

# Examiner Variability in Clinical Assessments: Do Examiner Pairings Influence Candidate Ratings?

Aileen Faherty (✉ [aileen.faherty@nuigalway.ie](mailto:aileen.faherty@nuigalway.ie))

National University of Ireland Galway Library <https://orcid.org/0000-0002-7016-0099>

Yvonne Finn

National University of Ireland Galway

Tim Counihan

National University of Ireland Galway

---

## Research article

**Keywords:** Clinical assessments, Reliability, Examiner Variability, Examiner Factors

**Posted Date:** August 20th, 2019

**DOI:** <https://doi.org/10.21203/rs.2.13210/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at BMC Medical Education on May 11th, 2020.  
See the published version at <https://doi.org/10.1186/s12909-020-02009-4>.

# Abstract

**Background** The reliability of clinical assessments is known to vary considerably and inter-examiner variability is a key contributor. This may result in significant differences in scores between comparable candidates, a serious challenge in medical education. An approach frequently adopted to avoid this and improve reliability is to pair examiners and ask them to come to an agreed score. Little is known however, about what occurs when these paired examiners interact to generate a score. **Methods** A fully-crossed design was employed with each participant examiner observing and scoring. A quasi-experimental research design used candidate's observed scores in a mock clinical assessment as the dependent variable. The independent variables were examiner numbers, demographics and personality. Demographic and personality data was collected by questionnaire. A purposeful sample of medical doctors who examine in the Final Medical examination at our institution was recruited. **Results** Variability between scores given by examiner pairs (N=6) was less than the variability with individual examiners (N=12). 75% of examiners (N=9) scored below average for neuroticism and 75% also scored high or very high for extroversion. Two thirds scored high or very high for conscientiousness. The higher an examiner's personality score for extroversion, the lower the amount of change in his/her score when paired up with a co-examiner; reflecting possibly a more dominant role in the process of reaching a consensus score. **Conclusions** While the variability between scores given by examiner pairs (N=6) was less than the variability with individual examiners (N=12), the reliability statistics for both assessments were comparable. Using paired examiners resulted in a more accurate and robust score than simply averaging two independent examiners scores. The higher an examiner's personality score for extroversion, the lower the amount of change in his/her score when paired up with a co-examiner; reflecting possibly a more dominant role in the process of reaching a consensus score. These findings could have implications for the organisation and administration of clinical assessments. Further studies with larger numbers of participants might establish if personality testing before choosing examiner pairs could be utilised to help pair examiners and improve examiner variability.

## Background

To become a competent physician, undergraduate medical students must be assessed not only on factual knowledge but also on communication and clinical skills. The reliability of clinical assessments to test these skills however, is known to be compromised by high levels of variability i.e. different results on repeated testing<sup>1,2</sup>.

Candidate variability, case variability (case specificity) and examiner variability all contribute to the overall variability of a clinical assessment. Candidate variability reflects the difference between candidates and in the absence of other variables (or error) candidate variability represents the true variability. Case specificity refers to the phenomenon that a candidate's performance can vary from one case to the next due to differing levels of difficulty or content<sup>2,3</sup>. Examiner variability refers to the fact that two examiners observing the same performance may award different scores. Many studies have shown that examiner variability is the most significant factor contributing to variability in clinical

examinations<sup>4,5</sup> and may even exceed the variability accounted for by differences in candidates<sup>6</sup>. The degree of examiner variability which is deemed acceptable is generally a minimum of 0.6 with 0.8 being the gold standard (where 0 shows no relationship between two examiners scores and 1 is a perfect agreement)<sup>7</sup>.

Variability in how examiners score candidates may be consistent, for example, an examiner who always marks candidates stringently (often referred to as a hawk) or an examiner who is consistently lenient (a dove)<sup>3</sup>. This kind of consistent examiner behavior can often be adjusted for when analyzing results. However, examiner variability may not always be so consistent and predictable.

Examiners in clinical assessments are subject to many forms of bias<sup>8</sup>. The 'Halo effect' refers to the phenomenon where an examiner's overall first impression of a candidate ("*he seems like he knows his stuff*") leads to failure to discriminate between discrete aspects of performance when awarding scores<sup>9</sup>. In addition, familiarity with candidates, the mood of the examiner, personality factors, and seeing information in advance have all also been found to affect examiners judgments<sup>10,11,12</sup>.

Variability may result in a borderline candidate achieving a score in the pass range in one assessment and the same candidate failing a comparable assessment testing the same/similar competencies. In high stakes examinations, such as medical licensing examinations, this can have serious implications for both the candidate, the medical profession and even society in general. Moreover, pass/fail decisions are now increasingly being challenged<sup>13</sup>.

Efforts to reduce variability in clinical assessments have ranged from utilising higher numbers of stations in Objective Structured Clinical Examinations (OSCEs), to employing objective checklists<sup>2,14</sup>. Many of these approaches have not been found to make any meaningful improvements to reliability<sup>15</sup>. However, increasing the number of observations in an assessment (by involving more examiners in the observation of many performances) *has* been shown to improve reliability<sup>16</sup>. In their evaluation of the mini-clinical exercise used in US medical licensing examinations, Margolis and colleagues stated that having a small number of raters rate an examinee multiple times was not as effective as having a larger number of raters rate the examinee on a smaller number of occasions and more raters enhanced score stability<sup>6</sup>. Consequently, an approach frequently adopted to improve reliability and limit the impact of inter-examiner variability is to pair examiners and ask them to come to an agreed score for a candidate's performance. Little is known however, about what occurs when these paired examiners interact to generate a score.

## Summary of existing literature

Although the hawk-dove effect was described by Osler as far back as 1913<sup>17</sup> its impact on the reliability of clinical examinations was only explored in recent years. In 1974 Fleming et al. described a major revision of the Membership of the Royal College of Physicians (MRCP) UK clinical examination and

identified one examiner as a hawk<sup>18</sup>. There was a significantly lower pass rate in the group of candidates where this examiner examined compared with the remainder (46.3% and 66.0% respectively).

In 2006, an analysis of the reliability of the MRCP UK clinical examination that existed at that time, the Practical Assessment of Clinical Examination Skills (PACeS) exam, found that 12% of the variability in this examination was due to the hawk-dove effect<sup>19</sup>. Examiners were more variable than stations.

In 2008 Harasym et al.<sup>20</sup> found an even greater effect due to the hawk-dove phenomenon in an OSCE evaluating communication skills. Forty four percent of the variability in scores was due to differences in examiner stringency/leniency; over four times the variance due to student ability (10.3%).

As mentioned above, many types of rater-bias are known to be at play when human judgement comprises part of any assessment process (halo effect, the mood of the rater, familiarity with candidates, personality factors etc<sup>8,9,10,11</sup>). Yeates and colleagues in 2013 proposed three themes to explain how examiner-variability arises<sup>21</sup>. They termed these: differential salience (what was important to one examiner differed to another); criterion uncertainty (assessors' conceptions of what equated to competence differed and were uncertain); information integration (assessors tend to judge in their own unique descriptive language forming global impressions rather than discrete numeric scores).

Govaerts suggests that some examiner-variability may simply arise from individual examiners' peculiarities in approach and idiosyncratic judgements made as a result, of the interaction between social and cognitive factors<sup>12</sup>.

Earlier reports had suggested that employing objective checklists would help overcome examiner-variability by regulating subjectivity<sup>2</sup>. More recently however, several lines of evidence suggest that global judgements produce more reliable results than highly structured tools<sup>4, 14</sup>. Furthermore, measurement instruments have been shown to account for less than 8% of the variance in performance ratings<sup>22</sup>.

Other proposals to improve reliability have involved increasing the number of items used per station. However, Wilkinson et al analysed examiners marks over a four-year period in New Zealand and found that while items-per-station increased over the four years, there was no correlation between items-per-station and the station inter-rater reliability<sup>4</sup>.

The impact of examiner training has also been looked at in many studies<sup>23</sup>. Cooke et al.<sup>24</sup> found no significant effect and while Holmboe et al.<sup>25</sup> showed that training produced an increase in examiner stringency, this increase was inconsistent.

In a recent literature review on rater cognition in competency based education Gauthier et al.<sup>15</sup> summarised the situation stating: "*attempts to address this variability problem by improving rating forms and systems, or by training raters, have not produced meaningful improvements*".

## **Aims:**

1. To explore the difference in candidate markings when using individual versus paired examiners
2. To analyse how, an examiners marks vary from when s/he examines alone to when s/he examines as one member of a pair of examiners
3. To explore if there is a correlation between examiner personality factors and examiner behaviour in marking candidates' performances.

## **Methods**

### **Design**

A fully-crossed design was employed with each participant examiner observing and scoring recordings of candidates' performances. A quasi-experimental research design was used. The dependent variable was candidate's observed scores in a mock clinical assessment. The independent variables were examiner number (single or paired), examiner demographics and examiner personality. There was no control group; examiner participants served as their own control i.e. control was exercised through more than one observation of the same phenomenon<sup>26</sup>.

### **Setting and characteristics of participants**

The study population consisted of qualified medical doctors who examine in the final medical short-case examination at our institution. Participants were invited by email and each received a participant information leaflet, electronic consent form and demographic questionnaire.

### **Description of all processes, interventions and comparisons**

In the final medical examination at our school, medicine and surgery are assessed together in a short-case examination. Each candidate is assessed over 6 short-cases, a mixture of medical and surgical cases, each lasting 6 minutes using a real or simulated patient. Candidates are observed by pairs of examiners, usually a surgeon paired with a physician. After each candidates' performance, examiners discuss and come to an agreed score using a domain based marking sheet. Our data collection exercise was set up to mimic as closely as possible this real-world examination scenario using recordings of simulated patients.

Participants were stratified to mimic the examiner pairings usually employed (a surgeon with a physician). The participants did not assess a real students' performance; instead we used video recordings of standardised student performances (using actors) that were previously created for the purposes of examiner training. We selected 3 videos as follows: one example each of a weak, average and good performance. Examiners were not aware what level of performance they would be watching.

Different case types were selected (one medical, one surgical and one general medical/surgical) to avoid one examiner being more familiar than the other examiners with the content of the selected cases. Each participant viewed, initially on their own individual screens, the three recordings and graded them independently. The total possible score at each station was 50 marks—with ten marks each allocated to five separate domains; attitude and professionalism, communication skills, clinical skills, knowledge and lastly management. Our schools OSCE Management Information System Software—OMIS by Qpercom Ltd was used to enter marks. Utilising this software examiners were blinded to their individual scores of a given performance. When the examiners scored the performance across the individual five domains, the scores were on a slider and the examiner did not see what their resultant overall mark was from combining the 5 domains.

After the examiners had scored the videos independently there was a break for refreshments. Examiners then completed a validated 60 item personality questionnaire - the NEO Five Factor Index (NEO-FFI)<sup>27</sup>. In this personality index, no single cut-off point separates those who “have” a particular personality-trait from those who do not, rather individual scores represent degrees of each of the five main personality traits—neuroticism, extroversion, openness to experience, agreeableness and conscientiousness. Score results are usually expressed as a T score and can be further described as being very low, low, average, high and very high for each of the domains. After completing the personality questionnaire, examiners were moved to a neutral location and paired up with another examiner to review and discuss the same three performances again and this time devise a joint mark which was entered on OMIS. The order of the videos when watched as individual examiners compared with observing in pairs was counterbalanced to control for an order effect<sup>28</sup>. Blinding the participant as to the overall original scores given and changing the order of videos from the previous observation was particularly important to maintain internal validity. We looked for a correlation between the total amount of change in an examiners marks from when they examined individually to when they examined in a pair, and their personality scores (See figure 1 for scatter plots).

## Statistical analysis

Data collected on candidate scores was analysed using the OMIS OSCE management software and SPSS 24 (IBM corp). Preliminary analyses confirmed that the data were not normally distributed and, therefore, non-parametric methods were employed in the statistical analysis. Descriptive statistics were generated using tables and charts. The OMIS OSCE management software allowed for psychometric analysis and provided support for generalisability analysis.

## Results

Fifty potential participants were contacted by email and invited to participate. Seventeen respondents accepted the invitation and twelve completed the study - 10 male and 2 female. They had an average of 13.6 years' experience examining in the final-medical short-case examination at our institution. Two

thirds were in posts that were combined clinical and academic. Two participants held any formal qualifications in medical education.

## Variability

Table 1 and figure 2 show the overall scores awarded by each examiner to the three candidates when examining alone and demonstrate considerable variability in examiners' scores.

Table 2 and figure 3 show the overall scores awarded by Examiners *when in pairs* to the three candidates. The ranges and standard deviations reveal that the variability between scores given by examiner pairs is, as might be expected, less than that in the assessment using 12 individual examiners.

Generalisability analysis allows for more in-depth analysis of the variance of our assessments, identifying the relative contribution of each of the components (or facets) of that assessment—the examiners (observations, O), the scenarios (S) and their interactions (SO). In the assessment using individual examiners, 87.1 % of variance was found to be due to examiners while 12.9% was due to the interaction between the examiner and the scenario (table 3).

## Reliability

Using Classical Test Theory Cronbach's alpha and intra-class correlation coefficients were calculated for the assessment using 12 single examiners and the second assessment using 6 examiner pairs. The reliability statistics for the two assessments were in fact comparable (table 4).

Using Generalisability theory, the G-coefficient of the assessment using 12 individual examiners was calculated as 0.95. The Standard Error of Measurement (SEM) was 4.5% (see table 5) which means the candidates' true score lies between the observed score  $\pm 4.5\%$ . This is quite a high margin which would have significant consequences for marks around the pass/fail and honours/pass thresholds. However, our Decision-study (D-study) - which gives us an indication of what happens to the reliability and SEM of an assessment if we increase the number of scenarios, showed that increasing the number of scenarios from 3 to 12 would reduce the SEM to a more acceptable level of 2% (see table 6).

## Impact of pairing up on Candidates' score/outcome

We compared candidates scores when they were examined by 12 individual examiners with their scores when they were examined by 6 examiner pairs (see tables 1 and 2). The 'good' performance was awarded an honour by all 12 individual examiners and all 6 examiner pairs. Similarly, the weak performance was failed by all examiners. However, when examined by individual examiners, the average performance was awarded 4 passes, 6 borderline results (between 40 and 49%) and failed by 2 examiners. When assessed by examiner pairs the average performance was not failed on any occasion but received 4 borderline

marks and 2 passes. Wilcoxon signed rank test showed a statistically significant difference between mean scores for the average student ( $p = 0.0430$ ).

## How each examiners' marks changed when they were paired up

The marks given by each examiner when they examined singly were compared with the agreed mark given by the same examiner to each candidate when examining *in a pair*. The amount of change in each examiner's overall mark for the three candidates was calculated. Table 7 shows the change in examiners marks and the direction of that change (a minus sign indicated their mark reduced when they paired up). The amount of change (regardless of whether positive or negative) for each examiner was then summed across the three candidates to devise a figure representing the total amount of change in marks per examiner.

There was a statistically significant negative correlation ( $-0.808$ ) between extroversion and change in examiners score - the higher an examiners' score for extroversion the lower the degree of change in his or her score when paired up with a co-examiner ( $p = 0.001$ ) (see table 8).

## Discussion

Our study shows that there is less examiner variability and therefore improved reliability in using examiner pairs. Using paired examiners resulted in a more accurate and robust score than simply averaging two independent examiners scores. The average performance was passed by all examiner pairs however two examiners failed this candidate when examining individually ( $p = 0.0430$ ). This has implications for candidate outcomes. The correlation between degree of change of examiners mark and score for extroversion suggests personality traits do have an impact on examiner behaviour and candidate outcomes.

Comparing the marks given by examiners in pairs to the marks they previously gave when examining alone proved revealing (see table 7). In no instance did the new mark simply equate to the mean of, or midpoint between the two individual examiners marks. Instead, in each case, the marks awarded by examiner pairs tended towards one examiner's previous original mark rather than the other, the 'dominant' examiner if you will. In 5/6 pairs this 'dominant' examiner was a physician. All of the physicians scored high or very high for extroversion and we found a statistically significant correlation between change in examiner score and extroversion - the higher an examiners score for extroversion the lower the amount of change in his or her score when paired up ( $p = 0.001$ ). This is perhaps not surprising as extroverts are described as assertive and talkative, two characteristics which would certainly enable an examiner to "stand their ground" as it were.

Our sample confirmed the findings of previous studies that in personality testing, doctors tend to score low for neuroticism and high for extraversion<sup>29</sup>. We did not find any relationship between examiner personality and stringency as was found in a previous study in our school<sup>17</sup>.



Our findings support the opinion that the score of examiner pairs may be a more accurate and robust score than simply averaging two independent examiners scores. This could have implications for the organisation and administration of clinical assessments. Further study with a larger number of participants might establish if personality testing before choosing examiner pairs is warranted.

## **Limitations:**

Recruitment of participants proved difficult and so our sample was small. There was a small number of female participants. It could be argued that there was a learning or testing effect in the set-up of our mock examination whereby the examiners assessed the same performances twice. Ideally, we would have used a larger number of video recordings to avoid compromising the internal validity of this study in this way however, increasing the length of the process would have made recruitment even more difficult. Some investigators raised concerns about the recording of participants' discussions giving rise to "the Hawthorne effect" where the awareness of being observed impacts on research participants' behaviour<sup>30</sup> however, a review of the literature found very little empirical support for this effect in medical education<sup>31</sup>.

## **Conclusions**

Our study shows that the practice of using paired examiners in clinical assessments is to be recommended. While using paired examiners may use more resources, in the case of high stakes assessments and an increasingly litigious society, grades are awarded by examiner pairs after robust discussion and therefore can be more easily defended in the case of appeals.

## **List Of Abbreviations**

CTT: Classical Test Theory

MRCP: Membership of the Royal College of Physicians

NEO-FFI: Neuroticism—Extroversion-Openness to experience Five Factor Index

OMIS: OSCE Management Information System Software (by Qpercom Ltd)

OSCE: Objective Structured Clinical Examinations

PACE: Practical Assessment of Clinical Examination Skills

SEM: Standard Error of Measurement

## **Declarations**

## **Ethics approval and consent to participate**

Ethical Approval was sought from and granted by the College of Medicine, Nursing and Health Sciences research ethics committee at the National University of Ireland Galway.

## **Consent for publication**

Not applicable

## **Availability of Data and Material**

The datasets used and/or analysed during the current study are available from the corresponding author on reasonable request

## **Competing interests**

The authors declare that they have no competing interests

## **Funding**

This research was funded wholly by the corresponding author

## **Authors contributions**

YF and TC conceived the original idea for this research. AF was principal investigator. AF and TC were involved in recruitment of participants. AF acquired and analysed the data. YF and AF interpreted the data. AF drafted the manuscript and all authors revised the manuscript and approved the final version for publication..

## **Acknowledgements**

The authors would like to thank Professor Peter Cantillon and Dr Paul O'Connor for their guidance and feedback during the carrying out of this research and Dr Thomas Kropmans for his assistance with data analysis.

## **Authors information**

AF is a General Practitioner and Lecturer in Clinical Practice with the Discipline of General Practice, National University of Ireland Galway.

YF is a Lecturer Above the Bar with the School of Medicine, National University of Ireland, Galway

## References

1. Downing S. Reliability: on the reproducibility of assessment data. *Medical education*. 2004;38(9):1006–1012.
2. Crossley J, Davies H, Humphris G, Jolly B. Generalisability: a key to unlock professional assessment. *Medical education*. 2002;36(10):972–978.
3. Crossley J, Russell J, Jolly B, Ricketts C, Roberts C, Schuwirth L, Norcini J. 'I'm pickin'up good regressions': the governance of generalisability analyses. *Medical education* 2007;41(10):926–934.
4. Wilkinson T, Frampton C, Thompson-Fawcett M, Egan T. Objectivity in objective structured clinical examinations. *Academic Medicine*. 2003;78(2):219–223.
5. McGill DA, Van der Vleuten CP, Clarke MJ. Supervisor assessment of clinical and professional competence of medical trainees: a reliability study using workplace data and a focused analytical literature review. *Advances in health sciences education*. 2011 Aug 1;16(3):405–25.
6. Margolis MJ, Clauser BE, Cuddy MM, Ciccone A, Mee J, Harik P, Hawkins RE. Use of the mini-clinical evaluation exercise to rate examinee performance on a multiple-station clinical skills examination: a validity study. *Academic Medicine*. 2006 Oct 1;81(10):S56–60.
7. Rushforth HE. Objective structured clinical examination (OSCE): review of literature and implications for nursing education. *Nurse education today*. 2007 Jul 1;27(5):481–90.
8. Saal FE, Downey RG, Lahey MA. Rating the ratings: Assessing the psychometric quality of rating data. *Psychological bulletin*. 1980 Sep;88(2):413.
9. Wood TJ. Exploring the role of first impressions in rater-based assessments. *Advances in Health Sciences Education*. 2014 Aug 1;19(3):409–27.
10. Stroud L, Herold J, Tomlinson G, Cavalcanti RB. Who you know or what you know? Effect of examiner familiarity with residents on OSCE scores. *Academic Medicine*. 2011 Oct 1;86(10):S8–11.
11. Gingerich A, Regehr G, Eva KW. Rater-based assessments as social judgments: rethinking the etiology of rater errors. *Academic Medicine*. 2011 Oct 1;86(10):S1–7.
12. Govaerts MJ, Van der Vleuten CP, Schuwirth LW, Muijtjens AM. Broadening perspectives on clinical performance assessment: rethinking the nature of in-training assessment. *Advances in health sciences education*. 2007 May 1;12(2):239–60.
13. Tweed M, Miola J. Legal vulnerability of assessment tools. *Medical Teacher*. 2001 Jan 1;23(3):312–4.
14. Regehr G, MacRae H, Reznick R, Szalay D. Comparing the Psychometric Properties of Checklists and Global Rating Scales for Assessing Performance on an GSCE-format Examination. *Acad. Med*. 1998;73:993–7.
15. Gauthier G, St-Onge C, Tavares W. Rater cognition: review and integration of research findings. *Medical education*. 2016 May;50(5):511–22.

16. Tavares W, Ginsburg S, Eva KW. Selecting and simplifying: rater performance and behavior when considering multiple competencies. *Teaching and learning in medicine*. 2016 Jan 2;28(1):41–51.
17. Finn Y, Cantillon P, Flaherty G. Exploration of a possible relationship between examiner stringency and personality factors in clinical assessments: a pilot study. *BMC medical education*. 2014 Dec;14(1):1052.
18. Fleming PR, Manderson WG, Matthews MB, Sanderson PH, Stokes JF. Evolution of an examination: MRCP (UK). *Br Med J*. 1974 Apr 13;2(5910):99–107.
19. McManus IC, Thompson M, Mollon J. Assessment of examiner leniency and stringency ('hawk-dove effect') in the MRCP (UK) clinical examination (PACES) using multi-facet Rasch modelling. *BMC Medical Education*. 2006 Dec;6(1):42.
20. Harasym PH, Woloschuk W, Cuning L. Undesired variance due to examiner stringency/leniency effect in communication skill scores assessed in OSCEs. *Advances in Health Sciences Education*. 2008 Dec 1;13(5):617–32.
21. Yeates P, O'Neill P, Mann K, Eva K. Seeing the same thing differently. *Advances in Health Sciences Education*. 2013 Aug 1;18(3):325–41
22. Landy FJ, Farr JL. Performance rating. *Psychological bulletin*. 1980 Jan;87(1):72.
23. Williams RG, Klamen DA, McGaghie WC. Cognitive, social and environmental sources of bias in clinical performance ratings. *Teaching and learning in medicine*. 2003 Oct 1;15(4):270–92.
24. Cook DA, Dupras DM, Beckman TJ, Thomas KG, Pankratz VS. Effect of rater training on reliability and accuracy of mini-CEX scores: a randomized, controlled trial. *Journal of general internal medicine*. 2009 Jan 1;24(1):74.
25. Holmboe ES, Hawkins RE, Huot SJ. Effects of training in direct observation of medical residents' clinical competence: a randomized trial. *Annals of internal medicine*. 2004 Jun 1;140(11):874–81.
26. DePoy E, Gitlin LN. Introduction to research-e-book: Understanding and applying multiple strategies. Elsevier Health Sciences; 2015 Mar 6.
27. Costa PT, McCrae RR. The revised neo personality inventory (neo-pi-r). *The SAGE handbook of personality theory and assessment*. 2008 Jun 24;2(2):179–98.
28. Cohen L, Manion L, Morrison K. Action research. In *Research methods in education* 2013 Mar 7 (pp. 368–385). Routledge.
29. Borges NJ, Savickas ML. Personality and medical specialty choice: a literature review and integration. *Journal of Career Assessment*. 2002 Aug;10(3):362–80.
30. Chiesa M, Hobbs S. Making sense of social research: How useful is the Hawthorne Effect?. *European Journal of Social Psychology*. 2008 Jan;38(1):67–74.
31. Paradis E, Sutkin G. Beyond a good story: from Hawthorne Effect to reactivity in health professions education research. *Medical education*. 2017 Jan;51(1):31–9.

## Tables

Table 1: Individual Examiners Overall Scores for each Candidate

Examiner Number	Good Candidate Overall Score	Average Candidate Overall Score	Weak Candidate Overall Score
1	64.00	44.00	34.00
3	74.00	50.00	36.00
5	64.00	38.00	20.00
6	64.00	44.00	24.00
7	68.00	42.00	34.00
9	80.00	44.00	28.00
10	80.00	34.00	28.00
11	82.00	48.00	26.00
12	70.00	56.00	40.00
14	94.00	58.00	12.00
16	90.00	48.00	30.00
17	86.00	50.00	34.00
Candidates' Mean	76.33 (10.54)	46.33 (6.86)	28.83 (7.69)
Range	30	24	28

Table 2: Paired Examiners Overall Scores for each Candidate

Examiner Pair	Good Candidate	Average Candidate	Weak Candidate
Examiners 1+5	64.00	48.00	26.00
Examiners 3+11	78.00	46.00	24.00
Examiners 6+14	82.00	56.00	18.00
Examiners 7+12	64.00	52.00	34.00
Examiners 9+16	88.00	48.00	28.00
Examiners 10+17	80.00	44.00	30.00
Candidates' Mean	76 (9.87)	49 (4.33)	34 (5.46)
Range	24	12	16

Table 3: Analysis of Variance of the main facets of assessment using 12 examiners (Observations=O, Scenarios=S and their interactions =SO)

Source	Components							
	SS	df	MS	Random	Mixed	Corrected	%	SE
O	1.4731	11	0.13392	0.04254	0.04254	0.04254	87.1	0.01752
S	0.0126	2	0.00630	0.00000	0.00000	0.00000	0.0	0.00040
OS	0.13860	22	0.00630	0.00630	0.00630	0.00630	12.9	0.00182
Total	1.62430	35					100	

Table 4: Reliability statistics for the assessments using single and paired examiners

	Cronbach's Alpha	Intraclass Correlation Co-efficient							
		Intraclass Correlation		95% Confidence Interval		F Test with True Value 0			
				Lower Bound	Upper Bound	Value	df1	df2	Sig
Single Examiners	0.99	Single Measures	0.887	.648	.997	98.97	2	22	.000
		Average Measures	0.990	.957	1.00	98.97	2	22	.000
Paired Examiners	0.983	Single Measures	0.925	.700	.998	60.533	2	10	.000
		Average Measures	0.987	.933	1.00	60.533	2	10	.000

Table 5: G-study table for assessment using 12 individual examiners

Sources of var.	Differ. Variance	Sources of var.	Relative err. Var.	% rel.	Absolute err. var.	% abs
O	0.04254		...	...	...	...
	...	S	...	...	0.000	0.0
	...	OS	0.0021	100	0.00210	100
Sum of variances	0.04254		0.00210	100%	0.00210	100%
Standard Deviation	0.20625		0.04583 (Relative SE)		0.04583 (Absolute SE)	

Table 6: D study analysis of assessment using 12 individual examiners

Opt 1	Opt 2	Opt 3	Opt 4	
Lev.Univ.	Lev.Univ.	Lev.Univ.	Lev.Univ.	Lev.Univ.
O	12 INF	12 INF	12 INF	12 INF
S	3 INF	6 INF	9 INF	15 INF

Observations	36	72	108	144	180
Coef_G rel.	0.953	0.976	0.984	0.988	0.990
rounded	0.95	0.98	0.98	0.99	0.99
Coef_G abs.	0.953	0.976	0.984	0.988	0.990
rounded	0.95	0.98	0.98	0.99	0.99
Rel. Err. Var.	0.00210	0.00105	0.00070	0.00052	0.00042
Rel.Std. Err.M.	0.04583	0.03240	0.02646	0.02291	0.02049
Abs. Err. Var.	0.00210	0.00105	0.00070	0.00053	0.00042
Abs.Std. Err.M.	0.04583	0.03240	0.02646	0.02291	0.02049

**Table 7: Changes in examiners marks when they moved from examining alone to examining in a pair.**

Examiners	Pair A		Pair B		Pair C		Pair D		Pair E		Pair F	
	1	5	3	11	7	12	6	14	9	16	10	17
Honours	0.0	0.0	4	-4	-4	-6	18	-12	8	-2	0	-6
Pass	4.0	10	-4	-2	10	-4	12	-2	4	0	10	-6
Fail	-8.0	6.0	-12	-2	0	-6	-6	6	0	-2	2	-4
Total	12	16	20	8	14	16	36	20	12	4	12	16

**Table 8: Relationship between the amount of change in examiners scores and personality**

	Spearman's Correlation co-efficient rho	P value
Neuroticism	0.352	0.262
Extraversion	-0.808**	0.001
Openness to Experience	-0.185	0.565
Agreeableness	-0.501	0.097
Conscientiousness	-0.451	0.141

**\*\***. Correlation is significant at the 0.01 level (2 tailed).

# Figures

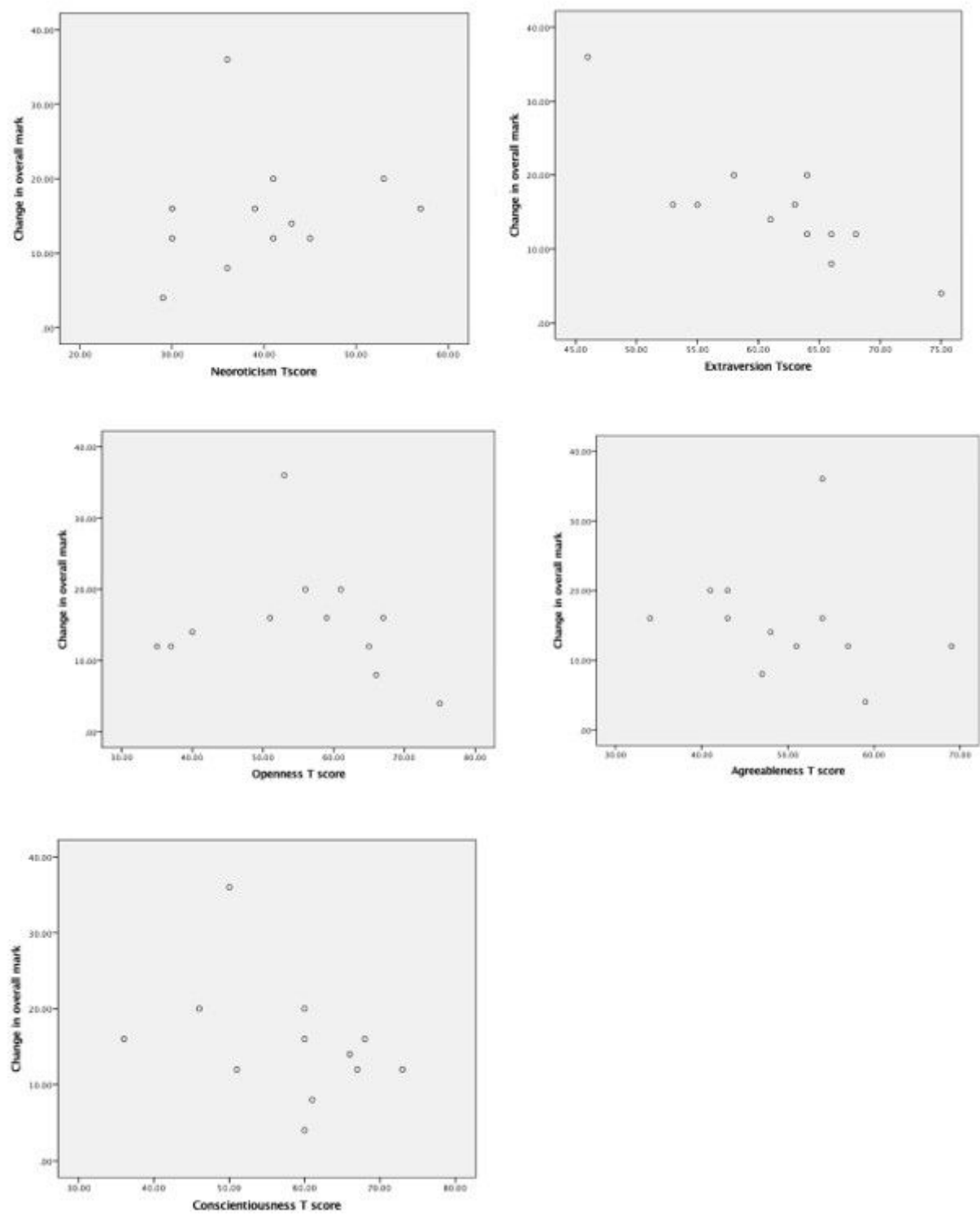
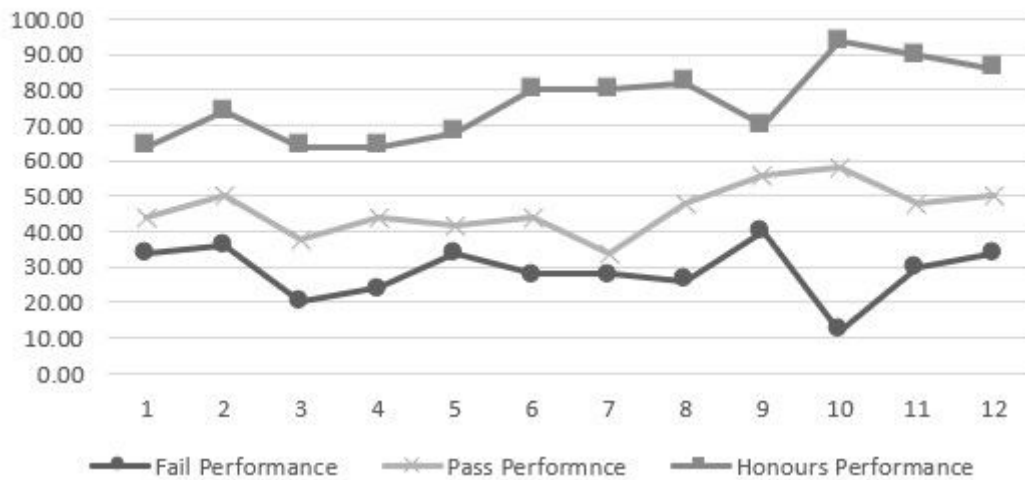


Figure 1

Scatter plots for correlation between change in examiners marks and personality



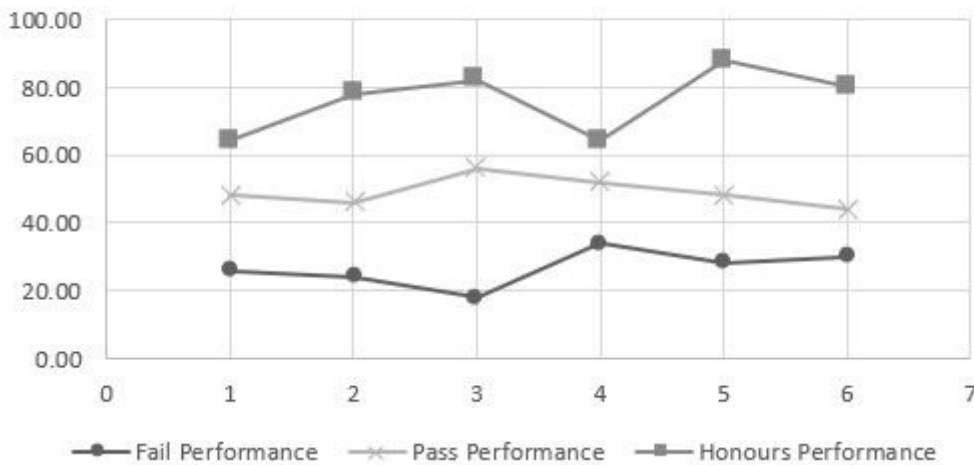
### Variability of Overall Scores - Individual Examiners



**Figure 2**

Variability of overall scores – Individual Examiners

### Variability of Overall Scores - Examiner Pairs



**Figure 3**

Variability of overall scores – Paired Examiners