

Article

An Inexpensive Retrospective Standard Setting Method Based on Item Facilities

Mclachlan, John Charles, Robertson, K. Alex, Weller, Bridget and Sawdon, Marina

Available at <http://clock.uclan.ac.uk/35832/>

Mclachlan, John Charles ORCID: 0000-0001-5493-2645, Robertson, K. Alex, Weller, Bridget and Sawdon, Marina (2020) An Inexpensive Retrospective Standard Setting Method Based on Item Facilities. BMC Medical Education .

It is advisable to refer to the publisher's version if you intend to cite from the work.

For more information about UCLan's research in this area go to <http://www.uclan.ac.uk/researchgroups/> and search for <name of research Group>.

For information about Research generally at UCLan please go to <http://www.uclan.ac.uk/research/>

All outputs in CLoK are protected by Intellectual Property Rights law, including Copyright law. Copyright, IPR and Moral Rights for the works on this site are retained by the individual authors and/or other copyright owners. Terms and conditions for use of this material are defined in the <http://clock.uclan.ac.uk/policies/>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

**An Inexpensive Retrospective Standard Setting Method Based on
Item Facilities**

John C McLachlan¹, K. Alex Robertson², Bridget Weller¹ and Marina Sawdon

- ¹University of Central Lancashire, Harrington Building, 11 Victoria St, Preston PR1 7QS.
- ²CNTW Trust, Hopewood Park Hospital, Ryhope, Sunderland, SR2 0NB.
- ³School of Medicine, University of Sunderland, Chester Rd, Sunderland SR1 3SD.

Correspondence: jcmclachlan1@uclan.ac.uk

20 **Abstract**

21 **Background**

22 Standard setting is one of the most challenging aspects of assessment in high-stakes
23 healthcare settings. The Angoff methodology is widely used, but poses a number of challenges,
24 including conceptualisation of the just-passing candidate, and the time-cost of implementing
25 the method. Cohen methodologies are inexpensive and rapid but rely on the performance of
26 an individual candidate. A new method of standard setting, based on the entire cohort and
27 every item, would be valuable.

28 **Methods**

29 We identified Borderline candidates by reviewing their performance across all assessments in
30 an academic year. We plotted the item scores of the Borderline candidates in comparison with
31 Facility for the whole cohort and fitted curves to the resulting distribution.

32 **Results**

33 It is observed that for any given Item, an equation of the form

34
$$y \approx C \cdot e^{Fx}$$

35 where y is the Facility of Borderline candidates on that Item, x is the observed Item

36 Facility of the whole cohort, and C and F are constants,

37 predicts the probable Facility for Borderline candidates over the test, in other words, the cut

38 score for Borderline candidates. We describe ways of estimating C and F in any given

39 circumstance, and suggest typical values arising from this particular study: that $C = 12.3$ and $F =$

40 0.021 .

41 **Conclusions**

42 C and F are relatively stable, and that the equation

43
$$y = 12.3.e^{0.021x}$$

44 can rapidly be applied to the item Facility for every item. The average value represents the cut
45 score for the assessment as a whole. This represents a novel retrospective method based on
46 test takers.

47 Compared to the Cohen method which draws on one score and one candidate, this method
48 draws on all items and candidates in a test. We propose that it can be used to standard set a
49 whole test, or a particular item where the predicted Angoff score is very different from the
50 observed Facility.

51

52 Keywords: standard-setting, retrospective, cost, rapid, exponential

53

54

55

56

57

58

59

60

61

62

63 **Background**

64 Standard setting is both important and problematic in medical education. The Angoff method¹
65 is widely used for standard setting selected-response items in high stakes settings such as the
66 General Medical Council tests for non-UK, non-EU doctors wishing to practice in the UK, and
67 USMLE Step 1, yet its use poses a number of challenges.

68 Perhaps the most significant of these is the requirement that assessors conceptualise a
69 particular kind of candidate, often described as the 'minimally competent' or 'Borderline'
70 candidate. In the context of Angoff standard setting, 'Borderline' generally represents a
71 'Borderline pass', and it is in this sense that we use it here.

72 Whichever form of words is used, assessors may have very different ideas of what that class of
73 candidates represents. This is compounded by the fact that subject specialists among the
74 assessors may lack generalist knowledge,² or lack awareness of what particular level
75 candidates would appropriately have achieved.

76 As a consequence, a minimum number of assessors may be required, and this in itself poses
77 practical problems in identifying a sufficient number of assessors with sufficient expertise in
78 the subject, and indeed experience in using the Angoff method. One safety-net option is to use
79 the Hofstee compromise method³ if any 'Angoffed' assessment fails a 'Reality Check'.⁴

80 A particular tendency of novice assessors is 'reversion to the mean', where they tend to award
81 Angoff scores of around 50% rather than using the full scale range. This results in a low
82 correlation between the predicted Angoff value and the observed Facility (where Facility is the
83 percentage of candidates answering correctly) of the items.

84 Some of the same considerations apply to Ebel standard setting⁵. Again, the just-passing
85 candidate is difficult to conceptualise, and a panel of experts is required to carry out the
86 required classification.

87 An inexpensive alternative is to use either the Cohen method⁶, which derives the cut score
88 from a multiple of the 95th centile candidate, or the similar modified Cohen method⁷, which
89 relies on the 90th centile candidate. These methods are quick to implement, and do not require
90 the input of expensive staff time. However, they may be criticised on the basis that they rely
91 on the score of an individual candidate (or in the case of ties, a small number of candidates).
92 We return to this issue in the Discussion.

93 However, it is possible that assessments vary more in difficulty than does the ability of the
94 cohort, since medical students are highly selected for academic ability prior to entry. In this
95 case, the difficulty of the assessment may be the key variable, and the cumulative Facility of
96 the items is a guide to this.

97 Of course, Facility represents the whole cohort performance, rather than the performance of
98 the Borderline candidates. We hypothesised that for good quality One-Best-of-Five MCQs, the
99 relationship between Facility for the whole cohort, and the Facility for Borderline candidates,
100 would be curvilinear in nature, with the difference between them approaching zero as the
101 Facility approaches 100% and 20%. This is because if the entire cohort scores an item correctly,
102 then so will the Borderline candidates, and if the best candidates do no better than guessing,
103 then neither will the Borderline candidates.

104 In this study we therefore attempted to explore the effect of classifying different numbers of
105 students as 'Borderline' in comparison with the cohort as a whole. Classification was carried
106 out based on performance across the whole range of modules undertaken by the students as
107 described in the Methods.

108 The exact nature of the relationship between whole cohort and Borderline Facility will depend
109 on the proportion of Borderline candidates in the class, and we discuss ways in which this
110 might be estimated.

111 Where such a relationship emerges, it would be of value in assisting novice Angoff assessors in
112 estimating the performance of Borderline candidates for an item which had been used before.

113 It could also be used for adjusting any items where the discrepancy between the predicted
114 Angoff value and the observed Facility for that item is greater than seems plausible.

115 More importantly, the relationship could be used by itself as a standard setting method in
116 conditions in which Angoff or similar methods were not practical: for instance, if too few
117 subject matter experts were available to form an assessor panel, or where the resource costs
118 of using the Angoff method were too high. This would then be a retrospective method based
119 on test takers, rather than a prospective method based on test items.

120 The purpose of this study is to show proof of concept and although the analyses were carried
121 out locally, we believe our results would be adaptable and of interest to other settings outside
122 our school.

123

124 **Methods**

125 The analyses were based on a cohort of students at a UK Medical School. The number of
126 students involved was in the region of one hundred, but the exact number is not disclosed
127 since this may enable the particular cohort to be identified. Student names were never used in
128 the analysis, and student numbers were re-coded automatically so anonymity was preserved.

129 The data were used retrospectively, and this analysis has played no part in summative
130 decisions.

131 All calculations were carried out, and graphs plotted, using Microsoft Excel®.

132 Ethical approval for the project on this basis and for publication of results was granted by the

133 relevant University Ethical Approval Committee (approval code STEMH 1058).

134 The First Year medical student course in question contains three modules each year. Modules

135 1 and 2 address declarative knowledge, and contained a total of three papers, and Module 3

136 involves an OSCE skills assessment. Standards are set for Modules 1 and 2 by the modified

137 Angoff method, and for Module 3 by Borderline Regression. Module 3 had an additional

138 conjunctive condition which was that candidates had to pass at least 75% of the OSCE stations.

139 The anonymised candidates were classified by their performance in each of their modules,

140 with reference to the Standard Error of Measurement (SEM) of the exam and given a

141 corresponding score as described in Table 1.

Description	Boundaries	Score
Possible Borderline	between 1 and 2 SEM <i>above</i> the cut score	0.5 points
Probable Borderline	within 1 SEM of the cut score	1 point
Definite Borderline	between 1 and 2 SEM <i>below</i> the cut score	2 points

142

143 Table 1. Boundaries and score allocations for various Borderline categories

144 For the skills module, candidates who had failed 25% of the OSCE stations were also

145 considered Borderline and scored 1 point. See Table 2 for the distribution of scores in this

146 particular cohort.

'Borderline' Points	% of Cohort
0.5	12

1	9
1.5	5
2	4
2.5	2
3	2
3.5	1
4	1
4.5	2
5	2
5.5	0
6	0
6.5	2
All	43

147

148 Table 2. Proportions of candidates scoring various numbers of 'borderline' points as calculated
 149 in the text. Those scoring 0.5 points lay between 1 and 2 standard errors of measurement
 150 above the cut score.

151 Obviously, a candidate could gather points from more than one module. Points ranged from
 152 0.5 for approximately 12% of the cohort, to 6.5 for a few individuals. In total, approximately
 153 43% of the cohort had points. However, a total score of 0.5 points represented a performance
 154 between one and two SEM *above* the cut score in one Module only, which is likely to be the
 155 result of chance for an otherwise satisfactory candidate.

156 The Facility of the Borderline candidates for each item was plotted against the cohort facility,
 157 first for all Borderline candidates, then for a variety of different score combinations. Curves
 158 were fitted to these plots using the trendline function in Excel. This allowed us to explore the
 159 stability of the curve in terms of it's constants.

160 A standardised 'exponential curve' showing the relationship between the Facility of the
161 Borderline candidates and that of the cohort as a whole was then developed. It was
162 retrospectively applied to a total of 26 previous MCQ-style assessments over the last four
163 years of the Undergraduate medical programme as a standard setting method. Cut scores
164 were calculated on the basis of this exponential curve and compared to those which had been
165 obtained by a full Angoff procedure. Cohen and Modified Cohen method cut scores were also
166 calculated for each exam, although in practice only Angoff methods had been used. From
167 these, the proportion of candidates who would have failed each assessment by each method
168 were calculated. These results were plotted against the average score in each assessment.

169 A further theoretical calculation showing the effect of varying the proportion of candidates
170 classed as Borderline in the cohort was also carried out.

171

172 **Results**

173 A plot was constructed of the Facility of (a) (shown in Table 3) each Item in the test compared
174 to the score of all Borderline candidates, and the trendline added (Figure 1). As can be seen,
175 as predicted a curved trendline, approaching zero at Item Facilities of 0 and 100 is indeed
176 observed. The equation for this curve is shown on Figure 1, and is of the form

$$177 \quad y \approx C.e^{Fx}$$

178 Where y is the Facility of Borderline candidates, x is the observed Facility of the cohort as a
179 whole, and C and F are constants.

180 This process was repeated for various combinations of possible Borderline candidates, to
181 explore how stable this curve was in terms of its constants. As listed in Table 3, these
182 combinations were (b) excluding those who had scored only 0.5 points (i.e. had scored

183 between 1 and 6.5 points) on the basis that a score of 0.5 (between 1 and 2 SEM above the
 184 cut-score in a single module) probably represents noise in the performance of otherwise
 185 capable students (c) students who fell between 1.5 and 6.5 points, a more stringent
 186 interpretation of Borderline (d) candidates scoring between 1 and 5.5 points (excluding those
 187 candidates who would be clear fails and (e) showing only scores on different assessments from
 188 that shown in the plot, so that there is no element of circularity in the reasoning. The results
 189 are shown in Table 3.

	Range of Borderline scores	% of cohort	C	F
(a)	All possible Borderlines	43	13.125	0.021
(b)	1 - 6.5	32	12.756	0.0208
(c)	1.5 - 6.5	23	13.562	0.0192
(d)	1 - 5.5	27	12.6	0.0218
(e)	0.5 – 6.5 (Excluding Source)	28	12.964	0.0209
Exponential			12.3	0.021

190

191 Table 3. The values observed for curves of the form of Equation 2. A family of curves could be
 192 selected for the 'Standard' values; this particular combination was chosen because the
 193 difference from the Facility is zero at 20% and 100%.

194

195 As can be seen, these curves are all relatively consistent in terms of their constants. On this
 196 basis, a standard exponential curve was calculated on the basis that it intercepts Facility
 197 exactly at 20% and 100%. This curve had the constant values

198

$$y = 12.3e^{0.021x}$$

199 This equation can therefore be applied to the Facility of any individual item in a test and gives
200 the expected score for a Borderline candidate for that item. The average of these values is
201 therefore the cut score for the test as a whole.

202 For the 26 assessments over the four-year period of this study, the proportion of candidates
203 who would have failed each assessment by Angoff, Cohen, Modified Cohen and use of the
204 exponential equation were calculated. Average values for these are shown in Table 4.

	Angoff	Exponential	Cohen	Modified Cohen
Mean	15.13484	14.08984	25.36721	13.39841
Standard Deviation	9.528302	5.759773	10.97905	7.06995

205

206 Table 4. Percentage of 'Fail' students over a total of 26 exams, using 4 different standard
207 setting methods.

208

209 These results were also plotted against the average score in each assessment, as shown in
210 Figure 2. By all four methods, there is a linear relationship between the average score in the
211 test, and the percentage of candidates who fail – as is to be expected, the higher the average
212 score, the fewer candidates fail. However, the exponential curve method gives a much more
213 stable result: the slope is much shallower than those of the other three methods. This accords
214 with the lower Standard Deviation for the exponential curve method, as shown in Table 4.

215 A reasonable question would be to ask if Facility is dependent on the proportion of Borderline
216 candidates in the test. We modelled the impact of changing this proportion, and the impact on
217 the Facility of the test as a whole was small: for instance, the difference in cohort Facility when
218 15% versus 35% of the cohort were classed as Borderline was 3%.

219 This suggests that the overall performance of a cohort of students may be relatively stable to
220 changes in the proportions of Borderline candidates, a point which we will return to in the
221 Discussion.

222 **Discussion**

223 For a context in which candidates have already undertaken multiple assessments previously
224 standard-set by some conventional means, repeating the approach described here is possible
225 and relatively straightforward. Candidates can be classified as Borderline on the basis of their
226 performance across all assessments, and the equivalent of Figure 1 plotted. We predict that a
227 curve of the same form, and with constant values close to that of the standard exponential
228 curve will be observed.

229 *Ways of using the exponential curve*

230 The exponential curve equation can be used as a rapid and inexpensive primary standard
231 setting method, in the same settings as Cohen and modified Cohen methodologies are
232 currently employed. The Facility of each item is calculated in most item-banking applications.
233 These Facilities can be exported to a spreadsheet and the exponential equation copied into the
234 adjacent cells. Using the 'fill down' command in Excel, this takes seconds to do. The average
235 value of the exponential equation outcomes is the cut score for the test as a whole.

236 We believe that the exponential equation is preferable to both Cohen methodologies, because
237 it is derived from the results of all items and all candidates, rather than the results of one
238 candidate in the test. In addition, it is more stable to changes in the average score than either
239 Cohen methodology, or even Angoff approaches.

240 Alternatively, it could be used in conjunction with Angoff methods, to adjust the cut score
241 value of individual items where there is a major discrepancy between the pre-calculated
242 Angoff value and the observed Facility.

243 Compared to the Angoff method itself, this method avoids the need for assembling an expert
244 reference group, and the time-consuming and contentious process of estimating an Angoff
245 value for every item in the test. It is very much less costly terms of staff time to carry out, and
246 may bring significant opportunity cost benefits.

247 The method may be useful in standard-setting new kinds of items, such as Very Short Answer
248 Items, which have been observed to have lower Facilities than MCQs⁸, and where Angoff
249 values calculated by the usual method may not be appropriate.

250 *Challenges to this approach*

251 The key issue is the stability of the constants C and F under different conditions.

252 Two conditions must be met for C and F to be relatively stable. The first is that the variance in
253 difficulty of the assessments should be greater than the variance in ability of the candidates.

254 It has indeed been demonstrated for medical students that “test-difficulty is a major source of
255 variation while cohort and education effects probably are minor”⁹. Similarly, Cohen-

256 Schotenaus and van der Vleuten concluded that “the most probable cause (of pass mark
257 variability) is variability in test difficulty across different tests, both within and across courses”.

258 This may be due to the fact that medical students are highly selected at entry to be at the top
259 end of the academic ability spectrum.

260 The second is that the proportion of Borderline candidates should be a relatively stable
261 proportion of the cohort as a whole. Again, the highly selected nature of medical students
262 suggests this is a reasonable expectation. In any case, we have observed that significant

263 variations in the proportion of Borderline candidates bring about only small changes in cohort
264 Facility.

265 As a consequence, it is not unreasonable to think that C and F may vary only within a narrow
266 range. Facility of items in a test as a whole may well be the most important variable in medical
267 exams as previous authors have indicated.

268 **Conclusions**

269 This novel standard setting method offers an inexpensive and easy to implement alternative to
270 existing methods, which takes account of all candidates and all items. It is more stable to
271 changes in mean score in the exam than alternative methods.

272

273 **Competing Interests**

274 The authors declare they have no competing interests with regard to this work.

275 **Ethical Approval**

276 Ethical Approval was granted by the Science, Technology, Engineering and Mathematics Ethics
277 Approval Committee of the University of Central approval code STEMH 1058. Student consent
278 to the appropriate sharing of anonymised examination results is given in the Learning
279 Agreement signed by each student, on the basis of GDPR Article 6(1)(e): Processing necessary
280 for the performance of a task in the public interest.

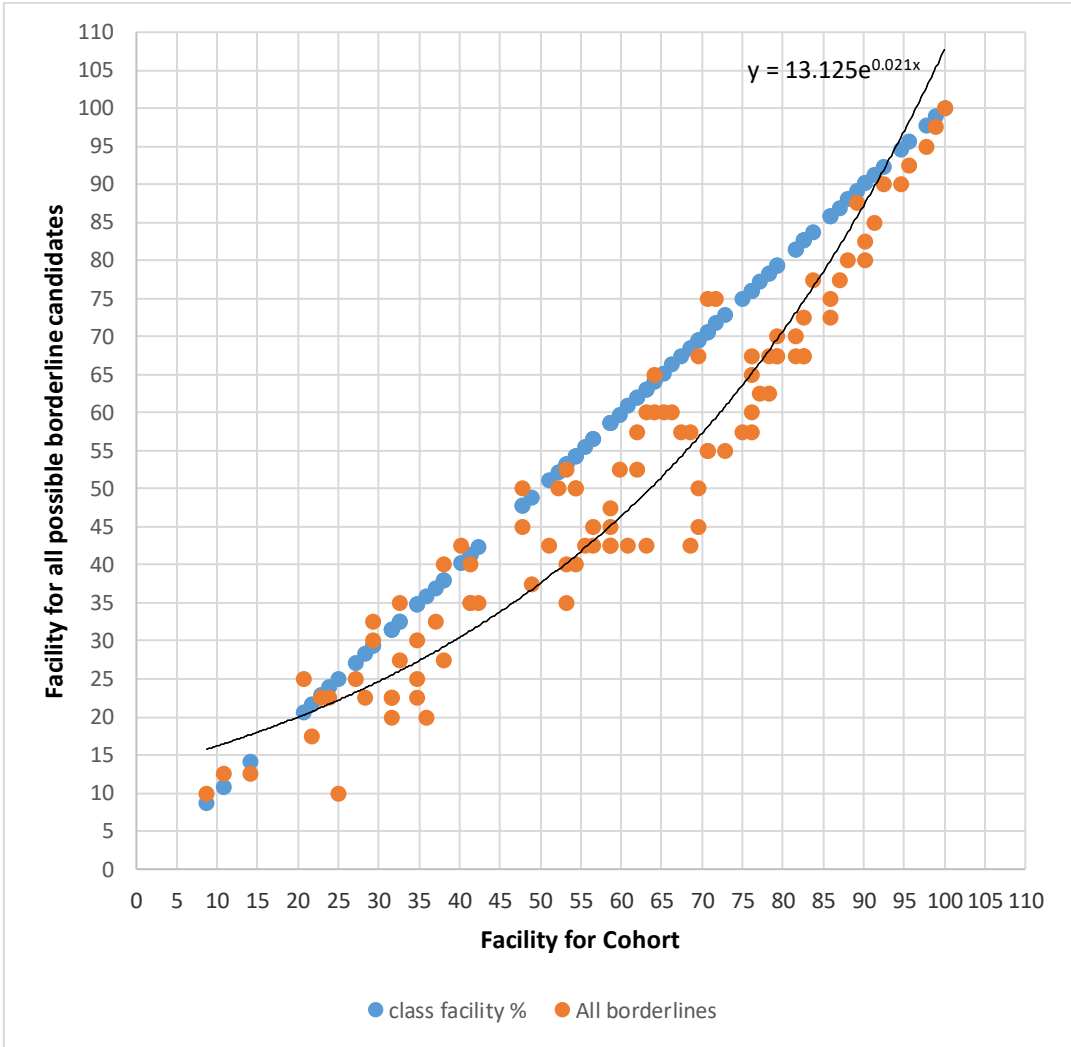
281 **Availability of Data and Materials**

282 The anonymised datasets analysed during the current study are available from the
283 corresponding author on reasonable request.

284 **Authors' Contributions**

285 JMcL led on the conception and design of the work and is the guarantor of the paper. KAR
286 contributed to drafting, revising and critically appraising the content of the paper. BW was
287 responsible for extracting the original data from the records, and contributed to drafting,
288 revising and critically appraising the content of the paper. MS contributed to critically
289 appraising the content of the paper, and running checks on data sets using the methodology.
290 All authors have approved the final version of the paper.

291



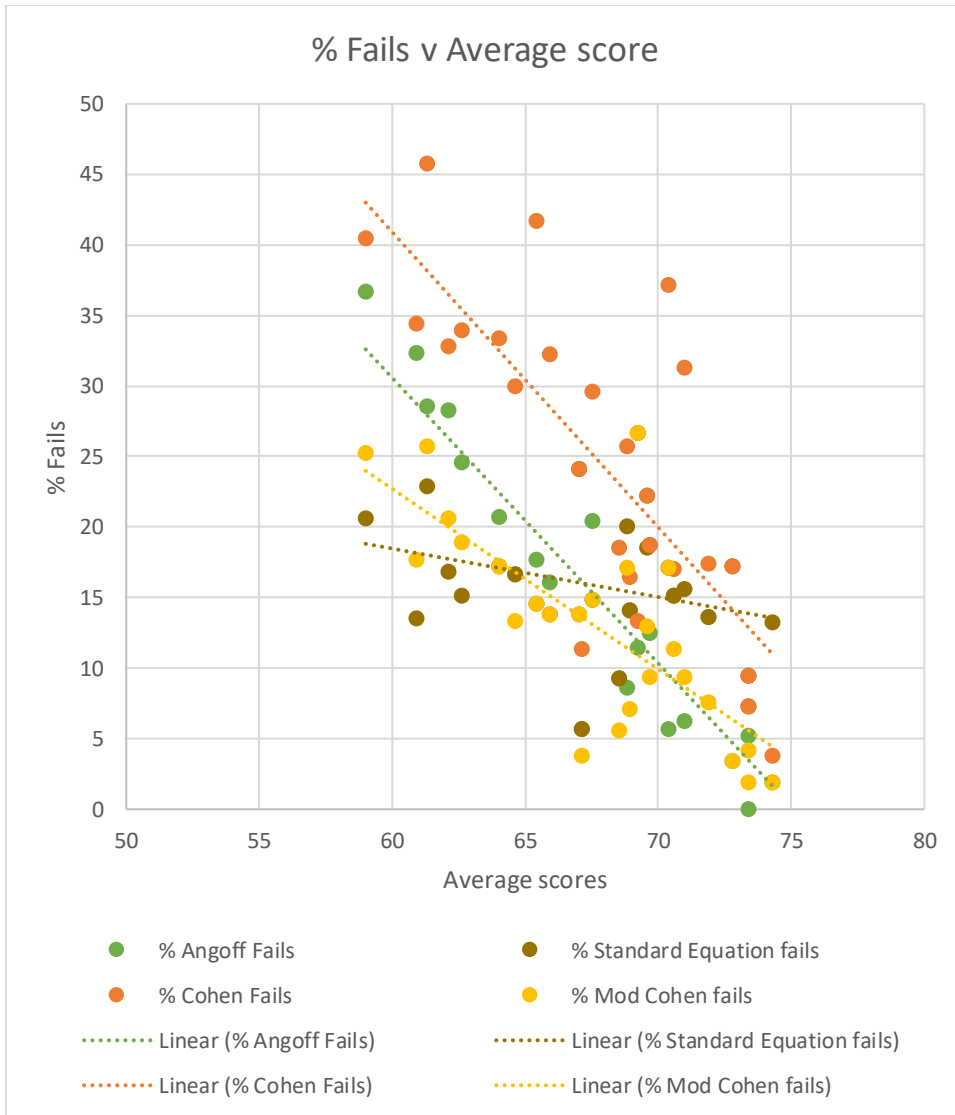
292

293

294 Fig. 1. Facility for all candidates plotted against all possible borderline candidates. As a
 295 reference, cohort Facility is plotted against itself as a 45° slope.

296

297



299

300

301 Fig. 2. The percentage fail rate for each of four different standard setting methods, plotted
 302 against the average score in that exam, for a total of 26 exams. The linear trendline for each
 303 has been added for clarity.

304

305

-
- ¹ Angoff W. H. Scales, norms, and equivalent scores. In: Thorndike RL, ed. *Educational Measurement* (2nd ed.). Washington, DC: American Council on Education. 1971;508-600.
- ² Clauser JC, Hambleton RK, Baldwin P. The effect of rating unfamiliar items on Angoff passing scores. *Educational and psychological measurement* 2017; **77**:901-916.
- ³ De Gruijter DN. Compromise models for establishing examination standards. *Journal of Educational Measurement* 1985; **22**: 263-269.
- ⁴ Kramer A, Muijtjens A, Jansen K, Düsman H, Tan L, Van Der Vleuten C. Comparison of a rational and an empirical standard setting procedure for an OSCE. *Medical Education* 2003; **37**:132-139.
- ⁵ Norcini J, Guille R. Combining tests and setting standards. In: Norman GR, van der Vleuten CPM, Newble D. Eds. *International Handbook of Research in Medical Education*. Springer, Dordrecht. 2002: 811-834.
- ⁶ Cohen-Schotanus J, van der Vleuten CP. A standard setting method with the best performing students as point of reference: practical and affordable. *Medical Teacher* 2010; **32**: 154-160.
- ⁷ Taylor CA. Development of a modified Cohen method of standard setting. *Medical teacher*. 2011; **33**:e678-82.
- ⁸ Sam AH, Field SM, Collares CF, van der Vleuten CP, Wass VJ, Melville C, Harris J, Meeran K. Very-short-answer questions: reliability, discrimination and acceptability. *Medical Education* 2018; **52**: 447-455.
- ⁹ Muijtjens AM, Hoogenboom RJJ, Verwijnen GM, Van der Vleuten CPM. Relative or absolute standards in assessing medical knowledge using progress tests. *Advances in Health Sciences Education* 1998; **3**:81-87.