



# Standard Practice for Statistical Assessment and Improvement of Expected Agreement Between Two Test Methods that Purport to Measure the Same Property of a Material<sup>1</sup>

This standard is issued under the fixed designation D6708; the number immediately following the designation indicates the year of original adoption or, in the case of revision, the year of last revision. A number in parentheses indicates the year of last reapproval. A superscript epsilon ( $\epsilon$ ) indicates an editorial change since the last revision or reapproval.

## 1. Scope\*

1.1 This practice covers statistical methodology for assessing the expected agreement between two standard test methods that purport to measure the same property of a material, and deciding if a simple linear bias correction can further improve the expected agreement. It is intended for use with results collected from an interlaboratory study meeting the requirement of Practice D6300 or equivalent (for example, ISO 4259). The interlaboratory study must be conducted on at least ten materials that span the intersecting scopes of the test methods, and results must be obtained from at least six laboratories using each method.

1.2 The statistical methodology is based on the premise that a bias correction will not be needed. In the absence of strong statistical evidence that a bias correction would result in better agreement between the two methods, a bias correction is not made. If a bias correction is required, then the *parsimony principle* is followed whereby a simple correction is to be favored over a more complex one.

NOTE 1—Failure to adhere to the parsimony principle generally results in models that are over-fitted and do not perform well in practice.

1.3 The bias corrections of this practice are limited to a constant correction, proportional correction or a linear (proportional + constant) correction.

1.4 The bias-correction methods of this practice are method symmetric, in the sense that equivalent corrections are obtained regardless of which method is bias-corrected to match the other.

1.5 A methodology is presented for establishing the 95 % confidence limit (designated by this practice as the *between methods reproducibility*) for the difference between two results where each result is obtained by a different operator using different apparatus and each applying one of the two methods

X and Y on identical material, where one of the methods has been appropriately bias-corrected in accordance with this practice.

NOTE 2—In earlier versions of this standard practice, the term “cross-method reproducibility” was used in place of the term “between methods reproducibility.” The change was made because the “between methods reproducibility” term is more intuitive and less confusing. It is important to note that these two terms are synonymous and interchangeable with one another, especially in cases where the “cross-method reproducibility” term was subsequently referenced by name in methods where a D6708 assessment was performed, before the change in terminology in this standard practice was adopted.

NOTE 3—Users are cautioned against applying the between methods reproducibility as calculated from this practice to materials that are significantly different in composition from those actually studied, as the ability of this practice to detect and address sample-specific biases (see 6.8) is dependent on the materials selected for the interlaboratory study. When sample-specific biases are present, the types and ranges of samples may need to be expanded significantly from the minimum of ten as specified in this practice in order to obtain a more comprehensive and reliable 95 % confidence limits for between methods reproducibility that adequately cover the range of sample specific biases for different types of materials.

1.6 This practice is intended for test methods which measure quantitative (numerical) properties of petroleum or petroleum products.

1.7 The statistical methodology outlined in this practice is also applicable for assessing the expected agreement between any two test methods that purport to measure the same property of a material, provided the results are obtained on the same comparison sample set, the standard error associated with each test result is known, the sample set design meets the requirement of this practice, and the statistical degree of freedom of the data set exceeds 30.

## 2. Referenced Documents

### 2.1 ASTM Standards:<sup>2</sup>

D5580 Test Method for Determination of Benzene, Toluene, Ethylbenzene, *p/m*-Xylene, *o*-Xylene, C<sub>9</sub>, and Heavier

<sup>1</sup> This practice is under the jurisdiction of ASTM Committee D02 on Petroleum Products, Liquid Fuels, and Lubricants and is the direct responsibility of Subcommittee D02.94 on Coordinating Subcommittee on Quality Assurance and Statistics.

Current edition approved Jan. 1, 2016. Published February 2016. Originally approved in 2001. Last previous edition approved in 2015 as D6708 – 15. DOI: 10.1520/D6708-16.

<sup>2</sup> For referenced ASTM standards, visit the ASTM website, www.astm.org, or contact ASTM Customer Service at service@astm.org. For *Annual Book of ASTM Standards* volume information, refer to the standard's Document Summary page on the ASTM website.

\*A Summary of Changes section appears at the end of this standard

Aromatics, and Total Aromatics in Finished Gasoline by Gas Chromatography

D5769 Test Method for Determination of Benzene, Toluene, and Total Aromatics in Finished Gasolines by Gas Chromatography/Mass Spectrometry

D6299 Practice for Applying Statistical Quality Assurance and Control Charting Techniques to Evaluate Analytical Measurement System Performance

D6300 Practice for Determination of Precision and Bias Data for Use in Test Methods for Petroleum Products and Lubricants

D7372 Guide for Analysis and Interpretation of Proficiency Test Program Results

2.2 ISO Standard:<sup>3</sup>

ISO 4259 Petroleum Products—Determination and application of precision data in relation to methods of test.

### 3. Terminology

#### 3.1 Definitions:

3.1.1 *between-method bias, n*—a quantitative expression for the mathematical correction that can statistically improve the degree of agreement between the expected values of two test methods which purport to measure the same property.

3.1.2 *between methods reproducibility (R<sub>XY</sub>)*, *n*—a quantitative expression of the random error associated with the difference between two results obtained by different operators using different apparatus and applying the two methods X and Y, respectively, each obtaining a single result on an identical test sample, when the methods have been assessed and an appropriate bias-correction has been applied in accordance with this practice; it is defined as the 95 % confidence limit for the difference between two such single and independent results.

3.1.2.1 *Discussion*—A statement of between methods reproducibility must include a description of any bias correction used in accordance with this practice.

3.1.2.2 *Discussion*—Between methods reproducibility is a meaningful concept only if there are no statistically observable sample-specific relative biases between the two methods, or if such biases vary from one sample to another in such a way that they may be considered random effects. (see 6.7.)

3.1.3 *closeness sum of squares (CSS)*, *n*—a statistic used to quantify the degree of agreement between the results from two test methods after bias-correction using the methodology of this practice.

3.1.4 *total sum of squares (TSS)*, *n*—a statistic used to quantify the information content from the inter-laboratory study in terms of total variation of sample means relative to the standard error of each sample mean.

#### 3.2 Symbols:

X, Y = single X-method and Y-method results, respectively

$X_{ijk}$	$Y_{ijk}$	= single results from the X-method and Y-method round robins, respectively
$X_{i\bar{p}}$	$Y_i$	= means of results on the $i^{\text{th}}$ round robin sample
$S$		= the number of samples in the round robin
$L_{Xi\bar{p}}$	$L_{Yi}$	= the numbers of laboratories that returned results on the $i^{\text{th}}$ round robin sample
$R_X$	$R_Y$	= the reproducibilities of the X- and Y-methods, respectively
$S_{RXi\bar{p}}$	$S_{RYi}$	= the reproducibility standard deviations, evaluated at the means of the $i^{\text{th}}$ round robin sample
$S_{rXi\bar{p}}$	$S_{rYi}$	= the repeatability standard deviations, evaluated at the means of the $i^{\text{th}}$ round robin sample
$s_{Xi\bar{p}}$	$s_{Yi}$	= standard errors of the means $i^{\text{th}}$ round robin sample
$\bar{X}$	$\bar{Y}$	= the weighted means of round robins (across samples)
$x_{i\bar{p}}$	$y_i$	= deviations of the means of the $i^{\text{th}}$ round robin sample results from $\bar{X}$ and $\bar{Y}$ , respectively.
$TSS_X$	$TSS_Y$	= total sums of squares, around $\bar{X}$ and $\bar{Y}$
$F$		= a ratio for comparing variances; not unique—more than one use
$v_X$	$v_Y$	= the degrees of freedom for reproducibility variances from the round robins
$w_i$		= weight associated with the difference between mean results (or corrected mean results) from the $i^{\text{th}}$ round robin sample
$CSS$		= weighted sum of squared differences between (possibly corrected) mean results from the round robin
$a, b$		= parameters of a linear correction: $\hat{Y} = a + bX$
$t_1, t_2$		= ratios for assessing reductions in sums of squares
$R_{XY}$	$\hat{Y}$	= estimate of between methods reproducibility
$\hat{Y}_i$		= Y-method value predicted from X-method result
$\hat{Y}_i$		= $i^{\text{th}}$ round robin sample Y-method mean, predicted from corresponding X-method mean
$\varepsilon_i$		= standardized difference between $Y_i$ and $\hat{Y}_i$ .
$L_X$	$L_Y$	= harmonic mean numbers of laboratories submitting results on round robin samples, by X- and Y- methods, respectively
$R_X$	$\hat{Y}$	= estimate of between methods reproducibility, computed from an X-method result only

### 4. Summary of Practice

4.1 Precisions of the two methods are quantified using inter-laboratory studies meeting the requirements of Practice D6300 or equivalent, using at least ten samples in common that span the intersecting scopes of the methods. The arithmetic means of the results for each common sample obtained by each method are calculated. Estimates of the standard errors of these means are computed.

NOTE 4—For established standard test methods, new precision studies generally will be required in order to meet the common sample requirement.

NOTE 5—Both test methods do not need to be run by the same laboratory. If they are, care should be taken to ensure the independent test

<sup>3</sup> Available from American National Standards Institute (ANSI), 25 W. 43rd St., 4th Floor, New York, NY 10036.

result requirement of Practice **D6300** is met (for example, by double-blind testing of samples in random order).

4.2 Weighted sums of squares are computed for the total variation of the mean results across all common samples for each method. These sums of squares are assessed against the standard errors of the mean results for each method to ensure that the samples are sufficiently varied before continuing with the practice.

4.3 The closeness of agreement of the mean results by each method is evaluated using appropriate weighted sums of squared differences. Such sums of squares are computed from the data first with no bias correction, then with a constant bias correction, then, when appropriate, with a proportional correction, and finally with a linear (proportional + constant) correction.

4.4 The weighted sums of squared differences for the linear correction is assessed against the total variation in the mean results for both methods to ensure that there is sufficient correlation between the two methods.

4.5 The most parsimonious bias correction is selected.

4.6 The weighted sum of squares of differences, after applying the selected bias correction, is assessed to determine whether additional unexplained sources of variation remain in the residual (that is, the individual  $Y_i$  minus bias-corrected  $X_i$ ) data. Any remaining, unexplained variation is attributed to sample-specific biases (also known as method-material interactions, or matrix effects). In the absence of sample-specific biases, the between methods reproducibility is estimated.

4.7 If sample-specific biases are present, the residuals (that is, the individual  $Y_i$  minus *bias-corrected*  $X_i$ ) are tested for randomness. If they are found to be consistent with a random-effects model, then their contribution to the between methods reproducibility is estimated, and accumulated into an all-encompassing between methods reproducibility estimate.

4.8 Refer to **Fig. 1** for a simplified flow diagram of the process described in this practice.

## 5. Significance and Use

5.1 This practice can be used to determine if a constant, proportional, or linear bias correction can improve the degree of agreement between two methods that purport to measure the same property of a material.

5.2 The bias correction developed in this practice can be applied to a single result ( $X$ ) obtained from one test method (method  $X$ ) to obtain a *predicted* result ( $\hat{Y}$ ) for the other test method (method  $Y$ ).

NOTE 6—Users are cautioned to ensure that  $\hat{Y}$  is within the scope of method  $Y$  before its use.

5.3 The between methods reproducibility established by this practice can be used to construct an interval around  $\hat{Y}$  that would contain the result of test method  $Y$ , if it were conducted, with about 95 % confidence.

5.4 This practice can be used to guide commercial agreements and product disposition decisions involving test methods that have been evaluated relative to each other in accordance with this practice.

5.5 The magnitude of a statistically detectable bias is directly related to the uncertainties of the statistics from the experimental study. These uncertainties are related to both the size of the data set and the precision of the processes being studied. A large data set, or, highly precise test method(s), or both, can reduce the uncertainties of experimental statistics to the point where the “statistically detectable” bias can become “trivially small,” or be considered of no practical consequence in the intended use of the test method under study. Therefore, users of this practice are advised to determine in advance as to the magnitude of bias correction below which they would consider it to be unnecessary, or, of no practical concern for the intended application prior to execution of this practice.

NOTE 7—It should be noted that the determination of this minimum bias of no practical concern is not a statistical decision, but rather, a subjective decision that is directly dependent on the application requirements of the users.

## 6. Procedure

NOTE 8—For an in-depth statistical discussion of the methodology used in this section, see **Appendix X1**. For a worked example, see **Appendix X2**.

6.1 Calculate sample means and standard errors from Practice **D6300** results.

6.1.1 The process of applying Practice **D6300** to the data may involve elimination of some results as outliers, and it may also involve applying a transformation to the data. For this practice, compute the mean results from data that have not been transformed, but with outliers removed in accordance with Practice **D6300**. The precision estimates from Practice **D6300** are used to estimate the standard errors of these means.

6.1.2 Compute the means as follows:

6.1.2.1 Let  $X_{ijk}$  represent the  $k^{\text{th}}$  result on the  $i^{\text{th}}$  common material by the  $j^{\text{th}}$  lab in the round robin for method  $X$ . Similarly for  $Y_{ijk}$ . (The  $i^{\text{th}}$  material is the same for both round robins, but the  $j^{\text{th}}$  lab in one round robin is not necessarily the same lab as the  $j^{\text{th}}$  lab in the other round robin.) Let  $n_{Xij}$  be the number of results on the  $i^{\text{th}}$  material from the  $j^{\text{th}}$   $X$ -method lab, after removing outliers that is, the number of results in *cell* ( $i, j$ ). Let  $L_{Xi}$  be the number of laboratories in the  $X$ -method round robin that have at least one result on the  $i^{\text{th}}$  material remaining in the data set, after removal of outliers. Let  $S$  be the total number of materials common to both round robins.

6.1.2.2 The mean  $X$ -method result for the  $i^{\text{th}}$  material is:

$$X_i = \frac{1}{L_{Xi}} \sum_j \frac{\sum_k X_{ijk}}{n_{Xij}} \quad (1)$$

where,  $X_i$  is the average of the cell averages on the  $i^{\text{th}}$  material by method  $X$ .

6.1.2.3 Similarly, the mean  $Y$ -method result for the  $i^{\text{th}}$  material is:

$$Y_i = \frac{1}{L_{Yi}} \sum_j \frac{\sum_k Y_{ijk}}{n_{Yij}} \quad (2)$$

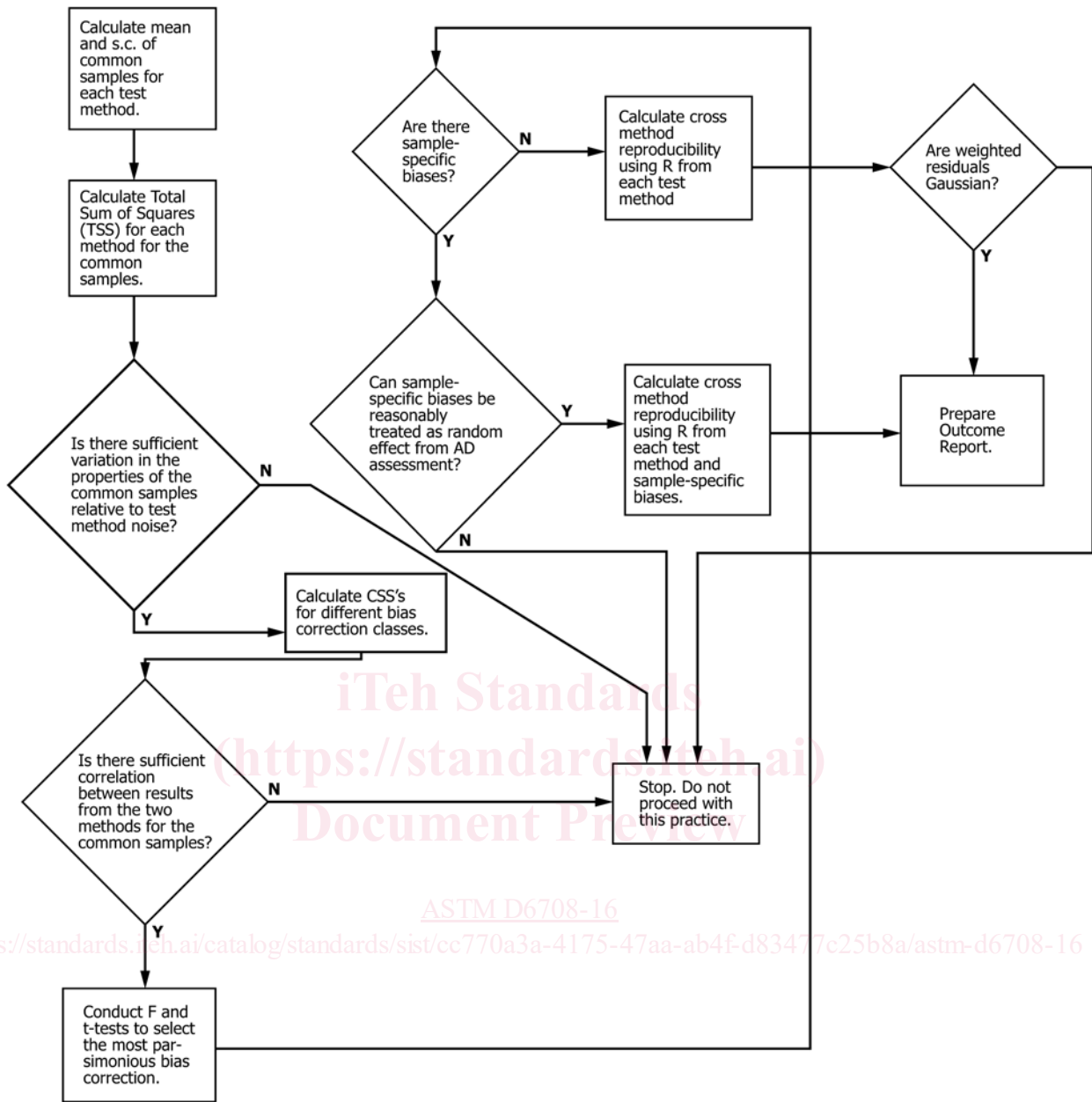


FIG. 1 Simplified Flow Diagram for this Practice

6.1.3 The standard errors (standard deviations of the means of the results) are computed as follows:

6.1.3.1 If  $s_{RX_i}$  is the estimated reproducibility standard deviation from the X-method round robin, and  $s_{rXi}$  is the estimated repeatability standard deviation, then an estimate of the standard error for  $X_i$  is given by:

$$s_{Xi} = \sqrt{\frac{1}{L_{Xi}} \left[ s_{RX_i}^2 - s_{rXi}^2 \left( 1 - \frac{1}{L_{Xi}} \sum_j \frac{1}{n_{Xij}} \right) \right]} \quad (3)$$

NOTE 9—Since repeatability and reproducibility may vary with X, even if the  $L_{Xi}$  were the same for all materials and the  $n_{Xij}$  were the same for all laboratories and all materials, the  $\{s_{Xi}\}$  might still differ from one material to the next.

6.1.3.2  $s_{Y_i}$ , the estimated standard error for  $Y_i$ , is given by an analogous formula.

6.2 Calculate the total variation sum of squares for each method, and determine whether the samples can be distinguished from each other by both methods.

6.2.1 The total sums of squares (TSS) are given by:

$$TSS_x = \sum_i \left( \frac{X_i - \bar{X}}{s_{Xi}} \right)^2 \quad \text{and} \quad TSS_y = \sum_i \left( \frac{Y_i - \bar{Y}}{s_{Yi}} \right)^2 \quad (4)$$

where:

$$\bar{X} = \frac{\sum_i \left( \frac{X_i}{s_{Xi}^2} \right)}{\sum_i \left( \frac{1}{s_{Xi}^2} \right)} \quad \text{and} \quad \bar{Y} = \frac{\sum_i \left( \frac{Y_i}{s_{Yi}^2} \right)}{\sum_i \left( \frac{1}{s_{Yi}^2} \right)} \quad (5)$$

are weighted averages of all  $X_i$ 's and  $Y_i$ 's respectively.

6.2.2 Compare  $F = TSS_X/(S-1)$  to the 95<sup>th</sup> percentile of Fisher's  $F$  distribution with  $(S-1)$  and  $\nu_x$  degrees of freedom for the numerator and denominator, respectively, where  $\nu_x$  is the degrees of freedom for the reproducibility variance (Practice D6300, paragraph 8.3.3.3) for the X-method round robin. If  $F$  does not exceed the 95<sup>th</sup> percentile, then the X-method is not sufficiently precise to distinguish among the  $S$  samples. Do not proceed with this practice, as meaningful results cannot be produced.

6.2.3 In a similar manner, compare  $F = TSS_Y/(S-1)$  to the 95<sup>th</sup> percentile of Fisher's  $F$  distribution, using the degrees of freedom of the reproducibility variance of the Y-method,  $\nu_y$ , in place of  $\nu_x$ . Similarly, do not proceed with this practice if  $F$  does not exceed the 95<sup>th</sup> percentile.

NOTE 10—If one or both of the conditions of 6.2.2 and 6.2.3 are satisfied only marginally, it is unlikely that this practice will produce meaningful results since in 6.4, the quantity  $(TSS_X + TSS_Y)$  will be compared to a closeness sum of squares computed in the next section, to determine whether the methods are sufficiently correlated. It will be difficult to meet that correlation requirement if the samples are too similar to one another.

6.3 Calculate the closeness sum of squares (CSS) statistic for each of the following classes of bias-correction methodology.

6.3.1 Class 0—No bias correction.

6.3.1.1 Compute the weights ( $w_i$ ) for each sample  $i$ :

$$w_i = \frac{1}{s_{Yi}^2 + s_{Xi}^2} \quad (6)$$

6.3.1.2 Compute CSS:

$$CSS_0 = \sum w_i (X_i - Y_i)^2 \quad (7)$$

6.3.2 Class 1a—Constant bias correction.

6.3.2.1 Using the weights ( $w_i$ ) from 6.3.1.1, compute the constant bias correction ( $a$ ):

$$a = \frac{\sum w_i (Y_i - X_i)}{\sum w_i} = \frac{\sum w_i Y_i}{\sum w_i} - \frac{\sum w_i X_i}{\sum w_i} \quad (8)$$

6.3.2.2 Compute CSS:

$$CSS_{1a} = \sum w_i (Y_i - (X_i + a))^2 \quad (9)$$

6.3.3 Class 1b—Proportional bias correction.

6.3.3.1 The computations of this subsection (6.3.3) are appropriate only if both of the following conditions apply: (1) the measured property assumes only non-negative values, and (2) a property value of zero has a physical significance (for example, concentrations of specific constituents). In addition, it is not mandatory but highly recommended that  $\max(Y_i) \geq 2 \min(Y_i)$ .

6.3.3.2 The computations involve iterative calculation of the weights ( $w_i$ ) and the proportional correction ( $b$ ).

6.3.3.3 Set  $b = 1$ .

6.3.3.4 Compute the weights ( $w_i$ ) for each sample  $i$ :

$$w_i = \frac{1}{S_{Yi}^2 + b^2 S_{Xi}^2} \quad (10)$$

6.3.3.5 Calculate  $b_0$ :

$$b_0 = \frac{\sum w_i X_i Y_i}{\sum w_i X_i^2 - \sum w_i^2 S_{Xi}^2 (Y_i - b X_i)^2} \quad (11)$$

6.3.3.6 If  $|b - b_0| > .001 b$ , replace  $b$  with  $b_0$  and go back to 6.3.3.4. Otherwise, the iteration can be stopped, as further iteration will not produce meaningful improvement. Replace  $b$  with  $b_0$  and go on to 6.3.3.7.

6.3.3.7 Calculate  $CSS_{1b}$ :

$$CSS_{1b} = \sum w_i (Y_i - b X_i)^2 \quad (12)$$

6.3.4 Class 2—Linear (proportional + constant) bias correction.

6.3.4.1 This involves iterative calculation of the weights ( $w_i$ ), the weighted means of  $X_i$ 's and  $Y_i$ 's, and the proportional term ( $b$ ).

6.3.4.2 Set  $b = 1$ .

6.3.4.3 Compute the weights ( $w_i$ ) for each sample  $i$ :

$$w_i = \frac{1}{s_{Yi}^2 + b^2 s_{Xi}^2} \quad (13)$$

6.3.4.4 Calculate the weighted means of  $\{X_i\}$  and  $\{Y_i\}$  respectively:

$$\bar{X} = \frac{\sum w_i X_i}{\sum w_i} \quad (14)$$

$$\bar{Y} = \frac{\sum w_i Y_i}{\sum w_i}$$

6.3.4.5 Calculate the deviations from the weighted means:

$$x_i = X_i - \bar{X} \quad (15)$$

$$y_i = Y_i - \bar{Y}$$

6.3.4.6 Calculate  $b_0$ :

$$b_0 = \frac{\sum w_i x_i y_i}{\sum w_i x_i^2 - \sum w_i^2 S_{Xi}^2 (y_i - b x_i)^2} \quad (16)$$

6.3.4.7 If  $|b - b_0| > .001 b$ , replace  $b$  with  $b_0$  and go back to 6.3.4.3, computing new values for the weights  $\{w_i\}$ ,  $\bar{X}$ ,  $\bar{Y}$ ,  $\{x_i\}$ ,  $\{y_i\}$ , and  $b_0$ . Otherwise, the iteration can be stopped, as further iteration will not produce meaningful improvement. Replace  $b$  with  $b_0$  and go to 6.3.4.8.

6.3.4.8 Calculate  $CSS_2$  and  $a$ :

$$CSS_2 = \sum w_i (y_i - b x_i)^2 \quad (17)$$

$$a = \bar{Y} - b \bar{X} \quad (18)$$

6.4 Test whether the methods are sufficiently correlated.

6.4.1 Calculate the  $F$ -statistic:

$$F = \frac{(TSS_X + TSS_Y - CSS_2)/S}{CSS_2/(S - 2)} \quad (19)$$

6.4.2 Compare  $F$  to the 95<sup>th</sup> percentile of Fisher's  $F$  distribution with  $S$  and  $S-2$  degrees of freedom in the numerator and denominator, respectively.

6.4.2.1 If  $F$  is less than the 95<sup>th</sup> percentile value, then, this practice concludes that the methods are too discordant to permit use of the results from one method to predict those of the other.

6.4.2.2 If  $F$  is greater than the tabled value, proceed to 6.5.

6.5 Conduct tests to select the most parsimonious bias correction class needed.

6.5.1 The closeness sums of squares for differences from each class of bias correction are used to select the most parsimonious bias correction class that can improve the expected degree of agreement between the  $\hat{Y}$  (the predicted Y-method result using X-method result) and the actual Y-method result on the same material. The classes of bias correction and the associated CSS as calculated earlier are repeated in the following table.

Bias Correction Class	CSS
Class 0—no correction	$CSS_0$
Class 1a—constant bias correction	$CSS_{1a}$
Class 1b—proportional bias correction (when appropriate)	$CSS_{1b}$
Class 2—linear (proportional + constant bias correction)	$CSS_2$

6.5.2 To determine whether any bias correction (Classes 1a, 1b or 2 above) can significantly improve the expected agreement between the two methods, calculate the following ratio:

$$F = \frac{(CSS_0 - CSS_2)/2}{CSS_2/(S - 2)} \quad (20)$$

6.5.2.1 Compare  $F$  to the upper 95th percentile of the  $F$  distribution with 2 and  $S-2$  degrees of freedom for the numerator and denominator, respectively.

6.5.2.2 If the calculated  $F$  is smaller, conclude that a bias correction of Class 1a, 1b, or 2 does not sufficiently improve the expected agreement between the two methods, relative to Class 0 (no bias correction). Proceed to 6.6.

6.5.2.3 If the calculated  $F$  is larger, conclude that a correction can improve the expected agreement between the two methods, and continue in 6.5.3.

6.5.3 If the  $F$ -value calculated in 6.5.2 is larger than the 95th percentile of  $F$ , compute the following  $t$ -ratios:

$$t_1 = \sqrt{\frac{CSS_0 - CSS_1}{CSS_2/(S - 2)}} \quad (21)$$

$$t_2 = \sqrt{\frac{CSS_1 - CSS_2}{CSS_2/(S - 2)}}$$

where,  $CSS_1$  is the lesser of  $CSS_{1a}$  or  $CSS_{1b}$ , provided the latter is appropriate and has been calculated.

6.5.3.1 Compare  $t_2$  to the upper 97.5th percentile of the  $t$  distribution with  $S-2$  degrees of freedom.

6.5.3.2 If  $t_2$  is larger, conclude that a bias correction of Class 2 (proportional + constant correction) can improve the expected agreement over that of a single term (constant or proportional) correction alone (Class 1). Proceed to 6.6.

6.5.3.3 If  $t_2$  is smaller than the  $t$ -percentile, compare  $t_1$  to the same upper 97.5th percentile of the  $t$  distribution with ( $S-2$ ) degrees of freedom.

6.5.3.4 If  $t_1$  is larger, conclude that a single term bias correction of Class 1 is preferred to a bias correction of Class 2. Use the constant correction unless  $CSS_{1b}$  is appropriate and is smaller than  $CSS_{1a}$ . Proceed to 6.6.

6.5.3.5 If  $t_1$  is smaller, then neither  $t_1$  nor  $t_2$  is statistically significant. A bias correction of Class 2 is preferred over single-term (constant or proportional) correction of Class 1.

6.6 Test for existence of sample-specific biases.

6.6.1 Compare the CSS of the bias-correction class selected in 6.5 to the 95th percentile value of a chi-square distribution with  $\nu$  degrees of freedom

where:

$\nu = S$  for Class 0 (-no bias) correction,

$\nu = S - 1$  for Class 1a or Class 1b (constant or proportional) correction

$\nu = S - 2$  for Class 2 (linear) correction

6.6.2 If the CSS is smaller than the chi-square percentile, it is reasonable to conclude that there are no sample-specific biases, that is, that there are no other sources of variation that are statistically observable above the measurement error. Perform the Anderson-Darling (A-D) assessment on the residuals as per 6.7.2.2 and 6.7.2.3. If the outcome is not significant at the 5 % level, calculate the between methods reproducibility ( $R_{XY}$ ) as per Eq 22 below. If the A-D assessment is significant, application of the practice is considered terminated with failure at this point, as the statistical evidence suggests that a single between-method reproducibility ( $R_{XY}$ ) cannot be found that is applicable to all materials covered by the intersecting scope of both test methods. It is reasonable to conclude that, at least for some materials, the test methods are not measuring the same property.

$$R_{XY} = \sqrt{\frac{R_Y^2 + b^2 R_X^2}{2}} \quad (22)$$

where:

$b =$  the coefficient of the appropriate bias correction. (For Class 0 and Class 1a bias corrections,  $b=1$ .)

6.6.3 If the CSS is larger than the chi-square percentile (see 6.6.1), there is strong evidence that biases between the methods have not been adequately corrected by the bias-corrections of 6.3. In other words, the relative biases are not consistent across the  $S$  common samples of the round robins. The user may wish to investigate whether the biases can be attributed to other observable properties of the samples. Or he or she may wish to restrict attention to a smaller class of materials for the purpose of establishing a between methods reproducibility. Such investigations are beyond the scope of this practice, as the issues typically are not statistical in nature. This practice does recommend investigating whether it is reasonable to treat the sample-specific biases as random effects, as described in 6.7.

6.7 Treatment of Sample-Specific Relative Bias as a Variance Component:

6.7.1 If the CSS exceeds the 95th percentile value of the appropriate chi-square distribution (see 6.6.1), there is strong evidence that sources other than measurement error are contributing towards the variation of the expected agreement between the two methods. In this practice, these sources are attributed to sample-specific effects (also known as matrix effects or method-material interactions). In some cases these sample-specific effects can be treated as random effects, and hence can be incorporated as an additional source of variation into a between methods reproducibility as described in this section. Note that, even when it is appropriate to treat these sample-specific effects as random, the additional variation may