

1 published studies reporting test-retest designs as part of tool development. Because of its
2 predominance and importance in preceding other essential stages in tool development (e.g.,
3 scale validation, development of reference/cut-off scores, etc.), the present paper will focus
4 on test-retest reliability, and the statistical problems still observed in this area. While other
5 discussions exist in relation to test-retest methodology, such as the choice of appropriate
6 retest time frames (Chmielewski & Watson, 2009), the focus of this paper will be the analysis
7 of data and presentation of results, rather than study design or data collection. Furthermore,
8 the paper will focus specifically on test-retest of numeric outcomes since these are common
9 outcomes in psychological practice (e.g., test scores, rating scales, etc.) and hence, are often
10 the focus of analysis difficulties. However, some of the key principles also apply to
11 categorical and diagnostic measures.

12 Many of the topics that will be discussed in the current paper were initially examined
13 by Altman & Bland (1983) within medical research, and issues around types and assessments
14 of reliability have been discussed by numerous authors since that time (e.g., Baumgartner,
15 2000; Bedard, Martin, Krueger, & Brazil, 2000; Ludbrook, 2002; Streiner, Norman, &
16 Cairney, 2014; Weir, 2005). However, practical changes have been slow to translate into
17 many research domains, including psychology, and this was the rationale for the current
18 paper. The aim of this paper is not to present a comprehensive review of all statistical
19 methods for assessing measurement reliability, or even test-retest reliability specifically, but
20 to discuss some of the fundamental aspects of test-retest reliability analysis that may not be
21 well-understood by researchers undertaking this type of study. The paper will summarise
22 some of the common errors that continue to be observed in published studies, and offer an
23 introduction to relatively simple methods for assessing and reporting test-retest analyses that
24 avoid such errors.

25 **What is test-retest reliability?**

1 Whilst there are many different meanings ascribed to the term ‘reliability’ across
2 scientific disciplines, ‘test-retest’ reliability refers to the systematic examination of
3 consistency, reproducibility, and agreement among two or more measurements of the same
4 individual, using the same tool, under the same conditions (i.e., when we don’t expect the
5 individual being measured to have changed on the given outcome). Test-retest studies help us
6 to understand how dependable our measurement tools are likely to be if they are put into
7 wider use in research and/or clinical practice. When a measurement tool is used on a single
8 occasion, we want to know that it will provide an accurate representation of the patient or
9 participant so that the outcome may be used for practical purposes (e.g., diagnostics,
10 differentiation of individuals or groups). When a measurement tool is used on multiple
11 occasions (e.g., to compare baseline and follow-up) we want to know that the tool will give
12 accurate results on all occasions, so that observed changes in outcome can be attributed to
13 genuine change in the individual, rather than instability in the measurement tool; this is
14 particularly relevant when assessing the efficacy of treatments and interventions. Finally,
15 when a measurement tool is used to assess different groups (e.g., patients receiving different
16 treatments, different characteristics), we want to know that the tool is accurately measuring
17 all individuals so that any group differences may be considered genuine and not an artifact of
18 measurement. Although demonstrating validity is the key to knowing that the right *thing* is
19 being assessed with any given tool, assessing validity is only truly possible once it has been
20 established that a tool is measuring *something* in the same way each time it is used.

21 In the context of test reliability studies, there are two approaches to understanding the
22 comparability/reliability of test scores – we’ll refer to them in this paper as ‘relative
23 consistency’ and ‘agreement’ – that hold very different definitions of what it means for
24 measurements to be ‘reliable’. Relative consistency, also termed ‘rank-order stability’
25 (Chmielewski & Watson, 2009), means that the relative position or rank of an individual

1 within a sample is consistent across raters/times, but systematic differences in the raw scores
2 given to individuals by different raters or at different times are unimportant. For example, one
3 assessor may score the first three people in a sample as 100, 105, and 107 for IQ, and the
4 second may score the same three people, at the same time, as 105, 110, and 112. Even though
5 the raw scores given by the two raters are not the same, the difference in rating is consistent
6 across all three participants and they maintain the same rank relative to one another;
7 therefore, the IQ measure would be considered to have relative reliability across raters. In
8 contrast, agreement is concerned with the extent to which the raw observed scores obtained
9 by the measurement tool match (or, agree) between raters or time-points, when measuring the
10 same individual in the absence of any actual change in the outcome being measured.

11 If the relative ordering of individuals within a given sample is of greater importance
12 or use than the observed differences between individuals (e.g., finishing position in a race)
13 then assessing the relative consistency between measurements may be suitable. However, this
14 is not typically the case when assessing the test-retest reliability of standardized measuring
15 tools such as psychometric questionnaires. In this case, the aim is to try and make objective
16 measurements that are unaffected by the time or place of measurement, or by attributes of the
17 individual making the measurement. Once the tool is applied in practice, we want to be
18 confident that any given measurement is accurate, and that any differences in outcome
19 observed within a study or clinical practice, are due to real changes in an individual, or
20 genuine differences between individuals/groups. Therefore, the purpose of reliability studies
21 in these contexts is to determine the extent to which repeated measurements agree (i.e., are
22 the same), over time, rater, or context (i.e., test-retest), when used to assess the same
23 unchanged individual. In such a case, it is necessary to assess absolute differences in scores,
24 since these provide a direct measure of score stability at an individual level. Aside from the
25 mere presence/absence of stability, absolute score differences also permit the assessment of

1 additional estimates relevant to test-retest reliability, such as the size and homogeneity of
2 differences across sample ranges. References to test-retest reliability are prevalent across
3 psychological questionnaire-based studies, thus acknowledging the perceived importance of
4 accuracy and repeatability in these tools. However, the methods reported to assess ‘test-
5 retest’ often neglect absolute score differences, and so in many cases they are unsuitable for
6 quantifying the intended form of reliability, and make limit use of the data.

7 **Problems with current methods**

8 Problems in analyses of test-retest reliability most often relate to either unsuitable
9 choice of analysis methods or insufficient reporting of methods. This means, in the first
10 instance, that potentially invalid results are obtained and published (and perhaps trusted in
11 wider practice) and in the second instance, that appraisal and accurate replication of methods
12 and results is precluded. The use of unsuitable statistical methods to assess test-retest
13 reliability may arise from a lack of understanding around different types of reliability,
14 inaccurate understanding of statistical tests/techniques and the results that they produce, or,
15 more pessimistically, as a means for arriving as a desired result where more appropriate or
16 conservative approaches may not (Streiner, 2007). Replication of methodology from
17 published research and resources, some of which may be outdated by the time of use or of
18 poorer quality, can further perpetuate less-than-optimal statistical choices. Inferential
19 statistics are frequently misused in this context and often supersede direct examination and
20 interpretation of the observed differences between measurements. It is very common to see
21 test-retest reliability assessed using bivariate correlation, and non-significant inferential tests
22 of difference, such as paired t-tests, used as evidence of similarity between measurements;
23 neither of which are able to quantify the equality/similarity of repeated scores (Bland &
24 Altman, 1986; Hole, 2014). The following sections summarize the features of these methods
25 that make them unsuitable for assessing agreement-based reliability such as test-retest.

1 **Correlation**

2 **Correlation is not agreement.** A common misconception is that high correlation
3 between two measurements equates to agreement between them. In reality, quite incongruent
4 paired measurements can produce strong and highly statistically significant correlation
5 coefficients, despite the observed agreement between these measurements being very poor.
6 Parametric correlation coefficients (Pearson’s product moment correlations), which are
7 frequently presented in reliability studies, use a -1 to 1 coefficient to quantify how
8 consistently one variable increases (or decreases) relative to another variable increasing,
9 according to how close points lie to any straight line. This can be seen by plotting the
10 measurements against one another and adding a line of best fit. In contrast, agreement in
11 scores means that the 2+ results produced for each individual are the same. To illustrate
12 agreement on a scatter plot the points must lie, not on any straight line, but on the line of
13 equality specifically, where the intercept is 0 and the slope of the line is 1 (Streiner et al.,
14 2014). The difference between correlation and agreement is demonstrated in the data in table
15 1 taken from a laboratory study of adult food preference. This data shows the ratings given by
16 participants when presented with the same food on two occasions. Despite relative stability of
17 food preferences in adulthood, we see that, even relative to the measurement scale, there are
18 large differences (range 15.7 to 226.7) between ratings given on the two occasions.

19

20

Table 1

21

22 The scatterplot in figure 1 illustrates just how far away paired ratings are from
23 agreement, since very few points lie on or close to the dashed line of equality. Despite this
24 clear disparity in ratings, highlighted in both the plot and the absolute score differences, the
25 Pearson’s correlation coefficient for this data is 0.93 ($p < 0.001$), which would undoubtedly be

1 reported as a very strong association.

2

3

Figure 1

4

5 **Correlation conceals systematic bias.** Correlation coefficients are standardized
6 statistics that always fall between -1 (perfect negative association) and 1 (perfect positive
7 association). The units and magnitude of the variables being compared are irrelevant in the
8 calculation of the coefficient, and coefficients are not sensitive to mean differences or
9 changes in scores; as such, coefficients will mask systematic biases (the amount that one
10 measurement differs from another) between measurements/measurers. What this means for
11 test-retest reliability is that even very large differences between test and retest values, which
12 may represent significant intra-rater instability or inter-rater difference, will not be detected
13 by correlation analysis if the differences are consistent in a sample. In practice, this means
14 that critical factors affecting measurement reliability such as order effects (practice, boredom,
15 fatigue, etc.) and user interpretation may never be identified. The values in table 2 can be
16 used as an example here; this table presents scores given by two teachers double marking a
17 computer-based exam task and the differences between the scores for each of the 14 students.
18 The table also presents a third set of transformed scores used to exemplify a large difference
19 (bias) in marking. If the way that pairs of measurements increase and decrease relative to one
20 another is constant, the correlation coefficients between measurements will be exactly the
21 same whether there is no bias, a small bias (e.g., around 1.5 points on average), or a very
22 large bias is present (e.g., around 46 points on average). Whilst we are unlikely to see
23 repeated measures differing by such a margin as teachers A and C in real-life data, this more
24 extreme example is used to illustrate an important point. In real world contexts systematic
25 bias can occur if a measurement tool is open to interpretation by the specific user, or where

1 learning and practice effects influence the outcomes of repeated measurements. These are
2 serious flaws for psychometric/psychological assessment tools, which should be unbiased and
3 standardized to permit comparison of measurements within and between samples. Given the
4 substantial negative impact that systematic bias has on agreement-based reliability (values
5 can be far from agreement), methods that conceal and are unaffected by such problems are
6 not appropriate for test-retest reliability analyses.

7

8

Table 2

9

10 Although it is uncommon in test-retest reliability, occasionally the measurements that
11 we want to compare have different outcome scales/units, but which denote exactly the same
12 practical result. For example, comparing height measurements between two raters, one of
13 whom uses inches and the other centimeters; or, comparing a total score to a mean or a
14 percentage score. In such cases, the outcomes must be standardized prior to analysis to permit
15 appropriate examination of agreement-based reliability.

16 **Correlation is influenced by sample variability.** A final drawback of correlation
17 analyses is that the strength of a correlation coefficient is influenced by the spread of the
18 measurements in a given sample. More heterogeneous samples will produce stronger
19 correlation coefficients than less heterogeneous samples, in the absence of any disparity in
20 the within-pair measurement differences of each. This means that the resulting correlation
21 coefficient is relative to the sample on which the analysis was based. While absolute
22 differences in coefficients may be relatively small when differences in spread are small, this
23 factor means that it may not be appropriate to directly compare correlation coefficients
24 produced from different samples and populations, and, that coefficients derived from narrow
25 sample ranges may not be representative of the broader population. For example, if you were

1 to compare reliability of growth measurements from a sample of children aged 3-5 years with
2 a sample of children aged 3-10 years, the latter group would be far more variable than the
3 former, so a larger correlation coefficient would be produced for the 3-10 year olds even if
4 agreement in absolute growth measures was the same for both samples. This would also be
5 relevant when comparing reliability estimates from clinical and non-clinical populations,
6 where variation in psychological outcome measures maybe highly disparate between groups.
7 As such the researcher may find that the tool appears more reliable in the non-clinical group
8 than the clinical group (or vice versa), when in actual fact the absolute differences in scores
9 in each group are comparable. It is important to note that this specific issue for test-retest
10 analysis does not arise as a result of narrow or incomparable samples (though these have their
11 own inherent issues if they are unrepresentative), but as a direct result of the use of a relative
12 method (i.e., correlation) to estimate reliability; therefore, it can be overcome by examining
13 absolute differences in scores.

14 The above issues surrounding correlation analysis also apply to regression analyses
15 when used to assess agreement, since simple regression of one measurement onto another is
16 also based upon association. These problems are particularly hazardous when data are not
17 plotted and examined visually, and reliability is endorsed based on statistical output alone.
18 An expectation of high agreement between measures may also lead to less rigorous
19 consideration of raw data and statistical results.

20 **Statistical tests of difference**

21 Reliance on traditional statistical testing and p-values can be a hindrance to reliability
22 analysis; *“performing a test of significance for a reliability coefficient is tantamount to*
23 *committing a type III error – getting the right answer to a question no one is asking”*
24 (Streiner, 2007; Streiner et al., 2014). While there are ongoing debates around the use and/or
25 over-reliance on p-values in research generally, the specific issue in this context is that the

1 null hypotheses against which many such statistical tests are compared are relatively
2 meaningless when assessing reliability. Perhaps the greatest issue relevant to test-retest
3 reliability analysis is the use of hypothesis driven tests of difference, such as the paired t-test.
4 The common fallacy is that, if a test finds no significant difference between measurements
5 then the measurements agree, but this is not the case (Altman & Bland, 1995). Finding a
6 difference to be ‘significant’ simply means that systematic variability between the
7 measurements (i.e., between raters, conditions, or time-points) outweighs the variability
8 within measurements (i.e., between the individuals in the sample). Therefore, even large
9 differences between repeated measurements, which indicate very poor agreement, can be
10 statistically non-significant if the sample being tested is heterogeneous. The inverse is also
11 true; very similar test-retest scores, which should be seen as demonstrating high reliability,
12 may differ statistically significantly in a homogenous sample.

13 A related error in reliability analyses is the belief that the average (mean) difference
14 between two or more conditions is adequate to quantify agreement between individual pairs
15 of scores. This error is demonstrated by the data in table 3, which presents another example
16 of laboratory food (pizza) preference ratings (0-20 scale) from 34 participants assessed on
17 two occasions. Table 3 also includes the within-pair differences for scores, the mean score for
18 each time-point, and the mean within-pair difference.

19

20

Table 3

21

22 Relative to the scale of measurement, the absolute differences between ratings are
23 large and variable, ranging from -4.75 to 8.85; and yet, the average within-pair difference is
24 only 0.71. This value suggests far greater similarity in the data than is actually the case.

25 Calculating the mean difference in scores can mask notable disparity between paired

1 measurements. This is particularly true when some scores increase from test to retest and
2 others decrease; whatever the reason for such a pattern in the within-pair differences (e.g.,
3 random error of measurement, heteroscedasticity, etc.), this leads to a combination of positive
4 and negative differences that cancel each other out and result in a mean close to zero. When
5 this data is assessed using a paired samples t-test ($t(33) = 1.27, p = 0.21$) or a Wilcoxon signed
6 rank test to take account of mild skew ($Z = 258.00, p = 0.50$), or worse still, an independent
7 samples t-test ($t(66) = 0.91, p = 0.37$), the difference between measurements is found to be
8 irrefutably non-significant.

9 Alongside widespread use of correlations and t-tests, a relatively small number of
10 studies in psychology report alternative, more direct methods of analysis for test-retest
11 reliability that utilise absolute differences in scores (e.g., Viglione, Blume-Marcovici, Miller,
12 Giromini, & Meyer, 2012), suggesting gradual improvement in the field. However, where
13 alternative approaches are taken, it can be difficult to determine the validity of the analyses
14 for the given context, or to replicate the methods, because limited methodological detail is
15 reported.

16 **Ways to improve test-retest reliability analysis in psychology**

17 **Meet the aims and requirements of test-retest reliability analysis**

18 Typically, test-retest reliability studies are undertaken to see if repeated measurements
19 are ‘similar enough’ for the tool to be considered reliable, which denotes assessment of
20 agreement between raw observed values (i.e., does each individual receive the same value
21 each time they are measured?). The most important outcome for this type of reliability is the
22 size of the differences between related measurements for each individual, rather than whether
23 a difference is seen on average, or whether a ‘significant’ result is obtained. Traditional
24 hypothesis driven tests assess whether an observed average difference or association is
25 statistically different from zero or no difference/association, rather than indicating how

1 similar/different the observed scores obtained from a tool are. In contrast, suitable methods
2 for analysing test-retest reliability examine the difference(s) between measurements for each
3 case in the sample at an individual level, and assess whether or not the absolute differences
4 between scores obtained by the tool fall within an acceptable range according to the tool's
5 specific clinical, scientific, or practical field of use. Unlike relative consistency, this relies on
6 having an agreed or directly observable unit of measurement for the outcome score. A
7 specific cut-off value (size of difference) up to which measurements may be considered to
8 agree, should be identified and justified by the researcher before viewing the data, to avoid
9 biasing the reliability analyses. Establishing reliability in this way facilitates more in-depth
10 examination of the data (e.g., the size and consistency of differences across a sample) and
11 hence more thorough evaluation of reliability. It also permits the creation and validation of
12 reference values and cut-off scores, for diagnosis and classification and for understanding a
13 single outcome score for an individual; something which is precluded in relative measures of
14 reliability since systematic scoring differences are permissible.

15 **Select suitable methods**

16 **Limits of Agreement.** Bland-Altman Limits of Agreement (LOA) (Bland & Altman,
17 1986) is a statistical method typically used to assess agreement between two repeated
18 numeric measurements (i.e., test-retest scores, or comparison of methods). LOA are based on
19 descriptive statistics for paired data and are typically accompanied by a plot of the data to aid
20 data checking and interpretation. The limits themselves represent the upper and lower
21 boundaries of the middle 95% range of the observed data (within-pair differences),
22 constructed around the mean within-pair difference as $\text{mean} \pm 1.96(\text{SD})$. For improved
23 interpretation and inference beyond the sample, confidence intervals are also constructed
24 around the upper and lower LOA. Confidence intervals around the LOA will be wider than
25 those around the mean by a factor of 1.71, when samples are not small (Bland & Altman,

1 1999). Assuming normality of the data, this gives a range of values in which we are 95%
2 confident the population limit should lie. The ‘population’ is a hypothetical scenario in which
3 all possible measurement differences could be measured, but it provides a practical indication
4 of the variability/precision of measurements that we might expect to see if the tool were
5 implemented widely (e.g., a new clinical or research assessment tool).

6 An associated Bland-Altman plot sees the average of the two paired measurements
7 plotted on the x axis, against the difference between the two measurements on the y axis. The
8 plot is used to examine the data distribution and to screen for outliers and heteroscedasticity
9 (where the level of agreement between measurements depends on, or is proportionate to the
10 size of the measurement). When constructing the LOA plot, a horizontal reference line is first
11 added to show the mean within-pair difference; we hope to see this line sitting as close to
12 zero as possible. Points spread evenly either side of zero show random error variability in
13 which measurements do not differ on average. Points lying around any other positive or
14 negative value would indicate systematic bias in the measurements, and if the amount that
15 points vary from the mean line differs across the range of measurements, this suggests that
16 the data are heteroscedastic. Heteroscedasticity may be dealt with via data transformation to
17 permit statistically valid calculation of LOA ((Bland & Altman, 1999). However, it would be
18 essential to try and determine the source of heterogeneity, and to discuss the implications of
19 this data pattern for reliability and wider application of the measurement tool.

20 If we use the food preference ratings presented in table 3 as an example, we saw
21 previously that the mean within-pair difference for this data was 0.71. Using the standard
22 deviation (3.29) and sample size ($n=34$) we can calculate the standard error (0.56) and a 95%
23 confidence interval for the mean (-0.39, 1.81). The LOA, which represent an interval
24 containing 95% of the observed differences, can be calculated as -5.73 (95% CI -7.64, -3.81)
25 to 7.16 (95% CI 5.24, 9.07); confidence intervals for the LOA are based on a standard error

1 of 0.96 (se mean \times 1.71). These key values are added to a Bland-Altman plot (figure 2) to
2 illustrate the extent of agreement, and hence reliability.

3

4

Figure 2

5

6 As we expect, the majority of data points fall within the LOA. If this is not the case, it
7 is likely that the data are skewed and thus, the validity of the LOA is questionable. Figure 2
8 shows that across the range of observed measurements, data points are randomly scattered
9 around a mean close to zero; this suggests that there is little systematic bias between the two
10 measurements and no obvious data heterogeneity. Negative differences represent a higher
11 score at measure two compared to measure one, while positive differences represent the
12 inverse.

13 To conclude about agreement, both the LOA plot and statistics should be examined to
14 ascertain how much measurements did (in the observed data) and could (according to the
15 confidence intervals) differ from one another, and if these differences are smaller than a
16 predetermined cut-off for reliability. If there is no systematic bias between measurements
17 (i.e., positive and negative differences are randomly distributed around zero), then either of
18 the limits of agreement (positive or negative) and the confidence interval around that limit
19 can be referenced to conclude about reliability in the wider population. In reality, the mean
20 within-pair difference may deviate a little from zero even from random variation alone, and
21 as shown in our example data, this will lead to an imbalance in the limits of agreement. To
22 make a conservative estimate of reliability, the larger of the two limits should be selected and
23 the confidence interval for this limit used to conclude about agreement. In our example data,
24 the larger of the limits of agreement was 7.16 (95% CI 5.24, 9.07), showing that 95% of the
25 paired measurements in the sample did not differ by more than 7.16 units. In addition, the

1 confidence interval tells us that we can be 95% confident that measurement differences
2 should not exceed 9.07 in the wider population of all measurements.

3 The disparity between the sample and confidence interval indexes of difference or
4 agreement presented above (7.16 vs. 9.07), illustrates how confidence intervals can alter our
5 conclusions about reliability beyond what is observed in the data, and highlight why it is so
6 important to quantify precision for all estimates. For example, if researchers working with the
7 data had chosen 10 as the maximum difference permitted for this tool to show agreement, we
8 would be confident that our tool was reliable; the chosen cut-off exceeds both the sample and
9 population limits. If instead the cut-off had been 5, we would be quite confident in
10 concluding that our tool was not reliable, since differences observed in the sample and
11 inferred for the population exceed this margin. The most difficult scenario is when the cut-off
12 lies between the two indexes. For example, if the cut-off had been 8, we would have to
13 discuss the implications of our uncertainty around reliability. The observed data do not
14 exceed this value, but reliability is not confidently supported in the context of the wider
15 population. The only way to minimize differences between sample and population estimates
16 is to study large samples, thus reducing the width of confidence intervals around the LOA.

17 **Intraclass Correlation.** While there remain frequent problems with reliability
18 analyses in psychology, the use of the Intraclass Correlation Coefficients (ICC) (Shrout &
19 Fleiss, 1979) has been seen in psychological literature for some time (e.g., Angold &
20 Costello, 1995; Egger et al., 2006; Grant et al., 2003; Kernot, Olds, Lewis, & Maher, 2015;
21 March, Sullivan, & Parker, 1999; Silverman, Saavedra, & Pina, 2001). Unlike Pearson's
22 (interclass) correlation, ICC is an acceptable measure of reliability between two or more
23 measurements on the same individual/case. Despite the name and the presence of a
24 coefficient to quantify reliability, ICC is actually based on a ratio of rater, participant, and
25 error sources of measurement variability (derived from ANOVA models). This does mean

1 that ICC coefficients are, like other inferential tests, influenced by sample homogeneity;
2 when variability between measurements is constant, the more alike the sample is, the lower
3 the ICC will be (Bland & Altman, 1990; Lee et al., 2012). Therefore, ICC coefficients
4 derived from samples whose outcome variances differ, such as non-clinical and clinical
5 samples, should not be compared directly. For example, if a depression measure was used in
6 a non-clinical sample we would expect a modest range of scores with many cases scoring
7 close to zero, but this same tool applied to a sample of depressed individuals would likely
8 produce a much greater range of scores. In this case, the clinical sample would obtain a
9 higher ICC coefficient than the more homogenous non-clinical sample, in the absence of any
10 difference in the tool's reliability. This factor does not discredit ICC as a method of reliability
11 analysis, but highlights the importance of evaluating reliability using a representative sample
12 drawn from a relevant population (i.e., in which the tool will be used) (Bland & Altman,
13 1990). It also emphasizes the need to consider sample variance when interpreting ICC
14 coefficients and differences in reliability observed between samples and populations.

15 ICC coefficients quantify the extent to which multiple ratings for each individual
16 (within-individual) are statistically similar enough to discriminate between individuals, and
17 should be accompanied by a confidence interval to indicate the precision of the reliability
18 estimate. Most statistical software will also present a p-value for the ICC coefficient. This p-
19 value is obtained by testing sample data against the null hypothesis that measurements
20 within-person are no more alike than between-people (i.e., there is no reliability). In contrast,
21 reliability studies aim to answer the functional question 'are the repeated measurements made
22 using a tool similar enough to be considered reliable'. As such, the p-value provided is, in
23 most cases, of little practical use or relevance.

24 Though many authors report simply that '*ICC was used*' there are in fact six different
25 ICC types to suit different theoretical and methodological study designs (Atkinson & Nevill,

1 1998; Shrout & Fleiss, 1979). ICC can be used when a single sample of raters is used to
2 assess every individual in a test sample (type 2 ICC), or when different, randomly selected
3 raters are used across the total sample (type 1 ICC; e.g., when the same individuals cannot
4 feasibly make all measurements across a sample, such as national/multi-center studies). ICC
5 types 1 and 2 quantify agreement. A third ICC type (type 3) is used to assess consistency
6 among a fixed group of raters. Type 3 ICC permits systematic differences between raters, and
7 so represents consistency rather than agreement; therefore, it is only suitable when relative
8 reliability is of primary importance; as discussed earlier in this paper, this is infrequently the
9 case when assessing test-retest reliability for psychometric measures.

10 For each of the three main ICC types outline above, the coefficient can be calculated
11 in two ways; the first reflects the contributions of each individual rater (e.g., presented as
12 *Single Measures* in SPSS, *Single_raters* in R, and *Individual* in STATA), while the second
13 uses an average of raters (*Average Measures* in SPSS, *Average_raters* in R, or *Average* in
14 STATA). Average options will always result in a larger coefficient, because averaging dilutes
15 the differences across raters/ratings and gives a false inflation of agreement.

16 Three ICC types and two methods of calculation for each, translates into six different
17 ICC coefficients that could be calculated for any given set of data. However, the coefficients
18 will differ in size, the meaning of ‘reliability’ that they represent, and validity for the
19 particular study. Valid choice of ICC type should be determined by the selection of raters in
20 the particular study, whether or not reliability needs to be generalized to a wider population
21 (e.g., inter-rater reliability generalized to other clinicians using a given measure), and whether
22 consistency or agreement is required. This decision should be clearly outlined in the methods
23 section of research reports (Atkinson & Nevill, 1998; Krebs, 1986).

24 As an applied example, we can revisit the bias data in table 2. This table presented
25 data from 3 teachers (A, B, and C); teacher B scored on average 1.5 points higher than

1 teacher A, while teaching C scored on average 46 points higher than A. We saw previously
2 that correlation fails to recognize systematic bias and as such the correlations for A with B
3 and A with C were both 0.99 and highly statistically significant. If we now assess this data
4 with ICC type 2 (assuming a sample of teachers were used to assess the random sample of 14
5 students) to look at agreement, we find that good reliability is demonstrated for teachers A
6 and B who marked similarly (ICC (single measures) = 0.97), and appropriately, very poor
7 reliability is shown for teachers A and C who marked differently (ICC (single measures) =
8 0.16). When all three teachers are added into the ICC model we see a negligible increase to
9 0.18, suitably reflecting poor reliability across all three raters. As expected, when ICC type 3
10 is run, which treats raters as fixed and allows for systematic variability between raters, the
11 result is a considerably higher ICC coefficient of 0.49, which would be higher still if the bias
12 between raters, however large, was consistent. This again highlights the deficiency of
13 assessing consistency rather than agreement for test-retest types of reliability.

14 An ICC coefficient can also be accompanied by an ICC plot, which sees the sample
15 cases plotted in the x axis, outcome scores on the y axis, and different point characters used
16 for each rater/rating. ICC plots illustrate the size and nature of observed differences between
17 raters/ratings, and the clustering of scores within person relative to variability across the
18 sample, which aid the practical interpretation of statistical results. For example, figure 3
19 presents the exam marking data from table 2 for teachers A and B; from this plot we see that
20 teacher A scores consistently lower than teacher B, indicating a small bias, but in most cases
21 the marks are similar. In contrast, figure 4 presents the table 2 data for all three teachers
22 together. This plot clearly shows that teacher C marks much higher than teachers A and B,
23 representing a large positive bias, and hence poor reliability.

24

25

Figure 3

Figure 4

Improved Reporting

In any study, the aims of the research and the methods used to meet those aims should be clearly outlined; it is insufficient to present vague conclusions verified only by statistical output (e.g., ‘good reliability was shown $r_{(100)} = 0.85$, $p = 0.006$ ’). The purpose of test-retest reliability studies is to provide evidence that a tool will measure the same thing, in the same way, each time it is used (validity assessment then tells us if it is measuring the right thing). If the methods used to evidence this reliability are not sufficiently explained to validate their use, or the evidence is not presented in the context of a wider population (i.e., no confidence intervals), then the evidence is compromised, or absent altogether. Statements such as ‘ICC was used to assess reliability’ are common, despite important differences in ICC models, and the implications of their selection. Such reporting provides no evidence of mindful selection of methods, and may lead the reader to infer that software default settings were used, which vary between packages and may not be appropriate. For example, the default ICC type in IBM SPSS Statistics (version 22) is a two-way mixed effects model for consistency (type 3 ICC). This model is liable to give the highest ICC coefficient of all three main types, but is only appropriate to use when a fixed group of raters is used, and consistent differences between those raters are unimportant. This is contrary to test-retest studies that aim to examine agreement between measurements. It should also be clearly specified and justified when an average of raters is used rather than assessing across individual raters, since this will always inflate the resulting reliability coefficient.

Problems regarding the justification of analytical choices in ICC also extend to correlations and inferential tests of difference. Often, the application of these tests is stated, but neither a rationale for their selection, nor an explanation of how the results demonstrate

1 reliability, are given by the author. These omissions should lead readers to distrust the results,
2 but this is not always the case. Reporting of reliability studies should follow the same
3 recommendations for reporting any research methodology; the information given should be
4 sufficient to allow the reader, in principle, to replicate the study. Complete evidence of
5 reliability, or indeed, unreliability, includes information relevant to the methods and results of
6 the analysis. This should include what will be examined (e.g., agreement between test and
7 retest scores), how this will be assessed (e.g., Bland Altman limits of agreement), and why
8 the method was chosen (e.g., because limits of agreement assesses the extent to which paired
9 measurement in a sample agree). Authors should also clearly document what the results of
10 the assessment indicate about the data and about subsequent use of the tool, relative to
11 practical/clinical parameters and requirements for reliability.

12 There are some good examples of analyses and reporting within the psychological
13 literature (e.g., Grant et al., 2003; Kernot et al., 2015; Tighe et al., 2015) that demonstrate
14 concise yet informative methodological information. The cited authors are clear about the
15 statistical methods they chose; for example, Tighe and colleagues (2015) reported that they
16 “*calculated the intra-class correlations (ICC) using a two-way mixed effects model for the A,*
17 *B, and total Alda Scale scores*”. Similarly, Grant and colleagues (2003) reported that “*For*
18 *continuous measures, intraclass correlation coefficients (ICC) are presented as measures of*
19 *reliability. Since our reliability design assumed that interviewers were randomly drawn from*
20 *a larger population of interviewers, we used a one-way random effects ANOVA model to*
21 *derive intraclass correlation coefficients (Shrout and Fleiss, 1979).*” As well as providing
22 important details regarding the specific ICC models applied to their data and, in the case of
23 Grant et al (2003), the justification for this choice, both papers also presented 95%
24 confidence intervals alongside their ICC coefficients. This permits a greater level of
25 interpretation regarding the precision and wider applicability of their results. This information

1 gives the reader a much better indication of what specific statistical procedures were carried
2 out, from which they can better judge the suitability and strength of the resulting evidence.

3 **Extending the principles to other tests of reliability**

4 Although categorical data has not been discussed in the current paper, the key
5 principles of reliability analysis that have been discussed within a test-retest design can be
6 directly translated to these types of outcomes. Firstly, data should be considered at an
7 individual, paired level, in both presentation and analysis. Secondly, analyses should assess
8 the agreement between measurements from multiple conditions, times, or raters. And finally,
9 inferential tests of difference/association, such as chi squared and McNemar's tests for
10 categorical outcomes, should be avoided in favor of specific tests of agreement such as kappa
11 (Cohen, 1960) and its extensions (e.g., weighted kappa, generalized kappa).

12 The current paper has by no means offered an exhaustive list of potential methods for
13 assessing measurement reliability. Instead, two example analyses have been used to illustrate
14 what an appropriate assessment of agreement-based reliability should comprise. The
15 fundamental messages of the current paper aim to help researchers choose a test or method
16 that actually quantifies reliability, draw conclusions about reliability as directly as possible
17 from the data, and recognize that in most cases a p-value, if given, will provide little practical
18 information about the use or reliability of a measurement tool.

19

20

References

- 1
2 Altman, D. G., & Bland, J. M. (1983). Measurement in Medicine - the Analysis of Method
3 Comparison Studies. *Statistician*, 32(3), 307-317. doi:10.2307/2987937
- 4 Altman, D. G., & Bland, J. M. (1995). Statistics Notes: Absence of Evidence is not Evidence
5 of Absence. *British Medical Journal*, 311(7003), 485-485.
6 doi:<http://dx.doi.org/10.1136/bmj.311.7003.485>
- 7 Angold, A., & Costello, E. J. (1995). A Test-Retest Reliability Study of Child-Reported
8 Psychiatric-Symptoms and Diagnoses Using the Child and Adolescent Psychiatric-
9 Assessment (Capa-C). *Psychological Medicine*, 25(4), 755-762.
- 10 Atkinson, G., & Nevill, A. M. (1998). Statistical methods for assessing measurement error
11 (reliability) in variables relevant to sports medicine. *Sports Medicine*, 26(4), 217-238.
- 12 Baumgartner, T. A. (2000). Estimating the Stability Reliability of a Score. *Measurement in*
13 *Physical Education and Exercise Science*, 4(3), 175-178.
- 14 Bedard, M., Martin, N. J., Krueger, P., & Brazil, K. (2000). Assessing reproducibility of data
15 obtained with instruments based on continuous measurements. *Experimental Aging*
16 *Research*, 26(4), 353-365. doi:Doi 10.1080/036107300750015741
- 17 Bennett, R. J., & Robinson, S. L. (2000). Development of a measure of workplace deviance.
18 *Journal of Applied Psychology*, 85(3), 349-360. doi:10.1037//0021-9010.85.3.349
- 19 Bland, J. M., & Altman, D. G. (1986). Statistical Methods for Assessing Agreement between
20 Two Methods of Clinical Measurement. *Lancet*, 1(8476), 307-310.
- 21 Bland, J. M., & Altman, D. G. (1990). A note on the use of the intraclass correlation
22 coefficient in the evaluation of agreement between two methods of measurement.
23 *Computers in Biology and Medicine*, 20(5), 337-340.

- 1 Bland, J. M., & Altman, D. G. (1999). Measuring agreement in method comparison studies.
2 *Statistical Methods in Medical Research*, 8(2), 135-160.
3 doi:10.1191/096228099673819272
- 4 Chmielewski, M., & Watson, D. (2009). What is being assessed and why it matters: The
5 impact of transient error on trait research. *Journal of Personality and Social*
6 *Psychology*, .97(1), pp. doi:10.1037/a0015618 19586248
- 7 Cohen, J. (1960). A Coefficient of Agreement for Nominal Scales. *Educational and*
8 *Psychological Measurement*, 20(1), 37-46. doi:10.1177/001316446002000104
- 9 Egger, H. L., Erkanli, A., Keeler, G., Potts, E., Walter, B. K., & Angold, A. (2006). Test-
10 retest reliability of the Preschool Age Psychiatric Assessment (PAPA). *Journal of the*
11 *American Academy of Child and Adolescent Psychiatry*, 45(5), 538-549.
12 doi:10.1097/01.chi.0000205705.71194.b8
- 13 Garner, D. M., Olmstead, M. P., & Polivy, J. (1983). Development and Validation of a
14 Multidimensional Eating Disorder Inventory for Anorexia-Nervosa and Bulimia.
15 *International Journal of Eating Disorders*, 2(2), 15-34. doi:10.1002/1098-
16 108x(198321)2:2<15::Aid-Eat2260020203>3.0.Co;2-6
- 17 Goodman, R. (1997). The Strengths and Difficulties Questionnaire: a research note. *Child*
18 *Psychology & Psychiatry & Allied Disciplines*, 38(5), 581-586. doi:10.1111/j.1469-
19 7610.1997.tb01545.x 9255702
- 20 Goodman, R. (2001). Psychometric properties of the Strengths and Difficulties
21 Questionnaire. *Journal of the American Academy of Child & Adolescent Psychiatry*,
22 .40(11), 1337-1345. doi:10.1097/00004583-200111000-00015
- 23 Grant, B. F., Dawson, D. A., Stinson, F. S., Chou, P. S., Kay, W., & Pickering, R. (2003).
24 The Alcohol Use Disorder and Associated Disabilities Interview Schedule-IV
25 (AUDADIS-IV): reliability of alcohol consumption, tobacco use, family history of

- 1 depression and psychiatric diagnostic modules in a general population sample. *Drug*
2 *and Alcohol Dependence*, 71(1), 7-16. doi:10.1016/S0376-8716(03)00070-X
- 3 Hole, G. (2014). Eight things you need to know about interpreting correlations. Retrieved
4 from
5 [http://www.sussex.ac.uk/Users/grahamh/RM1web/Eight%20things%20you%20need](http://www.sussex.ac.uk/Users/grahamh/RM1web/Eight%20things%20you%20need%20to%20know%20about%20interpreting%20correlations.pdf)
6 [%20to%20know%20about%20interpreting%20correlations.pdf](http://www.sussex.ac.uk/Users/grahamh/RM1web/Eight%20things%20you%20need%20to%20know%20about%20interpreting%20correlations.pdf)
- 7 Kernot, J., Olds, T., Lewis, L. K., & Maher, C. (2015). Test-retest reliability of the English
8 version of the Edinburgh Postnatal Depression Scale. *Archives of Women's Mental*
9 *Health*, .18(2), pp. doi:10.1007/s00737-014-0461-4 25209355
- 10 Krebs, D. E. (1986). Declare your ICC type. *Physical Therapy*, 66(9), 1431-1431.
- 11 Lee, K. M., Lee, J., Chung, C. Y., Ahn, S., Sung, K. H., Kim, T. W., . . . Park, M. S. (2012).
12 Pitfalls and important issues in testing reliability using intraclass correlation
13 coefficients in orthopaedic research. *Clinics in Orthopedic Surgery*, 4(2), 149-155.
14 doi:10.4055/cios.2012.4.2.149
- 15 Ludbrook, J. (2002). Statistical techniques for comparing measurers and methods of
16 measurement: A critical review. *Clinical and Experimental Pharmacology and*
17 *Physiology*, 29(7), 527-536. doi:DOI 10.1046/j.1440-1681.2002.03686.x
- 18 March, J. S., Sullivan, K., & Parker, J. (1999). Test-retest reliability of the multidimensional
19 anxiety scale for children. *Journal of Anxiety Disorders*, 13(4), 349-358. doi:Doi
20 10.1016/S0887-6185(99)00009-2
- 21 Meyer, T. J., Miller, M. L., Metzger, R. L., & Borkovec, T. D. (1990). Development and
22 validation of the Penn State Worry Questionnaire. *Behaviour Research and Therapy*,
23 28(6), 487-495.

- 1 Pliner, P., & Hobden, K. (1992). Development of a Scale to Measure the Trait of Food
2 Neophobia in Humans. *Appetite*, *19*(2), 105-120. doi:10.1016/0195-6663(92)90014-
3 W
- 4 Rust, J., & Golombok, S. (2009). *Modern psychometrics: The science of psychological*
5 *assessment., 3rd ed.* New York, NY: Routledge/Taylor & Francis Group; US.
- 6 Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: uses in assessing rater reliability.
7 *Psychological Bulletin*, *86*(2), 420-428.
- 8 Silverman, W. K., Saavedra, L. M., & Pina, A. A. (2001). Test-retest reliability of anxiety
9 symptoms and diagnoses with the anxiety disorders interview schedule for DSM-IV:
10 Child and parent versions. *Journal of the American Academy of Child and Adolescent*
11 *Psychiatry*, *40*(8), 937-944. doi:Doi 10.1097/00004583-200108000-00016
- 12 Steptoe, A., Pollard, T. M., & Wardle, J. (1995). Development of a measure of the motives
13 underlying the selection of food: the food choice questionnaire. *Appetite*, *25*(3), 267-
14 284. doi:10.1006/appe.1995.0061
- 15 Streiner, D. L. (2007). A shortcut to rejection: How not to write the results section of a paper.
16 *Canadian Journal of Psychiatry-Revue Canadienne De Psychiatrie*, *52*(6), 385-389.
- 17 Streiner, D. L., Norman, G. R., & Cairney, J. (2014). *Health measurement scales : a practical*
18 *guide to their development and use.* Oxford, United Kingdom: Oxford University
19 Press.
- 20 Tighe, S. K., Ritchey, M., Schweizer, B., Goes, F. S., MacKinnon, D., Mondimore, F., . . .
21 Potash, J. B. (2015). Test-retest reliability of a new questionnaire for the retrospective
22 assessment of long-term lithium use in bipolar disorder. *Journal of Affective*
23 *Disorders*, *174*, 589-593. doi:10.1016/j.jad.2014.11.021

1 Viglione, D. J., Blume-Marcovici, A. C., Miller, H. L., Giromini, L., & Meyer, G. (2012). An
2 inter-rater reliability study for the rorschach performance assessment system. *Journal*
3 *of Personality Assessment*, 94(6), 607-612. doi:10.1080/00223891.2012.684118

4 Weir, J. P. (2005). Quantifying test-retest reliability using the intraclass correlation
5 coefficient and the SEM. *Journal of Strength and Conditioning Research*, 19(1), 231-
6 240. doi:Doi 10.1519/15184.1

7