theory*. Newbury Park, CA: Sage.
- 489 Harlen, Wynne. 2005. "Trusting teachers' judgement: research evidence of the reliability and validity
 490 of teachers' assessment used for summative purposes." *Research Papers in Education* 20
 491 (3):245-270. doi: 10.1080/02671520500193744.

- 493 Hughes, Clair. 2013. "A case study of assessment of graduate learning outcomes at the programme,
494 course and task level." 38:492-506. doi: 10.1080/02602938.2012.658020.
- 495 Irwin, Brian, and Stuart Hepplestone. 2012. "Examining increased flexibility in assessment formats."
496 *Assessment & Evaluation in Higher Education* 37 (7):773-785. doi:
497 10.1080/02602938.2011.573842.
- 498 Jensen, Jamie L., Dane A. Berry, and Tyler A. Kummer. 2013. "Investigating the Effects of Exam Length
499 on Performance and Cognitive Fatigue." *PLoS ONE* 8 (8):1-9. doi:
500 10.1371/journal.pone.0070270.
- 501 Kuo, Bor-Chen, Muslem Daud, and Chih-Wei Yang. 2015. "Multidimensional Computerized Adaptive
502 Testing for Indonesia Junior High School Biology." *Eurasia Journal of Mathematics, Science
503 and Technology Education* 11 (5):1105-1118.
- 504 Lord, Frederic M. 1952. *A Theory of Test Scores* Vol. 7, *Psychometric Monograph*. Richmond, VA.
- 505 Lord, Frederic M. 1953. "The relation of test score to the trait underlying the test." *Educational and
506 Psychological Measurement* 13:517-548.
- 507 Lord, Frederic M. 1980. *Applications of Item Response Theory to Practical Testing Problems*. London:
508 Routledge.
- 509 Lord, Frederic M., and Melvin R. Novick. 1968. *Statistical theories of mental test scores*. Reading, MA:
510 Addison-Wesley.
- 511 Mumford, Christine L. 2010. "A multiobjective framework for heavily constrained examination
512 timetabling problems." *Annals of Operations Research* 180 (1):3-31. doi: 10.1007/s10479-
513 008-0490-3.
- 514 Muraki, E. . 1997. "A generalized partial credit model." In *Handbook of modern item response theory*,
515 edited by W. van der Linden and R. K. Hambleton, 153-164. New York: Springer.
- 516 PARSCALE (Version 4.1). Scientific Software International, Lincolnwood, IK.
- 517 Nijmegen, Radboud University. 2011. "Conversion of Grades." Accessed 17 July 2014.
518 http://www.ru.nl/io/english/general_0/document/.
- 519 Palmer, Edward J., Paul Duggan, Peter G. Devitt, and Rohan Russell. 2010. "The modified essay
520 question: its exit from the exit examination?" *Medical Teacher* 32 (7):e300-e307. doi:
521 10.3109/0142159X.2010.488705.
- 522 Rue, H. 2014. "The R-INLA project." Accessed 17 July 2014. <http://www.r-inla.org/>.
- 523 Russell, Jill, Lewis Elton, Deborah Swinglehurst, and Trisha Greenhalgh. 2006. "Using the online
524 environment in assessment for learning: a case-study of a web-based course in primary
525 care." *Assessment & Evaluation in Higher Education* 31 (4):465-478. doi:
526 10.1080/02602930600679209.
- 527 Sadler, D. Royce. 2009. "Grade Integrity and the Representation of Academic Achievement." *Studies
528 in Higher Education* 34 (7):807-826.
- 529 Simpson, Lucy, and Jo-Anne Baird. 2013. "Perceptions of trust in public examinations." *Oxford
530 Review of Education* 39 (1):17-35. doi: 10.1080/03054985.2012.760264.
- 531 Vu, Nv, A. Baroffio, P. Huber, C. Layat, M. Gerbase, and M. Nendaz. 2006. "Assessing clinical
532 competence: a pilot project to evaluate the feasibility of a standardized patient -- based
533 practical examination as a component of the Swiss certification process." *Swiss Medical
534 Weekly* 136 (25-26):392-399.
- 535 Vukasovic, Martina. 2013. "Change of higher education in response to European pressures:
536 conceptualization and operationalization of Europeanization of higher education." *Higher
537 Education* 66 (3):311-324. doi: 10.1007/s10734-012-9606-4.
- 538 Westerheijden, Don F., Bjørn Stensaker, Maria J. Rosa, and Anne Corbett. 2014. "Next Generations,
539 Catwalks, Random Walks and Arms Races: Conceptualising the development of quality
540 assurance schemes." *European Journal of Education* 49 (3):421-434. doi:
541 10.1111/ejed.12071.
- 542 William, Dylan. 1996. "Standards in examinations: a matter of trust?" *The Curriculum Journal* 7
543 (3):293-306.

545 Wittek, Line, and Tone Kvernbekk. 2011. "On the problems of asking for a definition of quality in
546 education." *Scandinavian Journal of Educational Research* 55 (6):671-684. doi:
547 10.1080/00313831.2011.594618.

Figure 1

[Download source file \(7.47 kB\)](#)

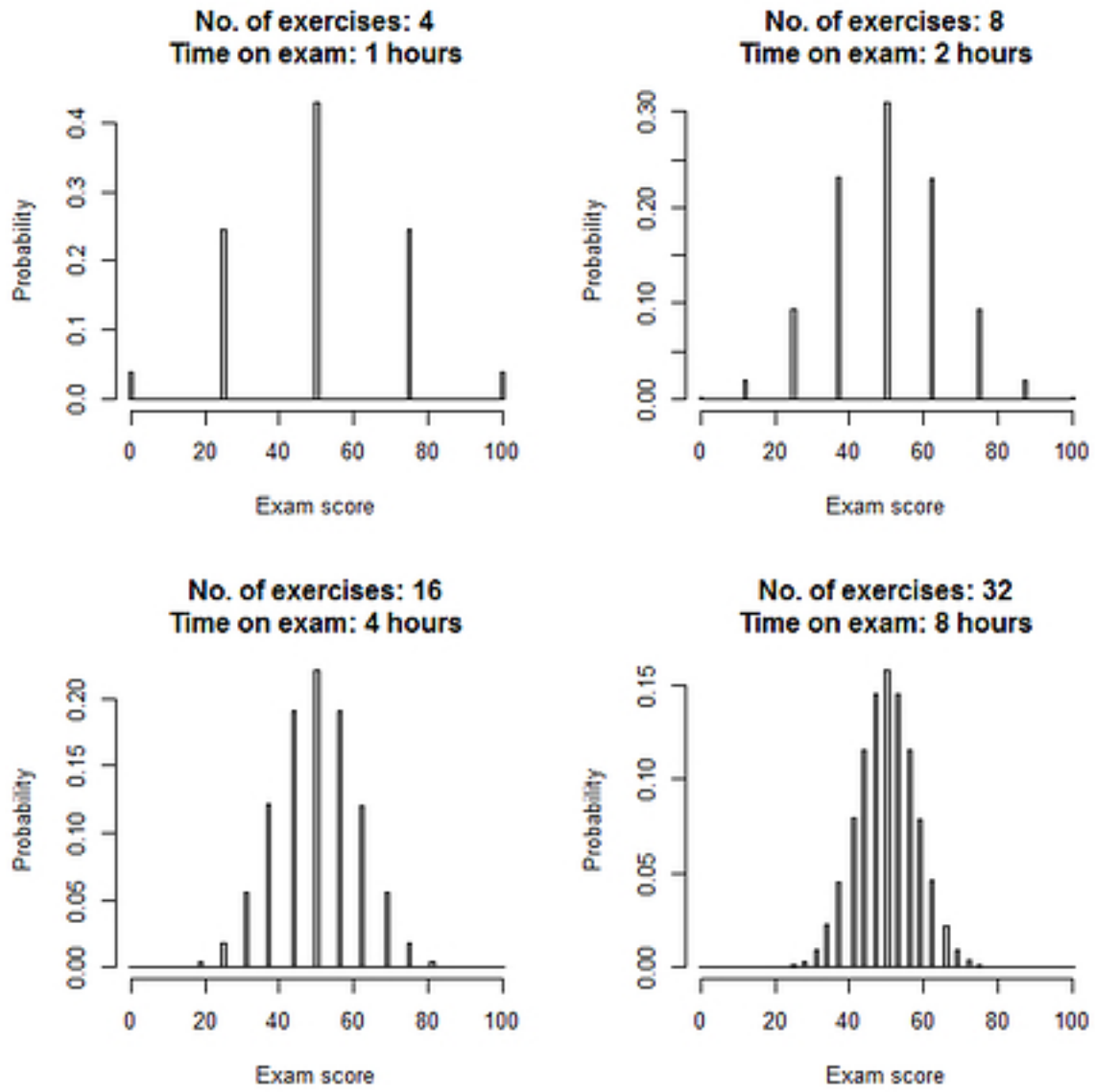


Figure 2

[Download source file \(6.97 kB\)](#)

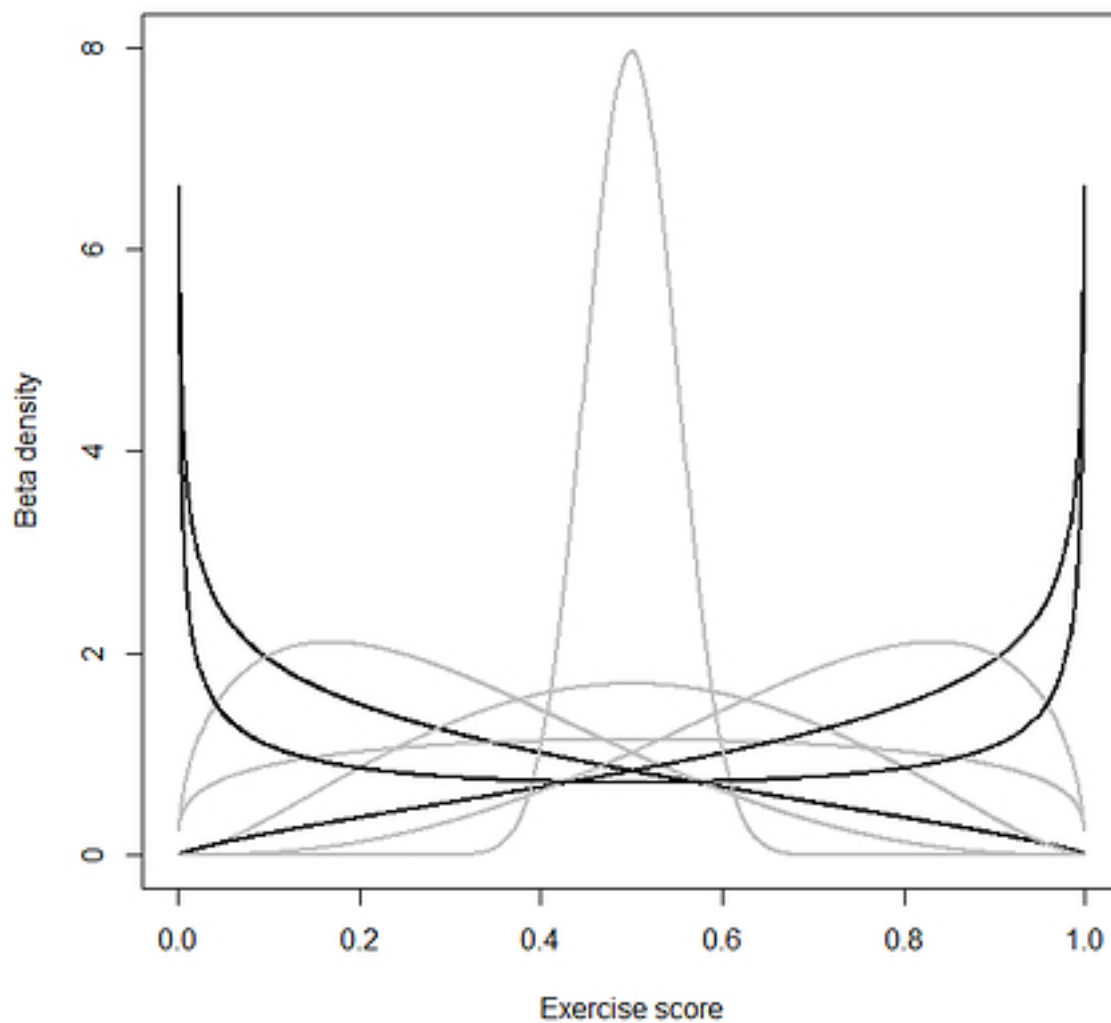


Figure 3

[Download source file \(4.1 kB\)](#)

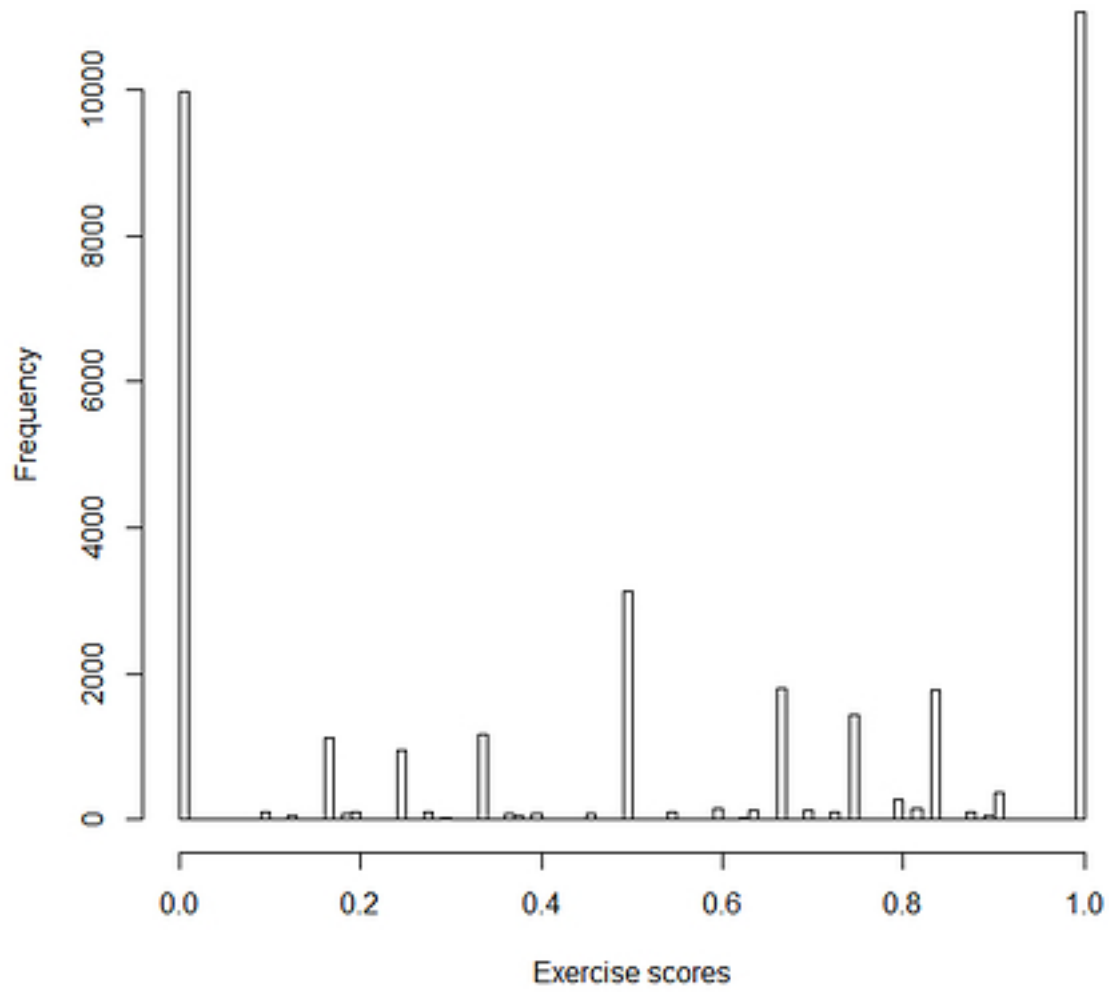


Figure 4

[Download source file \(6.95 kB\)](#)

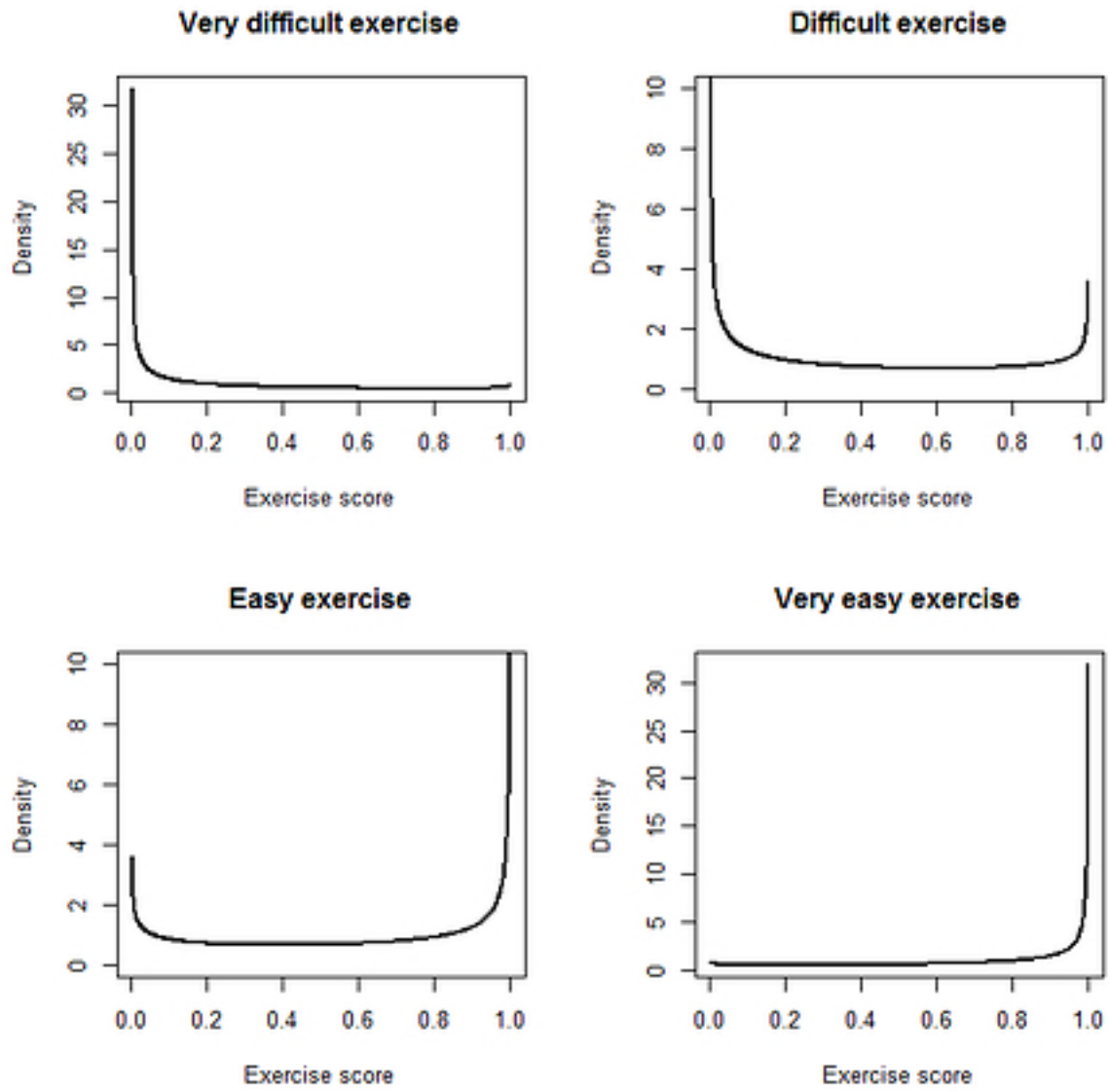


Figure 5

[Download source file \(10.98 kB\)](#)

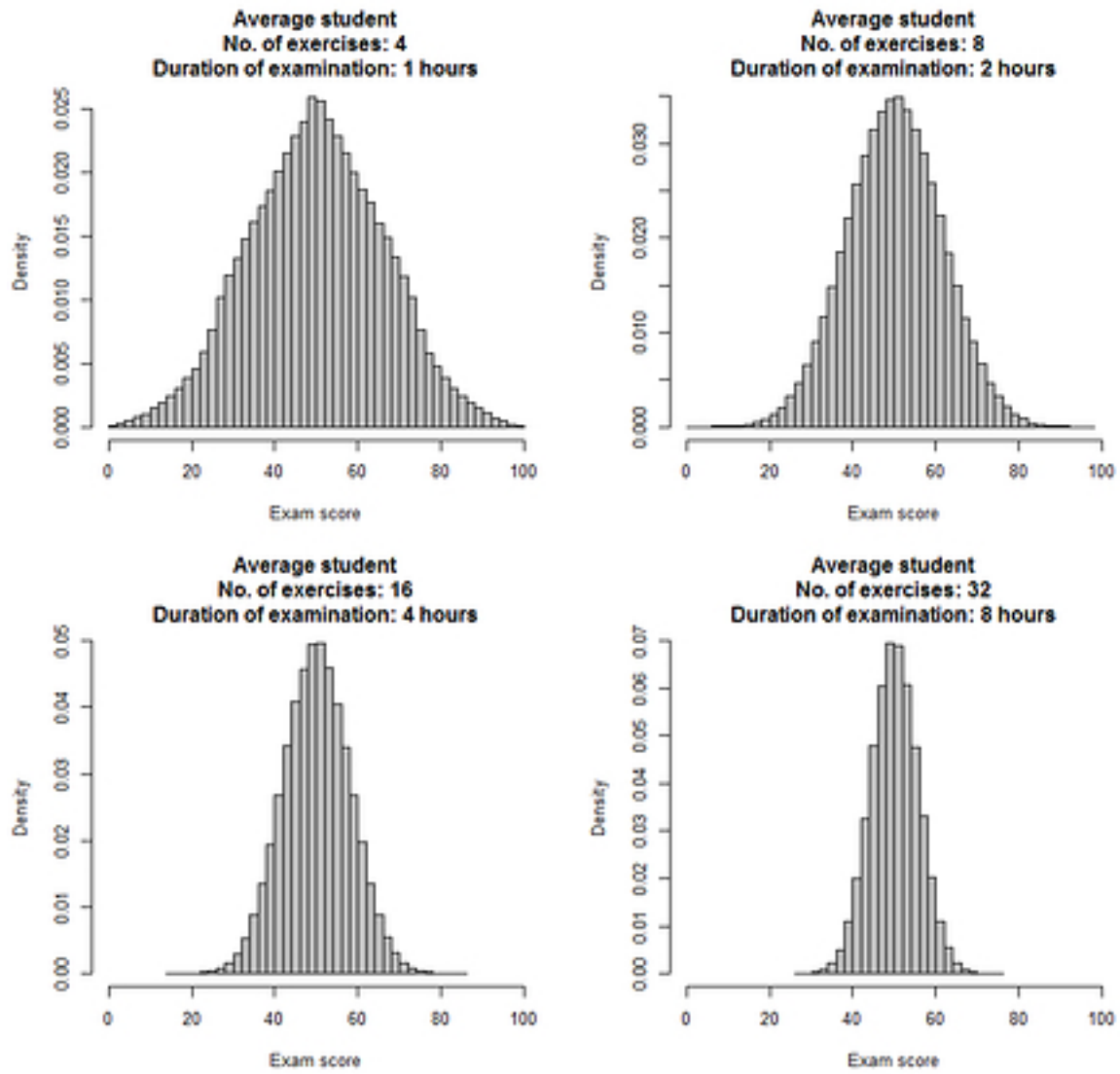


Figure 6

[Download source file \(10.72 kB\)](#)

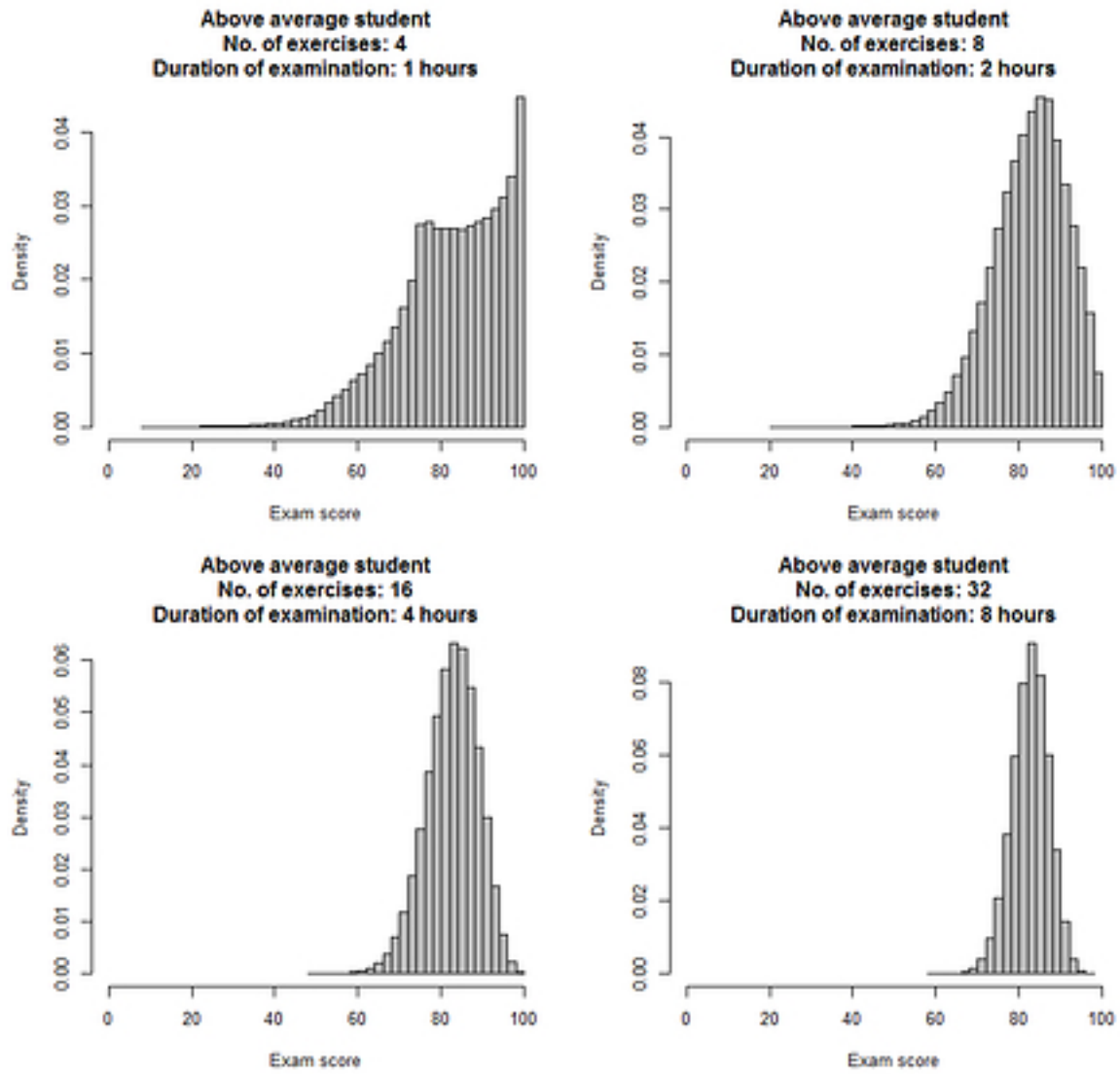
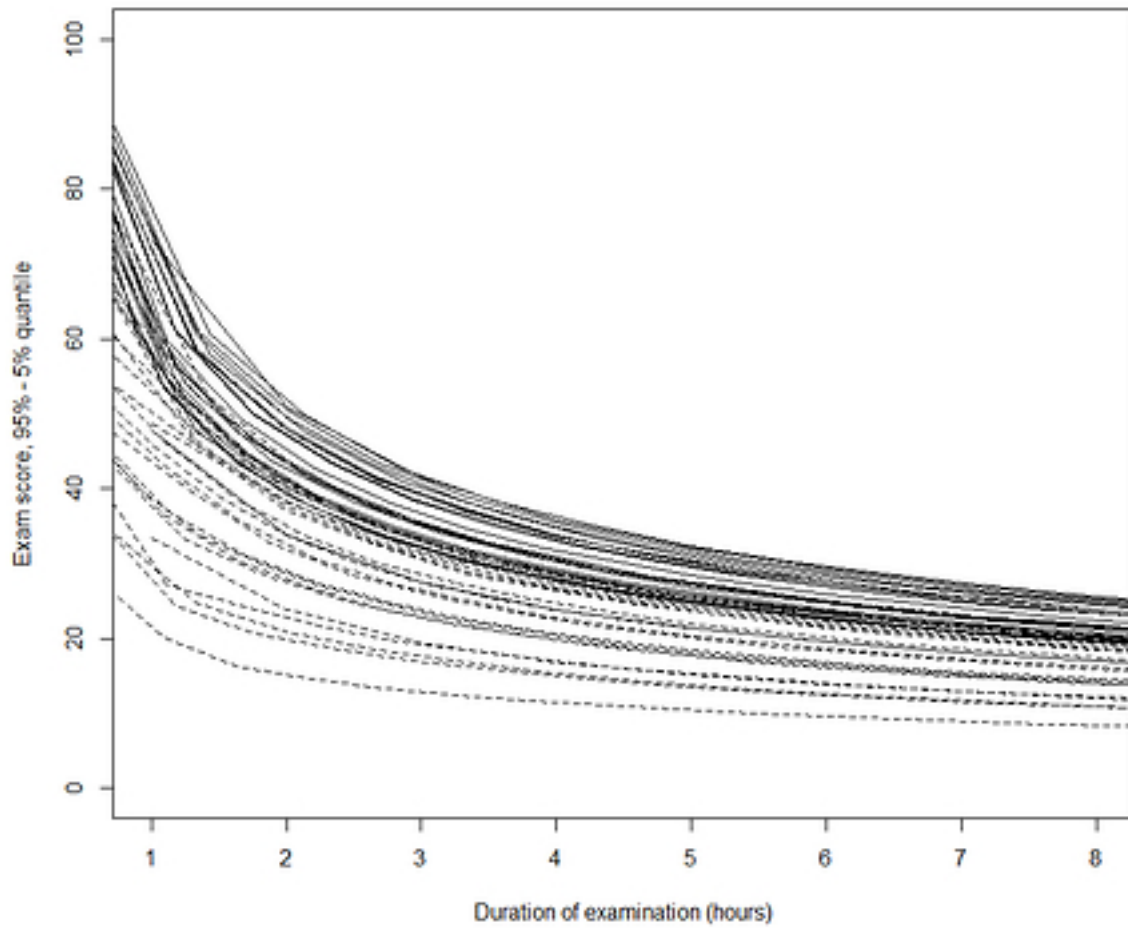


Figure 7

[Download source file \(13.57 kB\)](#)



Manuscript body

Manuscript body 1 - [Download source file \(86.36 kB\)](#)

Figures

Figure 1 - [Download source file \(7.47 kB\)](#)

Figure 2 - [Download source file \(6.97 kB\)](#)

Figure 3 - [Download source file \(4.1 kB\)](#)

Figure 4 - [Download source file \(6.95 kB\)](#)

Figure 5 - [Download source file \(10.98 kB\)](#)

Figure 6 - [Download source file \(10.72 kB\)](#)

Figure 7 - [Download source file \(13.57 kB\)](#)