

Analytical Methods

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1
2
3 **Is your ‘homogeneity test’ really useful?**
4
5
6

7 Michael Thompson
8 School of Biological and Chemical Sciences
9 Birkbeck University of London
10 Malet Street
11 London WC1E 7HX, UK
12

13
14 **Abstract**
15

16 ‘Homogeneity testing’ is a formal requirement in the preparation of reference materials for certification
17 and for proficiency testing. Few scientists, however, seem to be aware of the rather severe limitations
18 that apply to the outcome of any such test of a size that is economically feasible. Typically the tests
19 have low statistical power to detect significance, and the resulting estimates of between-bottle standard
20 deviation have wide confidence limits. Scientists should bear these limitations in mind and avoid being
21 over-prescriptive when drafting standards and guides.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
Reference materials are essential in chemical measurement for demonstrating comparability of results over space and time. Such materials are usually divided into many portions (here called ‘bottles’ to conform to ISO Guide 35¹), in activities such as proficiency testing, collaborative trials, internal quality control and the preparation of certified reference materials. Clearly all bottles of a reference material should have the same composition within margins that are optimally narrow in relation to the mass of a test portion and the cost of using the material. Some degree of testing is needed to show how closely that need is fulfilled. The test naturally adapted to the task is based on experiments with randomised replicated measurements, followed by analysis of variance (ANOVA).

14
15
16
17
18
19
20
21
22
23
24
25
26
27
ISO Guide 35 details the measures required to ensure that the test provides an unambiguous outcome. For instance, a random selection of bottles is essential to make the expectation of the test representative. A further requirement is the need to ensure that compositional variation among bottles of the test material is not confounded with systematic variation (drifts or saltations, for example) in the measurement procedure. For this reason the Guide recommends a random order in the sequence of measurements on the test portions or tests for time trends in the outcome. However, an unexpected consequence of the Guide is that it has encouraged the unwary to believe that its recommendations are sufficient as well as necessary. There is a tendency for scientists to believe that working according to the prescriptions of the Guide frees the user from the further responsibility of looking at all of the implications and examining the features that determine whether a particular design of homogeneity test is worthwhile.

28
29
30
31
32
33
This paper is designed to redress the balance. First let’s dwell on some ‘home-truths’ about homogeneity and ANOVA that impinge on our discussion. If we ignore these realities we become prone to misconceptions about why we should conduct homogeneity tests and how we should apply the results.

- 34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
- All reference materials are heterogeneous, at least in principle and, nearly always, in practice. The important criterion for analytical chemists is that the inevitable variation in composition between bottles, determined on appropriately-sized test portions, is sufficiently small in relation to an uncertainty consistent with fitness for purpose. This has given rise to the apt but clumsy designation ‘sufficiently [close to] homogeneous’ in preference to ‘homogeneous’.
 - There is no test for homogeneity: this is despite the almost universal misuse of the phrase ‘homogeneity testing’. We *can* test for a statistically significant degree of heterogeneity, which is not the same thing, but we must reconcile ourselves to the fact that an affordable test is unlikely to be sufficiently powerful.
 - The tests also provide estimates of between-bottle dispersion that (a) are likely to differ wildly from the true value and (b) are very sensitive to mildly unusual results unless robust ANOVA is used.

49 50 51 52 53 54 55 56 57 58 59 60

Significance, importance, and power

Tests for significant heterogeneity in chemical measurement that are on an economically-feasible scale are notoriously low in statistical power². This means in practice that the test is unlikely to detect consequential degrees of heterogeneity at 95% confidence in a large proportion, perhaps a substantial majority of instances. The bigger the experiment, and the more precise the analytical results, the better it will be for characterising an important degree of heterogeneity, but the more it will cost to conduct.

1
2
3 The overall message that emerges is that it is usually easy to render a material sufficiently close to
4 homogeneous, but difficult to find an analytical procedure with sufficient precision adequately to
5 explore the real variation between bottles by means of an experiment of economically-feasible size.
6 This is perhaps made clearer when the outcome of the test is focussed on the magnitudes of the
7 components of dispersion rather than on statistical significance. We find that the estimates of the
8 variance components other than analytical are unreliable in magnitude. For instance, proficiency test
9 providers often use a test with duplicate analysis on 10 bottles of the material, a (10×2) test.
10 Simulations can show us what to expect under the assumption that the analytical results were
11 independent, random, normally-distributed variables with repeatability standard deviation σ_r . By
12 conducting such a test on a hypothetical material that was indeed homogeneous (that is, the true
13 between-bottle standard deviation $\sigma_b = 0$), we find that the 95% confidence limits on the estimated
14 value $\hat{\sigma}_b$ would be $(0, 0.82\sigma_r)$. About half of the estimates would be zero, but the mean would be
15 about $0.24\sigma_r$ and results as high as $0.82\sigma_r$ would not be rare. The outcome is hardly better when the
16 between-bottle standard deviation is as large as $\sigma_r/2$.
17
18
19
20
21
22
23

24 **Information on between-bottle variance by using enhanced power—a unique example**

25
26 This shortcoming has recently been demonstrated in a study of heterogeneity in foodstuff proficiency-
27 test materials, namely products similar in composition to meat pies³. Two analytes were selected for
28 this pilot study, nitrogen (a proxy for protein) and fat. Nitrogen was considered because its
29 concentration in successive test materials was restricted to a small range, so each variance component
30 could reasonably be taken as homoscedastic (that is, of uniform variance). The mass fraction of fat, in
31 contrast, varied between about 2% and 30%, but it was found that the variances could be rendered close
32 to homoscedastic by considering standard deviations (SD) scaled to concentration (that is, as relative
33 standard deviations (RSD)). (Homoscedasticity is required to fulfil the requirements of ANOVA
34 model.)
35
36
37

38 Data from routine testing (10×2 experiments) in 20 successive rounds of the proficiency test were
39 analysed by hierarchical ANOVA. This boosted the degrees of freedom for the between-bottle variance
40 20-fold, greatly increasing the power of the statistical analysis and providing for the first time plausible
41 estimates of the corresponding standard deviations. For nitrogen the between-bottle SD estimate was
42 0.008% mass fraction (about one third of the analytical SD). For fat the between-bottle RSD was 0.006,
43 again about one third of the corresponding analytical RSD.
44
45

46 Bodies that organise the preparation of reference materials go to great lengths to ensure that the
47 resulting material is as close to homogeneous as can be reasonably achieved. As shown here, with a
48 carefully prepared material and with ten bottles analysed in duplicate, the routine test is usually unable
49 to detect the low level of heterogeneity. Moreover, the estimate of between-bottle variance is
50 unreliable, being inordinately disperse, and often substantially biased.
51
52
53

54 **A larger example, taken from ISO Guide 35, Appendix B3**

55
56
57
58
59
60

1
2
3 The example depicts results of a test for between-bottle heterogeneity of chromium in a certification
4 study. The study is of substantial size (20×3), that is, requiring the analysis of 60 separate test portions.
5 The data are shown graphically in Fig 1 and show a variation between bottles that is significant at 95%
6 confidence. Calculations displayed in the Guide show that $\hat{\sigma}_b = 3.93 \text{ mg kg}^{-1}$. However, even with this
7 large experiment and a statistically-significant outcome, the estimate $\hat{\sigma}_b$ is very variable, with 95%
8 confidence limits, estimated by the bootstrap⁴, of (1.78, 5.33).
9
10

11
12 The outcome is also sensitive to the presence of mild outliers. For example, if just one of the 60 data is
13 changed to a result of 140 mg/kg, as shown in Fig 1, the result of the ANOVA becomes not significant
14 at 95% confidence and the estimate is considerably changed, to $\hat{\sigma}_b = 1.82$ with 95% confidence limits
15 of (0.00, 4.05) mg/kg. The influence of a small proportion of analytical outliers can be largely
16 eliminated by the use of a robust ANOVA. In these applications, however, it is essential to ensure that
17 the robustification is applied only at the within-bottle level of the design: it then accommodates only
18 analytical outliers. Between-bottle variation represents real heterogeneity and variation at that level
19 must not be robustified. Application of such an ANOVA to the modified data above gives estimates
20 that are close to those of the original data, namely $\hat{\sigma}_{b,robust} = 3.85$ with 95% confidence limits of (2.34,
21 5.56) mg/kg.
22
23
24
25
26

27 **Heterogeneity and certification of a reference material**

28
29 The ‘official’ procedure for attaching an uncertainty to a certified value for a reference material
30 includes a term for between-bottle (and sometimes within-bottle) heterogeneity¹. This is formally
31 correct if we know the population parameter. But we have only an estimate of the parameter and, as we
32 have seen just above, the estimate is likely to be seriously in error. Using such a term for comparison
33 with a separate fitness-for-purpose criterion could be badly misleading. When the overall uncertainty
34 on the certified value is estimated, the difficulty can be obviated by subsuming the heterogeneity
35 contribution into the repeatability or reproducibility standard deviation.
36
37
38

39 We should also consider whether it is worthwhile separately to study the within-bottle heterogeneity.
40 Producers of certified reference materials are encouraged to do this presumably on the grounds that it
41 contributes to repeatability dispersion. But the estimates of this statistic are likely to be at least as
42 unreliable as between-bottle estimates and inflate the cost of the test twofold. Wouldn't the same
43 financial layout be better used by doubling the number of bottles in the experiment and ignoring
44 within-bottle variation? There is no general answer to this question. We must further remember that
45 many ‘homogenised’ materials have a tendency to segregate in their containers during transport and
46 use. Therefore even a hypothetically-reliable measure of within-bottle heterogeneity carried out by the
47 producer will be untrustworthy by the time the bottle is on the analyst's bench. We should also bear in
48 mind that, in routine practice, it is part of the analyst's job to ensure as far as possible that the test
49 portion taken is representative of the laboratory sample (using both words as strictly defined⁵). This
50 responsibility should apply equally to use of reference materials. Now risk of contamination obviously
51 prevents an analyst from removing a certified reference material from its container for further treatment
52 such a grinding. A thorough shaking within the closed vessel is the most that could be sanctioned, and
53 is always an essential preliminary before removal of the test portion. It would be clearly useful to know
54
55
56
57
58
59
60

1
2
3 that a preliminary shaking might restore the bottle's contents to sufficient homogeneity, but the cost of
4 that knowledge is likely to remain prohibitively large.
5
6

7 Finally we should consider whether in testing for heterogeneity it is worthwhile to use more elaborate
8 experimental designs, such as conducting the analysis in randomised blocks. This measure can be
9 useful when results are prone to suffer from drift in long runs of analysis or when a sufficiently long
10 sequence of analysis cannot be accommodated in a single run. (A 'run' is the period during which
11 repeatability conditions can be assumed to prevail.) This would tend effectively to improve the
12 precision of the analytical results. A minor disadvantage of a blocked design would be that a few
13 degrees of freedom are wasted on a feature that does not otherwise affect the question at issue, the
14 heterogeneity of the candidate material. Unfortunately including a run-to-run effect is apparently often
15 misinterpreted as an inherently desirable feature of the test for heterogeneity, whereas each instance
16 should be considered on its merits.
17
18
19

20 **Preliminary tests of datasets for deviations from procedure**

21 It should be regarded as essential practice to test datasets from heterogeneity studies to ensure that they
22 show no obvious sign that the analyst had deviated from the selected design so as to invalidate the
23 outcome of the ANOVA. Deviations from a strictly random ordering of the test portions in an
24 analytical procedural sequence are not uncommon. Features often encountered by the author include (a)
25 a trend or step-change in the results ordered by the bottle numbers or by order of analysis, (b) a
26 systematic difference between the first and second results on each bottle, (c) insufficient digit
27 resolution in the data for a reliable statistical analysis, and (d) outlying differences between duplicate
28 results. These features tend to emerge when the instructions for conducting the experiment are
29 insufficiently detailed. Many of these problems can be immediately spotted on a graph of the data,
30 organised in the order of analysis and (as in Fig 1,) organised by bottle number, and confirmed by a test
31 of significance.
32
33
34
35

36 **Conclusions**

37
38 'Homogeneity tests' usually fail to deliver what at first sight they seem to promise. Even with datasets
39 faultlessly produced according to an appropriate experimental design, the tests seldom detect
40 significant heterogeneity because experiments that are economically feasible have insufficient power.
41 They deliver estimates of between-bottle standard deviation that are wildly variable, biased, and very
42 sensitive to outlying analytical results. Awareness of these problems would be reinforced if the
43 statistical packages delivering ANOVA provided confidence limits on the estimated standard
44 deviations. ANOVA, robustified at the analytical level but giving these confidence limits, would be a
45 invaluable tool for 'homogeneity testing'. That would alert analysts to the possibility of improper
46 interpretation of the outcome.
47
48
49

50 A test for heterogeneity provides important reassurance to users and as such can hardly be dispensed
51 with. But it is suited better to screening materials for complete failures in the homogenisation process
52 or other type of mistake than to studying heterogeneity *per se*. In proficiency tests, where 10×2 tests are
53 at the limit of affordability, there is no question of doing more. In the preparation of certified reference
54 materials, a considerably greater expenditure may be justifiable, but even here the outcome may fall
55 short of providing all of the information that analysts would like to see.
56
57
58
59
60

There is an unfortunate relationship between the heterogeneity of a reference material and the precision of the analytical procedures that it is designed to control. To fulfil its purpose adequately, the material must have a between-bottle dispersion that is small in relation to the analytical variation of the procedure that it is used to monitor. But in that circumstance, that particular analytical procedure will be unable to detect the heterogeneity in a reasonably-sized experiment. In many, perhaps most instances, there will be no economically acceptable alternative analytical procedure with better precision. For between-bottle heterogeneity in a reference material, 'if you can measure it properly, it's too big'.

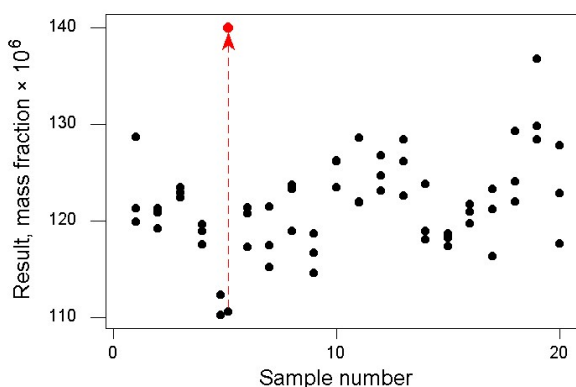


Fig 1. Results (black solid circles) for the triplicate analysis for chromium of 20 bottles of a candidate reference material. The red solid circle indicates a modification of the data to investigate the influence of a single result on the estimate of between-bottle standard deviation.

¹ ISO Guide 35: 2006. *Reference materials -- General and statistical principles for certification*.

² M Thompson. *Accred Qual Assur*, 2008, **13**, 581-584.

³ M Thompson and T Fearn. *Analytical Methods*, 2011, **3**, 2529-2533

⁴ Analytical Methods Committee. *AMC Technical Briefs No. 8 (2001). The bootstrap: a simple approach to estimating standard errors and confidence intervals when theory fails*. (Download gratis via www.rsc.org/amc)

⁵ M H Ramsey and S L R Ellison (eds), Eurachem/Eurochem/CITAC/ Nordtest/AMC Guide. *Measurement uncertainty arising from sampling. A guide to methods and approaches*. Eurachem 2007.