

OPEN

Corruption of the Pearson correlation coefficient by measurement error and its estimation, bias, and correction under different error models

Edoardo Saccenti^{1,4*}, Margriet H. W. B. Hendriks² & Age K. Smilde^{3,4}

Correlation coefficients are abundantly used in the life sciences. Their use can be limited to simple exploratory analysis or to construct association networks for visualization but they are also basic ingredients for sophisticated multivariate data analysis methods. It is therefore important to have reliable estimates for correlation coefficients. In modern life sciences, comprehensive measurement techniques are used to measure metabolites, proteins, gene-expressions and other types of data. All these measurement techniques have errors. Whereas in the old days, with simple measurements, the errors were also simple, that is not the case anymore. Errors are heterogeneous, non-constant and not independent. This hampers the quality of the estimated correlation coefficients seriously. We will discuss the different types of errors as present in modern comprehensive life science data and show with theory, simulations and real-life data how these affect the correlation coefficients. We will briefly discuss ways to improve the estimation of such coefficients.

The concept of correlation and correlation coefficient dates back to Bravais¹ and Galton² and found its modern formulation in the work of Fisher and Pearson^{3,4}, whose product moment correlation coefficient ρ has become the most used measure to describe the linear dependence between two random variables. From the pioneering work of Galton on heredity, the use of correlation (or co-relation as it was termed) spread virtually in all fields of research and results based on it pervade the scientific literature.

Correlations are generally used to quantify, visualize and interpret bivariate (linear) relationships among measured variables. They are the building blocks of virtually all multivariate methods such as Principal Component Analysis (PCA⁵⁻⁷), Partial Least Squares regression, Canonical Correlation Analysis (CCA⁸) which are used to reduce, analyze and interpret high-dimensional *omics* data sets and are often the starting point for the inference of biological networks such as metabolite-metabolite associations networks^{9,10}, gene regulatory networks^{11,12} and co-expression networks^{13,14}.

Fundamentally, correlation and correlation analysis are pivotal for understanding biological systems and the physical world. With the increase of comprehensive measurements (liquid-chromatography mass-spectrometry, nuclear magnetic resonance (NMR), gas-chromatography mass-spectrometry (MS) in metabolomics and proteomics; RNA-sequencing in transcriptomics) in life sciences, correlations are used as a first tool for visualization and interpretation, possibly after selection of a threshold to filter the correlations. However, the complexity and the difficulty of estimating correlation coefficients is not fully acknowledged.

Measurement error is intrinsic to every experimental technique and measurement platform, be it a simple ruler, a gene sequencer or a complicated array of detectors in a high-energy physics experiment, and already in the early days of statistics it was known that measurement error can bias the estimation of correlations¹⁵. This bias was first called attenuation because it was found that under the error condition considered, the correlation

¹Laboratory of Systems and Synthetic Biology, Wageningen University & Research, Wageningen, The Netherlands.

²DSM Biotechnology Center, Delft, The Netherlands. ³Biosystems Data Analysis, Swammerdam Institute for Life Sciences, University of Amsterdam, Amsterdam, The Netherlands. ⁴These authors contributed equally: Edoardo Saccenti and Age K. Smilde. *email: edoardo.saccenti@wur.nl

was attenuated towards zero. The attenuation bias has been known and discussed in some research fields^{16–19} but it seems to be totally neglected in modern *omics*-based science. Moreover, contemporary comprehensive *omics* measurement techniques have far more complex measurement error structures than the simple ones considered in the past on which early results were based.

In this paper, we intend to show the impact of measurement errors on the quality of the calculated correlation coefficients and we do this for several reasons. First, to make the *omics* community aware of the problems. Secondly, to make the theory of correlation up to date with current *omics* measurements taking into account more realistic measurement error models in the calculation of the correlation coefficient and third, to propose ways to alleviate the problem of distortion in the estimation of correlation induced by measurement error. We will do this by deriving analytical expressions supported by simulations and simple illustrations. We will also use real-life metabolomics data to illustrate our findings.

Measurement Error Models

We start with the simple case of having two correlated biological entities x_0 and y_0 which are randomly varying in a population. This may, *e.g.*, be concentrations of two blood metabolites in a cohort of persons or gene-expressions of two genes in cancer tissues. We will assume that these variables are normally distributed

$$\begin{pmatrix} x_0 \\ y_0 \end{pmatrix} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}_0) \quad (1)$$

with underlying mean

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_{x_0} \\ \mu_{y_0} \end{pmatrix} \quad (2)$$

and variance-covariance matrix

$$\boldsymbol{\Sigma}_0 = \begin{pmatrix} \sigma_{x_0}^2 & \sigma_{x_0 y_0} \\ \sigma_{x_0 y_0} & \sigma_{y_0}^2 \end{pmatrix}. \quad (3)$$

Under this model the variance components $\sigma_{x_0}^2$ and $\sigma_{y_0}^2$ describe the biological variability for x_0 and y_0 , respectively. The correlation ρ_0 , between x_0 and y_0 is given by

$$\rho_0 = \frac{\sigma_{x_0 y_0}}{\sqrt{\sigma_{x_0}^2 \sigma_{y_0}^2}}. \quad (4)$$

We refer to ρ_0 as the *true correlation*.

Whatever the nature of the variables x_0 and y_0 and whatever the experimental technique used to measure them there is always a random error component (also referred to as noise or uncertainty) associated with the measurement procedure. This random error is by its own nature not reproducible (in contrast with systematic error which is reproducible and can be corrected for) but can be modeled, *i.e.* described, in a statistical fashion. Such models have been developed and applied in virtually every area of science and technology and can be used to adjust for measurement errors or to describe the bias introduced by it. The measured variables will be indicated by x and y to distinguish them from x_0 and y_0 which are their errorless counterparts.

The correlation coefficient ρ_0 is sought to be estimated from these measured data. Assuming that N samples are taken, the sample correlation r_N is calculated as

$$r_N = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N s_x s_y}, \quad (5)$$

where (\bar{x}, \bar{y}) is the sample mean over N observations and s_x, s_y are the usual sample standard deviation estimators. This sample correlation is used as a proxy of ρ_0 . The population value of this sample correlation is

$$\rho = \frac{E[xy] - E[x]E[y]}{\sqrt{E[x^2] - E[x]^2} \sqrt{E[y^2] - E[y]^2}}, \quad (6)$$

and it also holds that

$$\lim_{N \rightarrow \infty} r_N = \rho. \quad (7)$$

We will call ρ the *expected correlation*. Ideally, $\rho_0 = \rho$ but this is unfortunately not always the case. In plain words: certain measurement errors do not cancel out if the number of samples increases.

In the following section we will introduce three error models and will show with both simulated and real data how measurement error impacts the estimation of the Pearson correlation coefficient. We will focus mainly on ρ_0 and ρ .

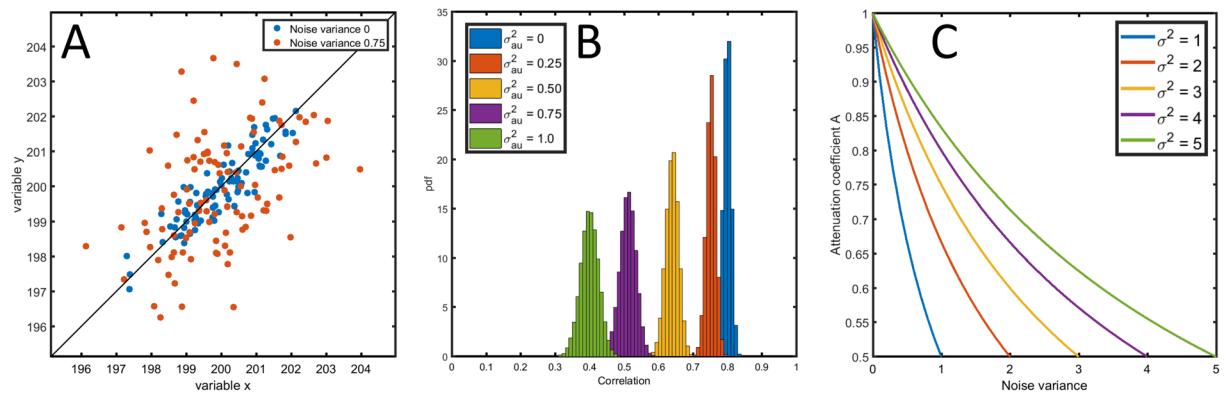


Figure 1. (A) Correlation plot of two variables x and y ($\sigma_{x_0}^2 = \sigma_{y_0}^2 = 1$) generated without ($\sigma_{au_x}^2 = \sigma_{au_y}^2 = 0$) and with uncorrelated additive error ($\sigma_{au_x}^2 = \sigma_{au_y}^2 = 0.75$) with underlying true correlation $\rho_0 = 0.8$ (model 8). (B) Distribution of the sample correlation coefficient for different levels of measurement error ($\sigma_{au_x}^2 = \sigma_{au_y}^2 = \sigma_{au}^2$) for a true correlation $\rho_0 = 0.8$. (C) The attenuation coefficient A from Eq. (10) as a function the measurement error for different level of the variance $\sigma^2 = \sigma_{x_0}^2 = \sigma_{y_0}^2$ of the variables x_0 and y_0 . See Material and Methods section 6.5.1 for details on the simulations.

Additive error. The most simple error model is the additive error model where the measured entities x and y are modeled as

$$\begin{cases} x = x_0 + \varepsilon_{au_x} \\ y = y_0 + \varepsilon_{au_y} \end{cases} \quad (8)$$

where it is assumed that the error components ε_{au_x} and ε_{au_y} are independently normally distributed around zero with variance $\sigma_{au_x}^2$ and $\sigma_{au_y}^2$ and are also independent from x_0 and y_0 . The subscripts au_x , au_y stand for *additive uncorrelated error* (ε) on variables x and y .

Variables x and y represent measured quantities accessible to the experimenter. This error model describes the case in which the measurement error causes within-sample variability, which means that p measurement replicates $x_{i,1}, x_{i,2}, \dots, x_{i,p}$ of observation x_i of variable x will all have slightly different values due to the random fluctuation of the error component ε_{au_x} ; the extent of the variability among the replicates depends on the magnitude of the error variance $\sigma_{au_x}^2$ (and similarly for the y variable). This can be seen in Fig. 1A where it is shown that in the presence of measurement error (i.e. $\sigma_{au_x}^2, \sigma_{au_y}^2 > 0$) the two variables x and y are more dispersed. Due to the measurement error, the expected correlation coefficient ρ is always biased downwards, i.e. $\rho < \rho_0$, as already shown by Spearman¹⁵ (see Fig. 1B) who also provided an analytical expression for the attenuation of the expected correlation coefficient as a function of the error components (a modern treatment can be found in reference²⁰):

$$\rho = A\rho_0, \quad (9)$$

where

$$A = \frac{1}{\sqrt{\left(1 + \frac{\sigma_{au_x}^2}{\sigma_{x_0}^2}\right)\left(1 + \frac{\sigma_{au_y}^2}{\sigma_{y_0}^2}\right)}}. \quad (10)$$

Equation (9) implies that in presence of measurement error the expected correlation is different from the true correlation ρ_0 which is sought to be estimated. The attenuation A is always strictly smaller than 1 and it is a decreasing function of the size of the measurement error relative to the biological variation (see Fig. 1C), as it can be seen from Eq. (10). The attenuation of the expected correlation, despite being known since 1904, has sporadically resurfaced in the statistical literature in the psychological, epidemiology and behavioral sciences (where it is known as attenuation due to intra-person or intra-individual variability, see¹⁹ and reference therein) but has been largely neglected in the life sciences, despite its relevance.

The error model (8) can be extended to include a correlated error term ε_{ac}

$$\begin{cases} x = x_0 + \varepsilon_{au_x} + \varepsilon_{ac} \\ y = y_0 + \varepsilon_{au_y} \pm \varepsilon_{ac} \end{cases} \quad (11)$$

with ε_{ac} normally distributed around zero with variance σ_{ac}^2 ; the correlated error term takes on exactly the same value for x and y in a given sample. The ‘ \pm ’ models the sign of the error correlation. When ε_{ac} has a positive sign

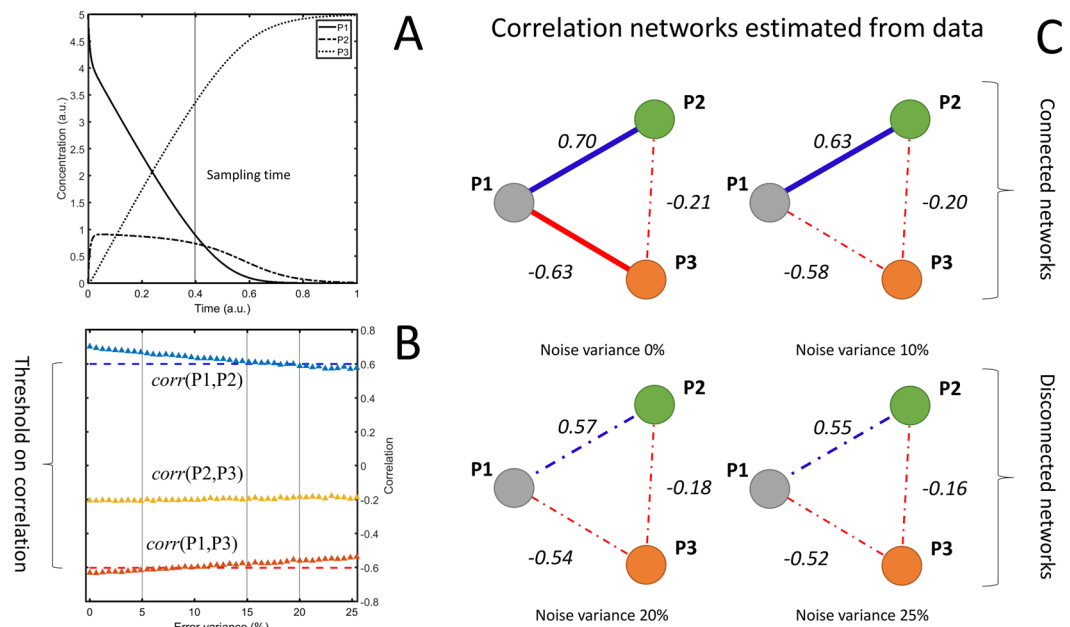


Figure 2. Consequences of measurement error when using correlation in systems biology. **(A)** Time concentration profile of three metabolites P1, P2 and P3 generated through a simple enzymatic metabolic model; 100 profiles are generated by randomly varying the kinetic parameters defining the model and sampled at time 0.4 (a.u.). **(B)** Average pairwise correlation of P1, P2 and P3 as a function of the variance of the additive uncorrelated error. **(C)** Inference of a metabolite-metabolite correlation network: two metabolites are associated if their correlation is above 0.6²³ (see threshold in **B**). The increasing level of measurement error hampers the network inference (compare the different panels). See Material and Methods section 6.5.2 for details on the simulations.

in both x and y the error is positively correlated; if the sign is discordant the error is negatively correlated. The subscript ac is used to indicate *additive correlated* error. The variance for x is given by

$$\sigma_x^2 = \sigma_{x_0}^2 + \sigma_{a_{u_x}}^2 + \sigma_{ac}^2 \quad (12)$$

and likewise for the variable y . In general, additive correlated error can have different causes depending on the type of instruments and measurement protocols used. For example, in transcriptomics, metabolomics and proteomics, usually samples have to be pretreated (sample work-up) prior to the actual instrumental analysis. Any error in a sample work-up step may affect all measured entities in a similar way²¹. Another example is the use of internal standards for quantification: any error in the amount of internal standard added may also affect all measured entities in a similar way. Hence, in both cases this leads to (positively) correlated measurement error. In some cases in metabolomics and proteomics the data are preprocessed using deconvolution tools. In that case two co-eluting peaks are mathematically separated and quantified. Since the total area under the curve is constant and (positive) error in one of the deconvoluted peaks is compensated by a (negative) error in the second peak, this may give rise to negatively correlated measurement error.

To show the effect of additive uncorrelated measurement error we consider the concentration profiles of three hypothetical metabolites P1, P2 and P3 simulated using a simple dynamic model (see Fig. 2A and Section 6.5.2) where additive uncorrelated measurement error is added before calculating the pairwise correlations among P1, P2 and P3: also in this case the magnitude of the correlation is attenuated, and the attenuation increases with the error variance (see Fig. 2B).

This has serious repercussions when correlations are used for the definition of association networks, as commonly done in systems biology and functional genomics^{10,22}: measurement error drives correlation towards zero and this impacts network reconstruction. If a threshold of 0.6 is imposed to discriminate between correlated and non correlated variables as usually done in metabolomics²³, an error variance of around 15% (see Fig. 2B, point where the correlation crosses the threshold) of the biological variation will attenuate the correlation to the point that metabolites will be deemed not to be associated even if they are biologically correlated leading to very different metabolite association networks (see Fig. 2C).

Multiplicative error. In many experimental situations it is observed that the measurement error is proportional to the magnitude of the measured signal; when this happens the measurement error is said to be multiplicative. The model for sampled variables in presence of multiplicative measurement error is

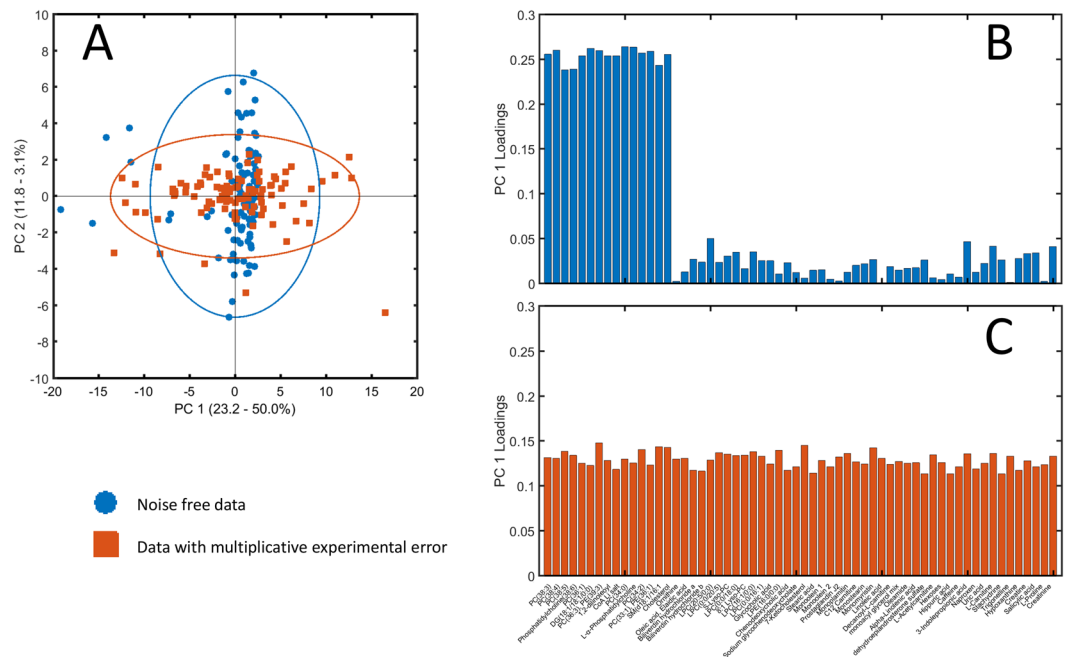


Figure 3. Consequences of multiplicative (correlated and uncorrelated) measurement error for data analysis. (A) Scatter plot of the overlaid view of the first two components of two PCA models of simulated data sets; one without multiplicative error and one with multiplicative error. For visualization purposes, the scores are plotted in the same graph, but the subspaces spanned by the first two principal components for the two data sets are of course different. The labels on both axes also present the percentage explained variation for the two analyses. (B) Loading plot for the error free data. (C) Loading plot for the data with multiplicative error. See Material and Methods section 6.5.3 for details on the simulations.

$$\begin{cases} x = x_0(1 + \varepsilon_{mu_x} + \varepsilon_{mc}) \\ y = y_0(1 + \varepsilon_{mu_y} \pm \varepsilon_{mc}) \end{cases} \quad (13)$$

where x_0 , y_0 , ε_{mu_x} , ε_{mu_y} and ε_{mc} have the same distributional properties as before in the additive error case, and the last three terms represent the *multiplicative uncorrelated errors* in x and y , respectively, and the *multiplicative correlated error*.

The characteristics of the multiplicative error and the variance of the measured entities σ_x^2 depend on the level μ_{x_0} of the signal to be measured (for a derivation of Eq. (14) see Section 6. 6.1.1):

$$\sigma_x^2 = \sigma_{x_0}^2 + \left(\sigma_{x_0}^2 + \mu_{x_0}^2 \right) \left(\sigma_{mu_x}^2 + \sigma_{mc}^2 \right), \quad (14)$$

while in the additive case the standard deviation is similar for different concentrations and does not depend explicitly on the signal intensity, as shown in Eq. (12). A similar equation holds for the variable y .

It has been observed that multiplicative errors often arises because of the different procedural steps like sample aliquoting²⁴; this is the case of deep sequencing experiments where the multiplicative error is possibly introduced by the pre-processing steps like, for example, linker ligation and PCR amplification which may vary from tag to tag and from sample to sample²⁵. In other cases the multiplicative error arises from the distributional properties of the signal, like in those experiments where the measurement comes down to counts like in the case of RNA fragments in an RNA-seq experiment or numbers of ions in a mass-spectrometer that are governed by Poisson distributions for which the standard deviation is equal to the mean. For another example, in NMR spectroscopy measured intensities are affected by the sample magnetization conditions: fluctuations in the external electromagnetic field or instability of the rf pulses affect the signal in a fashion that is proportional to the signal itself²⁶.

A multiplicative error distorts correlations and this affects the results of any data analysis approach which is based on correlations. To show the effect of multiplicative error we consider the analysis of a metabolomic data set simulated from real mass-spectrometry (MS) data, on which extra uncorrelated and correlated multiplicative measurement errors have been added. As it can be seen in Fig. 3A, the addition of error affects the underlying data structure: the error free data is such that only a subset of the measured variables contributes to explain the pattern in a low dimensional projection of the data, *i.e.* have PCA loadings substantially different from zero (3B). The addition of extra multiplicative error perturbs the loading structure to the point that all variables contribute equally to the model (3C), obscuring the real data structure and hampering the interpretation of the PCA model. This is not necessarily caused by the multiplicative nature of the error, but it is caused by the correlated error part. Since the term ε_{mc} is common to all variables, it introduces the same amount of correlation among all the varia-

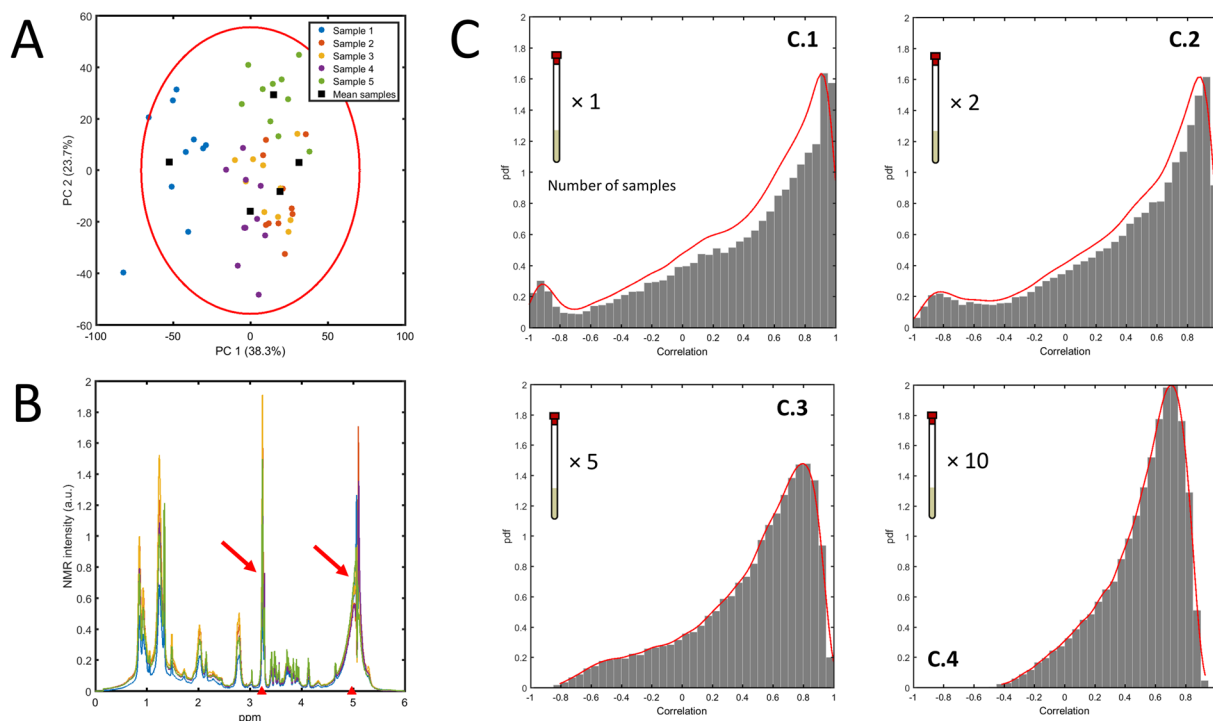


Figure 4. (A) PCA plot of 5 different samples of fish extract measured with technical replicates (10×) using NMR²⁹. (B) Overlap of the average binned NMR spectra of the 5 samples: the two resonances whose correlation is investigated are highlighted (3.23 and 4.98 ppm). (C) Distribution of the correlation coefficient between the two resonances calculated, taking as input the average over different numbers of technical replicates (see inserts). See Material and Methods section 6.5.4 for more details on the estimation procedure.

bles and this leads to all the variables contributing similarly to the latent vector (principal component). One may also observe that the variation explained by the first principal component increases when adding the correlated measurement error.

Realistic error. The measurement process usually consists of different procedural steps and each step can be viewed as a different source of measurement error with its own characteristics, which sum to both additive and multiplicative error components as is the case of comprehensive *omics* measurements²⁷. The model for this case is:

$$\begin{cases} x = x_0(1 + \varepsilon_{mu_x} + \varepsilon_{mc}) + \varepsilon_{au_x} + \varepsilon_{ac} \\ y = y_0(1 + \varepsilon_{mu_y} \pm \varepsilon_{mc}) + \varepsilon_{au_y} \pm \varepsilon_{ac} \end{cases} \quad (15)$$

where all errors have been introduced before and are all assumed to be independent of each other and independent of the true (biological) signals (x_0 and y_0).

This realistic error model has a multiplicative as well as an additive component and also accommodates correlated and uncorrelated error. It is an extension of a much-used error model for analytical chemical data which only contains uncorrelated error²⁸. From model (15) it follows that the error changes not only quantitatively but also qualitatively with changing signal intensity: the importance of the multiplicative component increases when the signal intensity increases, whereas the relative contribution of the additive error component increases when the signal decreases.

Since most of the measurements do not usually fall at the extremity of the dynamic range of the instruments used, the situation in which both additive and multiplicative error are important is realistic. For example, this is surely the case of comprehensive NMR and Mass Spectrometry measurements, where multiplicative errors are due to sample preparation and carry-over effect (in the case of MS) and the additive error is due to thermal error in the detectors²⁹. To illustrate this we consider an NMR experiment where a different number of technical replicates are measured for five samples (Fig. 4A,B). We are interested in establishing the correlation patterns across the (binned) resonances. For sake of simplicity we focus on two resonances, binned at 3.22 and 4.98 ppm. If one calculates the correlation using only one (randomly chosen) replicate per sample, the resulting correlation can be anywhere between -1 and 1 (see Fig. 4C.1). The variability reduces considerably if more replicates are taken and averaged before calculating the correlation (see Fig. 4C), but there is still a rather large variation, induced by the limited sample size. Averaging across the technical replicates reduces variability among the sample means: however this not accompanied by an equal reduction in the variability of the correlation estimation. This is because the error structure is not taken into account in the calculation of the correlation coefficient.

Estimation of Pearson's Correlation Coefficient in Presence of Measurement Error

In the ideal case of an error free measurement, where the only variability is due to intrinsic biological variation, ρ coincides with the true correlation ρ_0 . If additive uncorrelated error is present, then ρ is given by Eqs. (9) and (10) which explicitly take into account the error component; it holds that $\rho < \rho_0$.

In the next Section we will derive analytical expressions, akin to Eqs. (9) and (10), for the correlation for variables sampled with measurement error (additive, multiplicative and realistic) as introduced in Section 2.

Before moving on, we define more specifically the error components. The error terms in models (11), (13) and (15) are assumed to have the following distributional properties

$$\begin{pmatrix} \varepsilon_{au_x} \\ \varepsilon_{au_y} \end{pmatrix} \sim N(0, \Sigma_A) \text{ and } \begin{pmatrix} \varepsilon_{mu_x} \\ \varepsilon_{mu_y} \end{pmatrix} \sim N(0, \Sigma_M) \quad (16)$$

with variance-covariance matrices

$$\Sigma_A = \begin{pmatrix} \sigma_{au_x}^2 & 0 \\ 0 & \sigma_{au_y}^2 \end{pmatrix} \text{ and } \Sigma_M = \begin{pmatrix} \sigma_{mu_x}^2 & 0 \\ 0 & \sigma_{mu_y}^2 \end{pmatrix}, \quad (17)$$

and

$$\varepsilon_{mc} \sim N(0, \sigma_{mc}^2) \text{ and } \varepsilon_{ac} \sim N(0, \sigma_{ac}^2). \quad (18)$$

From definitions (16), (17) and (18) it follows that:

- (1) The expected value of the errors $E[\varepsilon_\alpha]$ is zero:

$$E[\varepsilon_\alpha] = 0 \quad \forall \alpha \text{ in } \{au, ac, mu, mc\}. \quad (19)$$

- (2) The covariance between $x_0(y_0)$ and the error terms is zero because $x_0(y_0)$ and errors are independent,

$$E[x_0\varepsilon_\alpha] - E[x_0]E[\varepsilon_\alpha] = 0 \quad \forall \alpha \text{ in } \{au, ac, mu, mc\}. \quad (20)$$

- (3) The covariance between the different error components is zero because the errors are independent from each other.

$$E[\varepsilon_\alpha\varepsilon_{\alpha'}] - E[\varepsilon_\alpha]E[\varepsilon_{\alpha'}] = 0 \quad \forall \alpha, \alpha' \text{ in } \{au, ac, mu, mc\}. \quad (21)$$

The Pearson correlation in the presence of additive measurement error. We show here a detailed derivation of the correlation among two variables x and y sampled under the additive error model (11). The variance for variable x (similar considerations hold for y) is given by

$$\text{var}(x) = E[x^2] - E[x]^2 \quad (22)$$

where

$$E[x] = E[x_0 + \varepsilon_{au_x} + \varepsilon_{ac}] = \mu_{x_0}. \quad (23)$$

and

$$\begin{aligned} E[x^2] &= E[x_0^2 + \varepsilon_{au_x}^2 + \varepsilon_{ac}^2 + 2x_0\varepsilon_{au_x} + 2x_0\varepsilon_{ac} + 2\varepsilon_{au_x}\varepsilon_{ac}] \\ &= \sigma_{x_0}^2 + \mu_{x_0}^2 + \sigma_{au_x}^2 + \sigma_{ac}^2. \end{aligned} \quad (24)$$

It follows that

$$\text{var}(x) = \sigma_{x_0}^2 + \sigma_{au_x}^2 + \sigma_{ac}^2. \quad (25)$$

The covariance of x and y is

$$\text{cov}(x, y) = E[xy] - E[x]E[y] \quad (26)$$

with

$$\begin{aligned} E[xy] &= E[x_0y_0 + x_0\varepsilon_{au_y} \pm x_0\varepsilon_{ac} + \varepsilon_{au_x}y_0 \\ &\quad + \varepsilon_{au_x}\varepsilon_{au_y} \pm \varepsilon_{au_x}\varepsilon_{ac} + \varepsilon_{ac}y_0 + \varepsilon_{ac}\varepsilon_{au_y} \pm \varepsilon_{ac}^2]. \end{aligned} \quad (27)$$

Considering (20) and (21), Eq. (27) reduces to

$$E[xy] = E[x_0y_0] \pm E[\varepsilon_{ac}^2] \quad (28)$$

with

$$\begin{aligned} E[x_0 y_0] &= \text{cov}(x_0, y_0) + E[x_0] E[y_0] \\ &= \sigma_{x_0 y_0} + \mu_{x_0} \mu_{y_0} \end{aligned} \quad (29)$$

and

$$\pm E[\varepsilon_{ac}^2] = \pm \sigma_{ac}^2, \quad (30)$$

with \pm depending on the sign of the measurement error correlation. From Eqs. (23), (28), (29) and (30) it follows

$$\text{cov}(x, y) = \sigma_{x_0 y_0} \pm \sigma_{ac}^2. \quad (31)$$

Plugging (25) and (31) into (6) and defining the attenuation coefficient A^a

$$A^a = \frac{1}{\sqrt{1 + \frac{\sigma_{au_x}^2}{\sigma_{x_0}^2} + \frac{\sigma_{ac}^2}{\sigma_{x_0}^2}} \sqrt{1 + \frac{\sigma_{au_y}^2}{\sigma_{y_0}^2} + \frac{\sigma_{ac}^2}{\sigma_{y_0}^2}}} = \frac{1}{\sqrt{1 + \xi_x^2 + \gamma_x^2} \sqrt{1 + \xi_y^2 + \gamma_y^2}}, \quad (32)$$

where $\xi_x^2 = \sigma_{au_x}^2 / \sigma_{x_0}^2$, $\xi_y^2 = \sigma_{au_y}^2 / \sigma_{y_0}^2$, $\gamma_x^2 = \sigma_{ac}^2 / \sigma_{x_0}^2$ and $\gamma_y^2 = \sigma_{ac}^2 / \sigma_{y_0}^2$; the superscript a in A^a stands for additive.

The Pearson correlation in presence of additive measurement error is obtained as:

$$\rho = A^a (\rho_0 \pm \gamma_x \gamma_y) \quad (33)$$

where the sign \pm signifies positively and negatively correlated error.

The attenuation coefficient A^a is a decreasing function of the measurement error ratios, that is, the ratio between the variance of the uncorrelated and the correlated error to the variance of the true signal. Compared to Eq. (9), in formula (33) there is an extra additive term related to the correlated measurement error expressing the impact of the correlated measurement error relative to the original variation. In the presence of only uncorrelated error (i.e. $\sigma_{ac}^2 = 0$), Eq. (33) reduces to the Spearman's formula for the correlation attenuation given by (9) and (10). As previously discussed, in this case the correlation coefficient is always biased towards zero (attenuated).

Given the true correlation ρ_0 , the expected correlation coefficient (33) is completely determined by the measurement error ratios. Assuming the errors on x and y to be the same ($\sigma_{au_x}^2 = \sigma_{au_y}^2$, $\sigma_{mu_x}^2 = \sigma_{mu_y}^2$, an assumption not unrealistic if x and y are measured with the same instrument and under the same experimental conditions during an *omics* comprehensive experiment) and taking for simplicity $\sigma_{x_0}^2 = \sigma_{y_0}^2$, then $\xi_x = \xi_y = \xi$ and $\gamma_x = \gamma_y = \gamma$ and Eq. (33) can be simplified to:

$$\rho = \frac{\rho_0 \pm \gamma^2}{1 + \xi^2 + \gamma^2}, \quad (34)$$

and ρ can be visualized graphically as a function of the uncorrelated and correlated measurement error ratios ξ and γ as shown in Fig. 5.

In the presence of positively correlated error, the correlation ρ is attenuated towards 0 if the uncorrelated error increases and inflated if the additive correlated error increases (Fig. 5A, which refers to Eq. (34)) when $\rho_0 > 0$. If $\rho_0 < 0$ the distortion introduced by the correlated error can be so severe that the correlation ρ can become positive. When the error is negatively correlated (Fig. 5B), the correlation ρ is biased towards 0 when $\rho_0 > 0$ (and can change sign), while it can be attenuated or inflated if $\rho_0 < 0$.

A set of rules can be derived to describe quantitatively the bias of ρ . For positively correlated measurement error (for negatively correlated measurement error see Section 6.2) if the true correlation ρ_0 is positive the correlation ρ is always strictly positive: this is shown on Fig. 6A where the relationship between ρ and ρ_0 is shown by means of Monte Carlo simulation (see Figure caption for more details). The magnitude of ρ ($\|\rho\|$) depends on how A^a (for readability in the following equations we will use A) and the additive term $\gamma_x \gamma_y > 0$ compensate each other. In particular when $\rho_0 > 0$

$$\rho \rightarrow \begin{cases} 0 < \rho < \rho_0 & \text{if } \rho_0 > \frac{A}{1-A} \gamma_x \gamma_y \\ \rho_0 & \text{if } \rho_0 = \frac{A}{1-A} \gamma_x \gamma_y \\ > \rho_0 & \text{if } \rho_0 < \frac{A}{1-A} \gamma_x \gamma_y \end{cases} \quad (35)$$

This means that ρ is always a biased estimator of the true correlation ρ_0 , with the exception of the second case which happens only for specific values of γ and ρ_0 . This is unlikely to happen in practice.

If $\rho_0 < 0$ it holds that

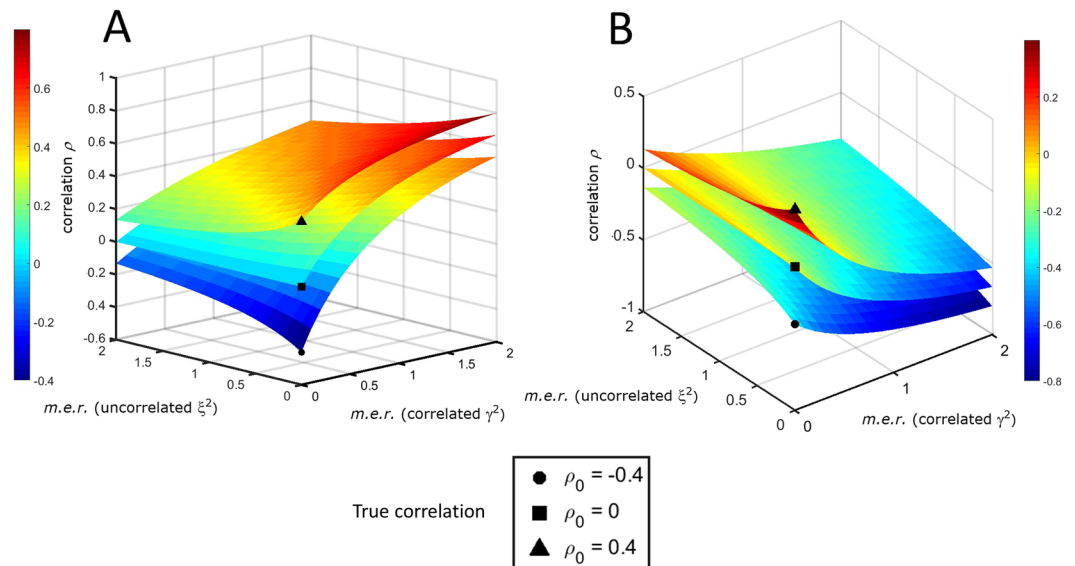


Figure 5. The expected correlation coefficient ρ in the presence of additive measurement error as a function of the uncorrelated (ξ^2) and correlated (γ^2) measurement error ratios (*m.e.r.*) for different values of the true correlation ρ_0 . **(A)** Positively correlated error. **(B)** Negatively correlated error.

$$\rho \rightarrow \begin{cases} 0 < \rho < \rho_0 & \text{if } \rho_0 > \frac{A}{1-A} \gamma_x \gamma_y \\ \rho_0 & \text{if } \rho_0 = \frac{A}{1-A} \gamma_x \gamma_y \\ > \rho_0 & \text{if } \rho_0 < \frac{A}{1-A} \gamma_x \gamma_y \end{cases} \quad (36)$$

The interpretation of Eq. (36) is similar to that of Eq. (35) but additionally, the correlation coefficient can even change sign. In particular, this happens when

$$|\rho_0| > \sqrt{\gamma_x \gamma_y}. \quad (37)$$

The terms $S = \frac{A}{1-A} \gamma_x \gamma_y$ and $S = \frac{A}{A+1} \gamma_x \gamma_y$ in Eqs. (35), (36), (71) and (72) describe limiting surfaces S of ρ_0 values delineating the regions of attenuation and inflation of the correlation coefficient ρ . As can be seen from Fig. 7, these surfaces are not symmetric with respect to zero correlation, indicating that the behavior of ρ is not symmetric around 0 with respect to the sign of ρ_0 and of the correlated error.

The Pearson correlation in presence of multiplicative measurement error. The correlation in the presence of multiplicative error can be derived using similar arguments and detailed calculations can be found in Section 6.1.1. Here we only state the main result:

$$\rho = \rho_0(1 \pm \sigma_{mc}^2)A^m \pm \delta_x \delta_y \sigma_{mc}^2 A^m \quad (38)$$

with $\delta_x = \mu_{x_0}/\sigma_{x_0}$, $\delta_y = \mu_{y_0}/\sigma_{y_0}$ (biological signal to biological variation ratios) and A^m is the attenuation coefficient (the superscript m stands for *multiplicative*):

$$A^m = \frac{1}{\sqrt{1 + \left(1 + \frac{\mu_{x_0}^2}{\sigma_{x_0}^2}\right) \left(\frac{\sigma_{mux}^2}{\sigma_{x_0}^2} + \frac{\sigma_{mc}^2}{\sigma_{x_0}^2}\right)} \sqrt{1 + \left(1 + \frac{\mu_{y_0}^2}{\sigma_{y_0}^2}\right) \left(\frac{\sigma_{myy}^2}{\sigma_{y_0}^2} + \frac{\sigma_{mc}^2}{\sigma_{y_0}^2}\right)}} \quad (39)$$

In this case, the correlation coefficient depends explicitly on the mean of the variables, as an effect of the multiplicative nature of the error component. Our simulations show that if the signal intensity is not too large, the correlation can change sign (as shown in Fig. 6B); if the signal intensity is very large the multiplicative error will have a very large effect and if the correlated error is positive the expected correlation ρ will also be positive, and will be negative if the error are negatively correlated. but simulations cannot be exhaustive (as shown in Fig. 6B).

The Pearson correlation in presence of realistic measurement error. When both additive and multiplicative error are present, the correlation coefficient is a combination of formula (33) and (38) (see Section 6.1.2 for detailed derivation):

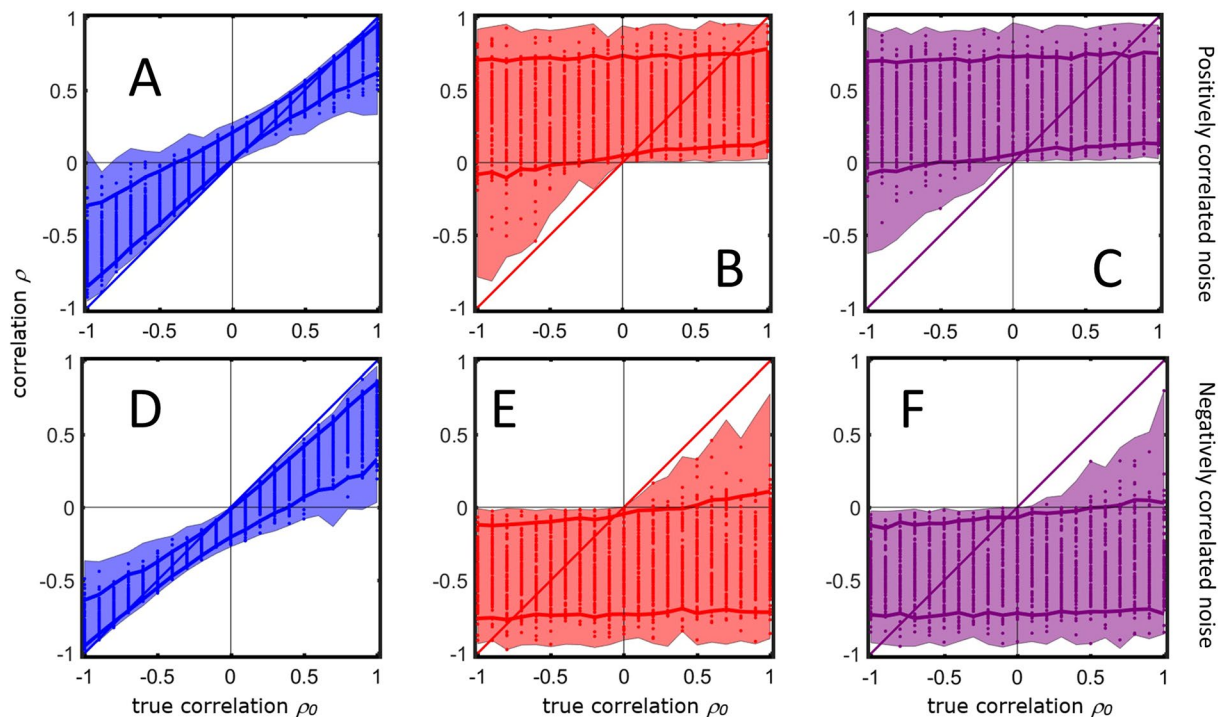


Figure 6. Calculations of the correlation coefficient ρ (40) as a function of the different realizations of the signal means and the size of the error components for different values of the true correlation ρ_0 . The shadowed area encloses the maximum and the minimum of the values of ρ calculated in the simulation using the different error models. The dots represent the realized values of ρ (only 100 of 10^5 Monte Carlo realizations for different values of the variances of error component are shown). The solid lines represent the 5-th and the 95-th percentiles of the observed values. (A) Additive measurement error with positive correlated error. (B) Multiplicative measurement error with positive correlated error. (C) Realistic case with both additive and multiplicative measurement error with positive correlated error. (D) Additive measurement error with negative correlated error. (E) Multiplicative measurement error with negative correlated error. (F) Realistic case with both additive and multiplicative measurement error with negative correlated error. For more details on the simulations see Material and Methods section 6.5.5.

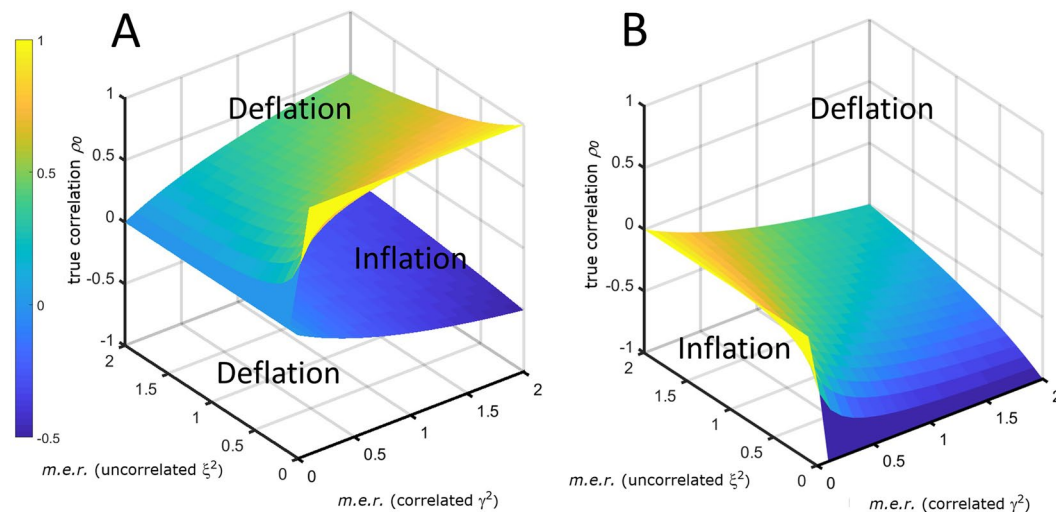


Figure 7. Limiting surfaces S for the inflation and deflation region of the correlation coefficient in presence of additive measurement error. The surfaces are a function of the uncorrelated (ξ^2) and correlated (γ^2) measurement error ratios ($m.e.r.$). (A) S in the case of positively correlated error. (B) S for negatively correlated error. The plot refers to ρ defined by Eq. (34) with $\xi_x^2 = \xi_y^2 = \xi^2$ and $\gamma_x^2 = \gamma_y^2 = \gamma^2$.

$$\rho = \rho_0(1 \pm \sigma_{mc}^2)A^r \pm (\gamma_x\gamma_y + \delta_x\delta_y\sigma_{mc}^2)A^r, \tag{40}$$

where the γ and δ parameters have been previously defined for the additive and multiplicative case. A^r is the attenuation coefficient (the superscript r stands for realistic):

$$A^r = \frac{1}{\sqrt{1 + \left(1 + \frac{\mu_{x0}^2}{\sigma_{x0}^2}\right)\left(\frac{\sigma_{mux}^2}{\sigma_{x0}^2} + \frac{\sigma_{mc}^2}{\sigma_{x0}^2}\right) + \frac{\sigma_{aux}^2}{\sigma_{x0}^2} + \frac{\sigma_{ac}^2}{\sigma_{x0}^2}} \sqrt{1 + \left(1 + \frac{\mu_{y0}^2}{\sigma_{y0}^2}\right)\left(\frac{\sigma_{myy}^2}{\sigma_{y0}^2} + \frac{\sigma_{mc}^2}{\sigma_{y0}^2}\right) + \frac{\sigma_{ayy}^2}{\sigma_{y0}^2} + \frac{\sigma_{ac}^2}{\sigma_{y0}^2}}. \tag{41}$$

General rules governing the sign of the numerator and denominator in Eq. (40) cannot be determined since it depends on the interplay of the six error components, the true mean and product thereof. Within the parameter setting of our simulations, the results presented in Fig. 6C show that the behavior of ρ under error model 15 is qualitatively similar to that in presence of only multiplicative error. However different behavior could be emerge with different parameter settings.

Generalized correlated error model. The error models presented in Eqs. (11), (13) and (15) assume a perfect correlation of the correlated errors, since the correlated error terms ε_{ac} appear simultaneously in both x and y ; the same hold true for ε_{mc} . A more general model that accounts for different degrees of correlation between the error components can be obtained by modifying the model (15) (other cases are treated in Section 6.3.) to

$$\begin{cases} x = x_0(1 + \varepsilon_{mux} + \varepsilon_{mc_x}) + \varepsilon_{aux} + \varepsilon_{ac_x} \\ y = y_0(1 + \varepsilon_{myy} + \varepsilon_{mc_y}) + \varepsilon_{aay} + \varepsilon_{ac_y} \end{cases} \tag{42}$$

where the correlated error components ε_{mc_x} , ε_{ac_x} , ε_{mc_y} and ε_{ac_y} are distributed as

$$\begin{pmatrix} \varepsilon_{ac_x} \\ \varepsilon_{ac_y} \end{pmatrix} \sim N(0, \Sigma_{AC}) \text{ and } \begin{pmatrix} \varepsilon_{mc_x} \\ \varepsilon_{mc_y} \end{pmatrix} \sim N(0, \Sigma_{MC}) \tag{43}$$

with variance-covariance matrices

$$\Sigma_{AC} = \begin{pmatrix} \sigma_{ac_x}^2 & \sigma_{ac_{xy}} \\ \sigma_{ac_{xy}} & \sigma_{ac_y}^2 \end{pmatrix} \text{ and } \Sigma_{MC} = \begin{pmatrix} \sigma_{mc_x}^2 & \sigma_{mc_{xy}} \\ \sigma_{mc_{xy}} & \sigma_{mc_y}^2 \end{pmatrix}, \tag{44}$$

where $\sigma_{ac_{xy}}$ is the covariance between error term ε_{ac_x} and ε_{ac_y} and $\sigma_{mc_{xy}}$ is the covariance between error term ε_{mc_x} and ε_{mc_y} .

It is possible to derive expression for the correlation coefficient under the model (43) as shown in Section 3.1 and in the Section 6.1.1 and 6.1.2. The only difference is that under this model the terms $E[\varepsilon_{ac}^2]$ and $E[\varepsilon_{mc}^2]$ in Eqs. (27), (58), (65) and (66) are replaced by $E[\varepsilon_{ac_x}, \varepsilon_{ac_y}] = \sigma_{ac_{xy}}$ and $E[\varepsilon_{mc_x}, \varepsilon_{mc_y}] = \sigma_{mc_{xy}}$, respectively.

From the definition of covariance it follows that

$$\sigma_{ac_{xy}} = \pi_{ac} \sqrt{\sigma_{ac_x}^2 \sigma_{ac_y}^2} \tag{45}$$

and

$$\sigma_{mc_{xy}} = \pi_{mc} \sqrt{\sigma_{mc_x}^2 \sigma_{mc_y}^2}, \tag{46}$$

where π_{ac} and π_{mc} are the correlations among the error terms for which it holds $-1 \leq \pi_{mc} \leq 1$ and $-1 \leq \pi_{ac} \leq 1$. If π_{ac} and π_{mc} are negative the errors are negatively correlated. Equation (40) becomes now:

$$\rho = \rho_0(1 + \pi_{mc}\sigma_{mc_x}\sigma_{mc_y})A^r + (\pi_{ac}\gamma_x\gamma_y + \delta_x\delta_y\pi_{mc}\sigma_{mc_x}\sigma_{mc_y})A^r, \tag{47}$$

with $\gamma_x = \sigma_{ac_x}/\sigma_{x_0}$, $\gamma_y = \sigma_{ac_y}/\sigma_{y_0}$, and

$$A^r = \frac{1}{\sqrt{1 + \left(1 + \frac{\mu_{x0}^2}{\sigma_{x0}^2}\right)\left(\frac{\sigma_{mux}^2}{\sigma_{x0}^2} + \frac{\sigma_{mc_x}^2}{\sigma_{x0}^2}\right) + \frac{\sigma_{aux}^2}{\sigma_{x0}^2} + \frac{\sigma_{ac_x}^2}{\sigma_{x0}^2}} \times \sqrt{1 + \left(1 + \frac{\mu_{y0}^2}{\sigma_{y0}^2}\right)\left(\frac{\sigma_{myy}^2}{\sigma_{y0}^2} + \frac{\sigma_{mc_y}^2}{\sigma_{y0}^2}\right) + \frac{\sigma_{aay}^2}{\sigma_{y0}^2} + \frac{\sigma_{ac_y}^2}{\sigma_{y0}^2}}. \tag{48}$$

This model generalizes the correlation coefficient among x and y from Eq. (40) to account for different strength of the correlation among the correlated error components. All considerations discussed in the previous sections do apply also to this model. Expressions for ρ in the case of additive and multiplicative error can be found in the Section 6.3.1 and 6.3.2.

By setting $\sigma_{ac_x}^2 = \sigma_{ac_y}^2 = \sigma_{ac}^2$, $\sigma_{mc_x}^2 = \sigma_{mc_y}^2 = \sigma_{mc}^2$, and $\pi_{ac} = \pi_{mc} = 1$ (perfect correlation), model (40) is obtained, and similarly models (33) and (38).

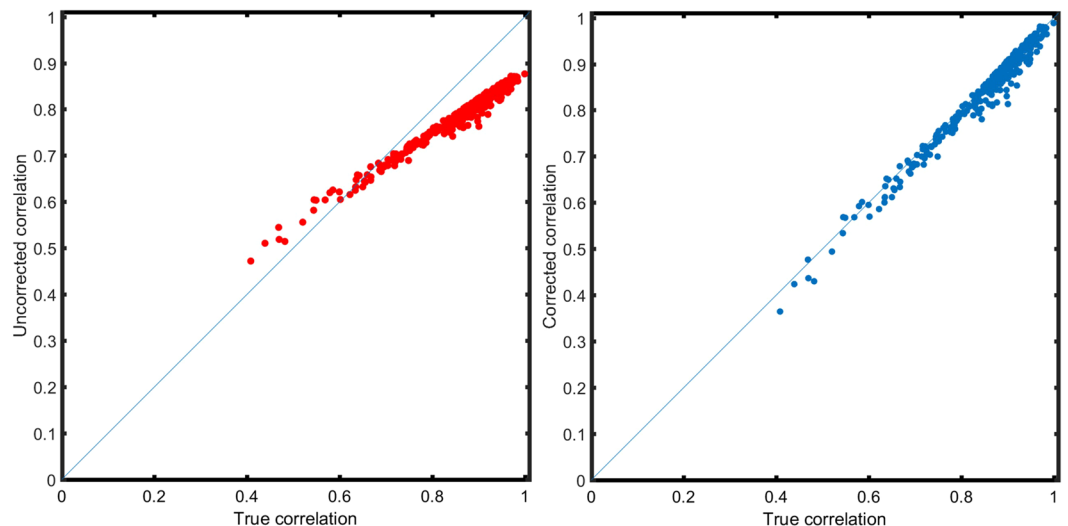


Figure 8. Correction of the distortion induced by the realistic measurement error (see Eq. (15)). **(A)** Pairwise correlations ρ among 25 metabolites calculated from simulated data with additive and multiplicative measurement error vs the true correlation ρ_0 . **(B)** Corrected correlation coefficients using Eq. (52) and using the known error variance components. See Section 6.5.6 for details on the data simulation.

Correction for Correlation Bias

Because virtually all kinds of measurement are affected by measurement error, the correlation calculated from sampled data is distorted to some degree depending on the level of the measurement error and on its nature. We have seen that experimental error can inflate or deflate the correlation and that ρ (and hence its sample realization r) is almost always a biased estimation of the true correlation ρ_0 . An estimator that gives a theoretically unbiased estimate of the correlation coefficient between two variables x and y taking into account the measurement error model can be derived. For simple uncorrelated additive error this is given by the Spearman's formula (49): this is a known results which in the past has been presented and discussed in many different fields^{16–19}. To obtain similar correction formulas for the error models considered here it is sufficient to solve for ρ_0 from the defining Eqs. (33), (38) and (40). The correction formulas are as follows (the \pm indicates positive and negatively correlated error):

1. Correction for simple additive error (only uncorrelated error):

$$\rho_0 = A^{-1}\rho. \quad (49)$$

2. Correction for additive error:

$$\rho_{\pm}^{\text{corrected}} = \frac{1}{A^a}\rho \mp \gamma_x\gamma_y. \quad (50)$$

3. Correction for multiplicative error:

$$\rho_{\pm}^{\text{corrected}} = \frac{1}{A^m(1 \pm \sigma_{mc}^2)}\rho \mp \frac{\sigma_{mc}^2}{1 \pm \sigma_{mc}^2}\delta_x\delta_y. \quad (51)$$

4. Correction for realistic error:

$$\rho_{\pm}^{\text{corrected}} = \frac{1}{A^c(1 \pm \sigma_{mc}^2)}\rho \mp \frac{\gamma_x\gamma_y + \delta_x\delta_y\sigma_{mc}^2}{1 \pm \sigma_{mc}^2}. \quad (52)$$

In practice, to obtain a corrected estimation of the correlation coefficient ρ_0 , the ρ is substituted by r in (50), (51) and (52), which is the sample correlation calculated from the data. The effect of the correction is shown, for the realistic error model (15), in Fig. 8 where the true know error variance components have been used. It should be noted that it is possible that the corrected correlation exceeds ± 1.0 . This phenomenon has already been observed and discussed^{16,30}: it is due to the fact that the sampling error of a correlation coefficient corrected for distortion is greater than would be that of an uncorrected coefficient of the same size (at least for the uncorrelated additive error^{4,18,31}). When this happens the corrected correlation can be rounded to ± 1.0 ^{19,31}.

Estimation of the error variance components. Simulations shown in Fig. 8 have been performed using the known parameters for the error components used to generate the data. In practical applications the error

components needs to be estimated from the measured data and the quality of the correction will depend on the accuracy of the error variance estimate.

The case of purely additive uncorrelated measurement error ($\sigma_{ac}^2 = 0$) has been addressed in the past^{18,19,32}: in this case the variance components $\sigma_{x_0}^2$ and $\sigma_{y_0}^2$ can be substituted with their sample estimates ($s_{x_0}^2$ and $s_{y_0}^2$) obtained from measured data, while the error variance components ($\sigma_{au_x}^2$ and $\sigma_{au_y}^2$) can be estimated if an appropriate experimental design is implemented, *i.e.* if n replicates are measured for each observation.

Unfortunately, there is no simple and immediate approach to estimate the error component in the other cases when many variance components need to be estimated (6 error variances in the case of error model (15) and 8 in the case of the generalized model (42), to which the estimations of π_{mc} and π_{ac} must be added).

Different approaches can be foreseen to estimate the error components which is not a trivial task, including the use of (generalized) linear mixed model^{33,34}, error covariance matrix formulation^{29,35,36} or common factor analysis factorization³⁷. None of these approaches is straightforward and require some extensive mathematical manipulations to be implemented; an accurate investigation of the simulation of the error component is outside the scope of this paper and will be presented in a future publication.

Discussion

Since measurement error cannot be avoided, correlation coefficients calculated from experimental data are distorted to a degree which is not known and that has been neglected in life sciences applications but can be expected to be considerable when comprehensive omics measurement are taken.

As previously discussed, the attenuation of the correlation coefficient in the presence of additive (uncorrelated) error has been known for more than one century. The analytical description of the distortion of the correlation coefficient in presence of more complex measurement error structures (Eqs. (33), (38) and (40)) has been presented here for the first time to the best of our knowledge.

The inflation or attenuation of the correlation coefficient depends on the relationship between the value of true correlation ρ_0 and the error component. In most cases in practice, ρ is a biased estimator for ρ_0 . In absence of correlated error, there is always attenuation; in the presence of correlated error there can also be increase (in absolute value) of the correlation coefficient. This has also been observed in regression analysis applied to nutritional epidemiology and it has been suggested that correlated error can, in principle, be used to compensate for the attenuation³⁸. Moreover, the distortion of the correlation coefficient also has implications for hypothesis testing to assess the significance of the measured correlation r .

To illustrate the counterintuitive consequences of correlated measurement error consider the following. Suppose that the true correlation is null. In that case, Eqs. (33), (38) and (40) reduce to

$$\rho = A^a \gamma_x \gamma_y, \quad (53)$$

$$\rho = A^m \delta_x \delta_y \sigma_{mc}^2, \quad (54)$$

and

$$\rho = \pm (\gamma_x \gamma_y + \delta_x \delta_y \sigma_{mc}^2) A^r, \quad (55)$$

which implies that the correlation coefficient is not zero. Moreover, in real-life situations there is also sampling variability superimposed on this which may in the end result in estimated correlations of the size as found in several omics applications (in metabolomics observed correlations are usually lower than 0.6^{10,23}; similar patterns are also observed in transcriptomics^{39,40}) while the true biological correlation is zero.

The correction equations presented need the input of estimated variances. Such estimates also carry uncertainty and the quality of these estimates will influence the quality of the corrections. This will be the topic of a follow-up paper. Prior information regarding the sizes of the variance components would be valuable and this points to new requirements for system suitability tests of comprehensive measurements. In metabolomics, for example, it would be worthwhile to characterize an analytical measurement platform in terms of such error variances including sizes of correlated error using advanced (and to be developed) measurement protocols.

Distortion of the correlation coefficient has implications also for experimental planning. In the case of additive uncorrelated error, the correction depends explicitly on the sample size N used to calculate r and on the number of replicates n_x, n_y used to estimate the intraclass correlation (*i.e.* the error variance components): since in real life the total sample size $N \times (n_x + n_y)$ is fixed, there is a trade off between the sample size and the number of replicates that can be measured and the experimenter has to decide whether to increase N or n_x .

The results presented here are derived under the assumption of normality of both measurement and measurement errors. If x_0 and y_0 are normally distributed, then x and y will be, in presence of additive measurement error, normally distributed, with variance given by (12). For multiplicative and realistic error the distribution of x and y will be far from normality since it involves the distribution of the product of normally distributed quantities which is usually not normal⁴¹. It is known that departure from normality can result in the inflation of the correlation coefficient⁴² and in distortion⁴³ of its (sampling) distribution and this will add to the corruption induced by the measurement error.

We think that in general correlation coefficients are trusted too much on face value and we hope to have triggered some doubts and pointed to precautions in this paper.

Material and Methods

Mathematical calculations. Derivation of ρ in presence of multiplicative measurement error. In presence of purely multiplicative error it holds

$$E[x] = E[x_0(1 + \varepsilon_{mu_x} \pm \varepsilon_{mc})] = \mu_{x_0} \quad (56)$$

and

$$\begin{aligned} E[x^2] &= E[x_0^2 + x_0^2(\varepsilon_{mu_x}^2 + \varepsilon_{mc}^2 + 2\varepsilon_{mu_x} + \varepsilon_{mc} \pm \varepsilon_{mc} \pm 2\varepsilon_{mu_x}\varepsilon_{mc})] \\ &= \sigma_{x_0}^2 + \mu_{x_0}^2 + \sigma_{mu_x}^2(\sigma_{x_0}^2 + \mu_{x_0}^2) + \sigma_{mc}^2(\sigma_{x_0}^2 + \mu_{x_0}^2), \end{aligned} \quad (57)$$

using (19)–(21) to calculate the expectation of the cross terms. For $E[xy]$ it holds

$$\begin{aligned} E[xy] &= E[x_0y_0 + x_0y_0(\varepsilon_{mc} \pm \varepsilon_{mc} \pm \varepsilon_{mc}^2 \pm \varepsilon_{mc}\varepsilon_{mu_y} \\ &= + \varepsilon_{mc}\varepsilon_{mu_x} + \varepsilon_{mu_x} + \varepsilon_{mu_y} + \varepsilon_{mu_x}\varepsilon_{mu_y})]. \end{aligned} \quad (58)$$

Because of the independence of x_0 , y_0 and the error terms, the expectations of all cross terms is null except

$$\begin{aligned} \pm E[x_0y_0\varepsilon_{mc}^2] &= \pm E[x_0y_0]E[\varepsilon_{mc}^2] \\ &= \pm \sigma_{mc}^2(\sigma_{x_0y_0}^2 + \mu_{x_0}\mu_{y_0}), \end{aligned} \quad (59)$$

where $E[x_0y_0]$ is given by Eq. (29). Plugging (56), (57) and (58) in (6), the expected correlation coefficient is

$$\rho = \frac{\sigma_{x_0y_0} \pm \sigma_{mc}^2(\sigma_{x_0y_0} + \mu_{x_0}\mu_{y_0})}{\sqrt{(\sigma_{x_0}^2 + (\sigma_{x_0}^2 + \mu_{x_0}^2)(\sigma_{mu_x}^2 + \sigma_{mc}^2))(\sigma_{y_0}^2 + (\sigma_{y_0}^2 + \mu_{y_0}^2)(\sigma_{mu_y}^2 + \sigma_{mc}^2))}}, \quad (60)$$

and it can re-written as (38) by setting $\gamma_x = \sigma_{ac}^2/\sigma_{x_0}^2$ and $\gamma_y = \sigma_{ac}^2/\sigma_{y_0}^2$, $\delta_x = \mu_{x_0}/\sigma_{x_0}$, $\delta_y = \mu_{y_0}/\sigma_{y_0}$ and defining the attenuation coefficient A^m (39).

Derivation of ρ in presence of realistic measurement error. To simplify calculations we set

$$\begin{cases} M_x = x_0(1 + \varepsilon_{mu_x} \pm \varepsilon_{mc}) \\ A_x = \varepsilon_{au_x} \pm \varepsilon_{ac} \end{cases} \quad (61)$$

and similarly we define M_y and A_y for variable y . It holds

$$E[A_x] = 0 \text{ and } E[M_x] = \mu_{x_0} \quad (62)$$

and

$$E[A_x^2] = \sigma_{au_x}^2 + \sigma_{ac}^2. \quad (63)$$

$E[M_x^2]$ is given by Eq. (57). Because error components are independent and with zero expectation (see Eqs. (19)–(21)) it holds

$$E[M_x A_x] = E[M_x A_y] = E[M_y A_x] = 0, \quad (64)$$

$$E[M_x M_y] = \sigma_{x_0y_0} + \mu_{x_0}\mu_{y_0} \pm (\sigma_{x_0y_0} + \mu_{x_0}\mu_{y_0})\sigma_{mc}^2, \quad (65)$$

$$E[A_x A_y] = \pm \sigma_{ac}^2. \quad (66)$$

It follows that

$$E[x] = \mu_{x_0}, \quad (67)$$

$$\begin{aligned} E[x^2] &= E[M_x^2] + E[A_x^2] + 2E[M_x A_x] \\ &= \sigma_{x_0}^2 + \mu_{x_0}^2 + \sigma_{mu_x}^2(\sigma_{x_0}^2 + \mu_{x_0}^2) + \sigma_{mc}^2(\sigma_{x_0}^2 + \mu_{x_0}^2) + \sigma_{au_x}^2 + \sigma_{ac}^2, \end{aligned} \quad (68)$$

and

$$E[xy] = E[M_x M_y] + E[A_x A_y] + E[M_x A_y] + E[M_y A_x]$$

$$= \sigma_{x_0 y_0} + \mu_{x_0} \mu_{y_0} \pm (\sigma_{x_0 y_0} + \mu_{x_0} \mu_{y_0}) \sigma_{\epsilon}^2 \pm \sigma_{ac}^2. \tag{69}$$

Plugging (67), (68), and (69) into (6) one gets the expression for the correlation coefficient in presence of additive and multiplicative measurement error:

$$\rho = \frac{\sigma_{x_0 y_0} \pm (\sigma_{x_0 y_0} + \mu_{x_0} \mu_{y_0}) \sigma_{mc}^2 \pm \sigma_{ac}^2}{\sqrt{(\sigma_{x_0}^2 + (\sigma_{x_0}^2 + \mu_{x_0}^2)(\sigma_{mu_x}^2 + \sigma_{mc}^2) + \sigma_{au_x}^2 + \sigma_{ac}^2) \times (\sigma_{y_0}^2 + (\sigma_{y_0}^2 + \mu_{y_0}^2)(\sigma_{mu_y}^2 + \sigma_{mc}^2) + \sigma_{au_y}^2 + \sigma_{ac}^2)}}, \tag{70}$$

that can re-written as (40) by using previously defined γ_x , γ_y , δ_x and δ_y and defining the attenuation coefficient A^e (41).

Behavior of ρ in the case of additive negatively correlated error. For negative correlated error, when the true correlation is positive

$$\rho \rightarrow \begin{cases} 0 < \rho < \rho_0 & \text{if } \rho_0 > \frac{A}{A-1} \gamma_x \gamma_y \\ \rho_0 & \text{if } \rho_0 = \frac{A}{A-1} \gamma_x \gamma_y \\ > \rho_0 & \text{if } \rho_0 < \frac{A}{A-1} \gamma_x \gamma_y \end{cases} \tag{71}$$

Since $\frac{A}{A-1} \gamma_x \gamma_y < 0$, ρ is always smaller than the true correlation. When the true correlation is negative ($\rho_0 < 0$) the expected correlation is always negative, but it can be, in absolute value, smaller or larger than the absolute value of the true correlation:

$$\rho \rightarrow \begin{cases} < \rho_0 & \text{if } \frac{A}{A-1} \gamma_x \gamma_y < \rho_0 < 0 \\ = \rho_0 & \text{if } \rho_0 = \frac{A}{A-1} \gamma_x \gamma_y \\ > \rho_0 & \text{if } \rho_0 > \frac{A}{A-1} \gamma_x \gamma_y \end{cases} \tag{72}$$

Correlation coefficient under the generalized error model. *Additive error.* Under the generalized additive correlated error model

$$\begin{cases} x = x_0 + \varepsilon_{au_x} + \varepsilon_{ac_x} \\ y = y_0 + \varepsilon_{au_y} + \varepsilon_{ac_y} \end{cases} \tag{73}$$

with ε_{ac_x} and ε_{ac_y} defined in Eq. (43), the correlation coefficient can be expressed as:

$$\rho = A^a(\rho_0 + \pi_{ac} \gamma_x \gamma_y), \tag{74}$$

with $\gamma_x = \sigma_{ac_x} / \sigma_{x_0}$, $\gamma_y = \sigma_{ac_y} / \sigma_{y_0}$, and

$$A^a = \frac{1}{\sqrt{1 + \frac{\sigma_{au_x}^2}{\sigma_{x_0}^2} + \frac{\sigma_{ac_x}^2}{\sigma_{x_0}^2}} \sqrt{1 + \frac{\sigma_{au_y}^2}{\sigma_{y_0}^2} + \frac{\sigma_{ac_y}^2}{\sigma_{y_0}^2}}}. \tag{75}$$

Multiplicative error. Under the generalized multiplicative error model

$$\begin{cases} x = x_0(1 + \varepsilon_{mu_x} + \varepsilon_{mc_x}) \\ y = y_0(1 + \varepsilon_{mu_y} + \varepsilon_{mc_y}) \end{cases} \tag{76}$$

with ε_{mc_x} and ε_{mc_y} defined in Eq. (43), the correlation coefficient can be expressed as:

$$\rho = \rho_0(1 + \pi_{mc} \sigma_{mc_x} \sigma_{mc_y}) A^m + \delta_x \delta_y \pi_{mc} \sigma_{mc_x} \sigma_{mc_y} A^m \tag{77}$$

with

$$A^m = \frac{1}{\sqrt{\left(1 + \left(\frac{\mu_{x0}^2}{\sigma_{x0}^2}\right)\left(\frac{\sigma_{m_{ux}}^2}{\sigma_{x0}^2} + \frac{\sigma_{m_{cx}}^2}{\sigma_{x0}^2}\right)\right)} \sqrt{\left(1 + \left(\frac{\mu_{y0}^2}{\sigma_{y0}^2}\right)\left(\frac{\sigma_{m_{uy}}^2}{\sigma_{y0}^2} + \frac{\sigma_{m_{cy}}^2}{\sigma_{y0}^2}\right)\right)}}. \quad (78)$$

General realistic error. Formulas for the correlation coefficient under the generalized realistic correlated error model are to be found in the main text in Eqs. (47) and (48).

Correction of the correlation coefficient under the generalized correlated error model. *Additive error.* Under the generalized additive correlated error model the corrected correlation coefficient is

$$\rho^{corrected} = \frac{1}{A^m} \rho - \pi_{ac} \gamma_x \gamma_y. \quad (79)$$

Multiplicative error. Under the generalized multiplicative correlated error model the corrected correlation coefficient is

$$\rho^{corrected} = \frac{1}{A^m (1 + \pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}})} \rho - \frac{\pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}}}{1 + \pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}}} \delta_x \delta_y. \quad (80)$$

Realistic error. Under the generalized realistic correlated error model the corrected correlation coefficient is

$$\rho^{corrected} = \frac{1}{A^m (1 + \pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}})} \rho - \frac{\pi_{ac} \gamma_x \gamma_y + \delta_x \delta_y \pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}}}{1 + \pi_{mc} \sigma_{m_{cx}} \sigma_{m_{cy}}}. \quad (81)$$

Simulations. We provide here details on the simulation performed and shown in Figs. 1–4, 6 and 8.

Simulations in Figure 1. $N = 100$ realizations of two variables x and y were generated under model with additive uncorrelated measurement error (11), with $\rho_0 = 0.8$, $\sigma_{x0}^2 = \sigma_{y0}^2 = 1$ and $\mu = (100, 100)$. Error variance components were set to $\sigma_{au_x}^2 = \sigma_{au_y}^2 = 0$ and to $\sigma_{au_x}^2 = \sigma_{au_y}^2 = 0.75$ (Panel A).

Simulations in Figure 2. The time concentrations profiles $P_1(t)$, $P_2(t)$ and $P_3(t)$ of three hypothetical metabolites P1, P2 and P3 are simulated using the following dynamic model

$$\begin{cases} \frac{d}{dt} P_1(t) = -k_1 P_1(t) (E_T - P_2(t)) + k_{-1} P_2(t) \\ \frac{d}{dt} P_2(t) = -k_{-1} P_1(t) + k_1 P_1(t) (E_T - P_2(t)) - k_2 P_2(t) \\ \frac{d}{dt} P_3(t) = +k_2 P_2(t) \end{cases} \quad (82)$$

which is the model of an irreversible enzyme-catalyzed reaction described by Michaelis-Menten kinetics. Using this model, $N = 100$ concentration time profiles for P1, P2 and P3 were generated by solving the system of differential equations after varying the kinetic parameters k_1 , k_{-1} and k_2 by sampling them from a uniform distribution. For the realization of the j th concentration profile

$$\begin{aligned} k_1^j &\approx U(0.9 \times k_1, 1.1 \times k_1) \\ k_{-1}^j &\approx U(0.9 \times k_{-1}, 1.1 \times k_{-1}) \\ k_2^j &\approx U(0.9 \times k_2, 1.1 \times k_2) \\ E_T^j &\approx U(0.9 \times E_T, 1.1 \times E_T) \end{aligned} \quad (83)$$

with population values $k_1 = 30$, $k_{-1} = 20$, $k_2 = 10$, and $E_T = 1$. Initial conditions were set to $(P_{10}, P_{20}, P_{30}) = (P_{10}^j, 0, 0)$ with $P_{10}^j \approx U(0.9 \times P_{10}, 1.1 \times P_{10})$ and $P_{10} = 5$. All quantities are in arbitrary units. Time profiles were sampled at $t = 0.4$ a.u. and collected in a data matrix \mathbf{X}_0 of size 100×3 . The variability in data matrix \mathbf{X}_0 is given by biological variation. The concentration time profiles of P1, P2 and P3 shown in Panel A are obtained using the population values for the kinetic parameters and for the initial conditions.

Additive uncorrelated and correlated measurement error is added on \mathbf{X}_0 following model (11) where P1, P2 and P3 in \mathbf{X}_0 play the role of x_0 , y_0 and of an additional third variable z_0 which follows a similar model. The variance of the error component was varied in 50 steps between 0 and 25% of the sample variance $s_{x_0}^2$, $s_{y_0}^2$ and $s_{z_0}^2$ calculated from \mathbf{X}_0 . The variance of the correlated error was set to $\sigma_{ac}^2 = 0.05$ in all simulations. Pairwise Pearson correlations $r_{i,j}$ with $i, j = \{P1, P2, P3\}$ were calculated for the error free case \mathbf{X}_0 and for data with measurement error added. 100 error realizations were simulated for each error value and the average correlation across the 100 realization is calculated and it is shown in Panel B.

The “mini” metabolite-metabolite association networks shown in Panel C are defined by first taking the Pearson correlation r_{ij} among P1, P2 and P3 and then imposing a threshold on r to define the connectivity matrix A_{ij}

$$A_{ij} = \begin{cases} 1 & \text{if } |r_{ij}| > 0.6 \\ 0 & \text{otherwise.} \end{cases} \quad (84)$$

For more details see reference¹⁰.

Simulations in Figure 3. Principal component analysis was performed on a 100×133 experimental metabolomic data set (see Section 6.6 for a description). The 15 variables with the highest loading (in absolute value) and the 45 variables with the smallest loading (in absolute value) on the first principal component were selected to form a 100×60 data set \mathbf{X}_0 (we call this now the error free data, as if it only contained biological variation). On this subset a new principal component analysis was performed. Then multiplicative correlated and uncorrelated measurement error was added on \mathbf{X}_0 . The variance of the additive error was set $\sigma_{mu_j}^2 = 0.05 \times s_{j_0}^2$ with $j = 1, 2, \dots, 60$ where $s_{j_0}^2$ is the variance calculated for the j th column of \mathbf{X}_0 , i.e., the biological variance. The variance of the correlated error was fixed to 5% of the average variance observed in the error free data ($\sigma_{mc}^2 = 0.045$).

Simulations in Figure 4. Let x_{ij} and y_{ij} denote the intensities of the resonances measured at 3.23 and 4.98 in the randomly drawn replicate j of sample F_i ($i = 1, 2, \dots, 5$) and define the 5×1 vectors of means

$$\mathbf{x}_j = \frac{1}{J} \begin{pmatrix} \sum_j x_{1j} \\ \vdots \\ \sum_j x_{5j} \end{pmatrix} \text{ and } \mathbf{y}_j = \frac{1}{J} \begin{pmatrix} \sum_j y_{1j} \\ \vdots \\ \sum_j y_{5j} \end{pmatrix}. \quad (85)$$

The correlation $r_j = \text{corr}(\mathbf{x}_j, \mathbf{y}_j)$ is calculated for $J = 1, 2, 5$, and 10; for each J the replicates used to calculate \mathbf{x}_j and \mathbf{y}_j are randomly and independently sampled, for each sample separately, from the total set of the 12 to 15 replicates available per sample. The procedure is repeated 10^5 times to construct the distributions of the correlation coefficient shown in Fig. 4C.

Simulations in Figure 6. Simulation results presented in Fig. 6 show the results from calculations of the sample correlation coefficient as a function of the true correlation ρ_0 and of the true means (μ_{x_0} and μ_{y_0}), the variances ($\sigma_{x_0}^2$ and $\sigma_{y_0}^2$) of the signals x_0 and y_0 and the measurement error variances as they appear in the definitions of ρ under the different error models (Eqs. (33), (38) and (40)). The calculations were done multiple times for varying values for μ_{x_0} and μ_{y_0} , which were randomly and independently sampled from a uniform distribution $U(0, \mu_0)$, where μ_0 was set to be equal to 23.4, which was the maximum values observed in Data set 1 (see Section 6.6). Values for $\sigma_{x_0}^2$ and $\sigma_{y_0}^2$ were randomly and independently sampled from a uniform distribution $U(0, \sigma_0^2)$, where σ_0^2 was set to be equal to the average variance observed in the experimental Data set 1. The values of the variance of all error components are randomly and independently sampled from $U(0, \frac{1}{4}\sigma_0^2)$. The overall procedure was repeated 10^4 for each value of ρ_0 in the range $[-1, 1]$ in steps of 0.1.

Simulations in Figure 8. The first 25 variables from Data set 1 have been selected and used to compute the means μ_0 and the correlation/covariance matrix Σ_0 used to generate error-free data $\mathbf{X}_0 \sim N(\mu_0, \Sigma_0)$ of size $10^4 \times 25$ on which additive and multiplicative measurement error (correlated and uncorrelated) is added (error model (15)) to obtain \mathbf{X} . All error variances are set to 0.1 which is approximately equal to 5% of the average variance observed in \mathbf{X}_0 . Pairwise correlations among the 25 metabolites are calculated from \mathbf{X} . The correlations are corrected using Eq. (52) using the known distributional and error parameters (μ_0, Σ_0) used to generate the data. The data generation is repeated 10^3 times and correlations (uncorrected and corrected) are averaged over the repetitions.

Data sets. *Data set 1.* A publicly available data set containing measurements of 133 blood metabolites from 2139 subjects was used as a base for the simulation to obtain realistic distributional and correlation patterns among measured features. The data comes from a designed case-cohort and a matched sub-cohort (controls) stratified on age and sex from the TwinGene project⁴⁴. The first 100 observation were used in the simulation described in Section 6.5.3 and shown in Fig. 3.

Data were downloaded from the Metabolights public repository⁴⁵ (www.ebi.ac.uk/metabolights) with accession number MTBLS93. For full details on the study protocol, sample collection, chromatography, GC-MS experiments and metabolites identification and quantification see the original publication⁴⁶ and the Metabolights accession page.

Data set 2. This data set was acquired in the framework of a study aiming to the “Characterization of the measurement error structure in Nuclear Magnetic Resonance (NMR) data for metabolomic studies²⁹”. Five biological

replicates of fish extract F1 - F5 were originally pretreated in replicates (12 to 15) and acquired using ^1H NMR. The replicates account for variability in sample preparation and instrumental variability. For details on the sample preparation and NMR experiments we refer to the original publication.

Software. All calculations were performed in Matlab (version 2017a 9.2). Code to generate data under the measurement error models (11), (13) and (15) is available at systemsbiology.nl under the SOFTWARE tab.

Received: 21 June 2019; Accepted: 22 December 2019;

Published online: 16 January 2020

References

1. Bravais, A. *Analyse mathématique sur les probabilités des erreurs de situation d'un point* (Impr. Royale, 1844).
2. Galton, F. Co-relations and their measurement, chiefly from anthropometric data. *Proceedings of the Royal Society of London* **45**, 135–145 (1889).
3. Pearson, K. Note on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London* **58**, 240–242 (1895).
4. Spearman, C. Demonstration of formulae for true measurement of correlation. *The American Journal of Psychology* 161–169 (1907).
5. Pearson, K. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**, 559–572 (1901).
6. Hotelling, H. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology* **24**, 417 (1933).
7. Jolliffe, I. *Principal component analysis* (Springer, 2011).
8. Härdle, W. & Simar, L. *Applied multivariate statistical analysis*, vol. 22007 (Springer, 2007).
9. Müller-Linow, M., Weckwerth, W. & Hütt, M.-T. Consistency analysis of metabolic correlation networks. *BMC Systems Biology* **1**, 44 (2007).
10. Jahagirdar, S., Suarez-Diez, M. & Saccenti, E. Simulation and reconstruction of metabolite-metabolite association networks using a metabolic dynamic model and correlation based-algorithms. *Journal of proteome research* (2019).
11. Dunlop, M. J., Cox, R. S. III., Levine, J. H., Murray, R. M. & Elowitz, M. B. Regulatory activity revealed by dynamic correlations in gene expression noise. *Nature genetics* **40**, 1493 (2008).
12. Marbach, D. *et al.* Wisdom of crowds for robust gene network inference. *Nature Methods* **9**, 796–804, <https://doi.org/10.1038/nmeth.2016> (2012).
13. Stuart, J. M., Segal, E., Koller, D. & Kim, S. K. A gene-coexpression network for global discovery of conserved genetic modules. *Science* **302**, 249–255 (2003).
14. Zhang, B. & Horvath, S. A general framework for weighted gene co-expression network analysis. *Statistical applications in genetics and molecular biology* **4** (2005).
15. Spearman, C. The proof and measurement of association between two things. *The American journal of psychology* **15**, 72–101 (1904).
16. Thouless, R. H. The effects of errors of measurement on correlation coefficients. *British Journal of Psychology* **29**, 383 (1939).
17. Beaton, G. H. *et al.* Sources of variance in 24-hour dietary recall data: implications for nutrition study design and interpretation. *The American journal of clinical nutrition* **32**, 2546–2559 (1979).
18. Rosner, B. & Willett, W. Interval estimates for correlation coefficients corrected for within-person variation: implications for study design and hypothesis testing. *American journal of epidemiology* **127**, 377–386 (1988).
19. Adolph, S. C. & Hardin, J. S. Estimating phenotypic correlations: correcting for bias due to intraindividual variability. *Functional Ecology* **21**, 178–184 (2007).
20. Fuller, W. A. *Measurement error models*, vol. 305 (John Wiley & Sons, 2009).
21. Moseley, H. N. Error analysis and propagation in metabolomics data analysis. *Computational and structural biotechnology journal* **4**, e201301006 (2013).
22. Rosato, A. *et al.* From correlation to causation: analysis of metabolomics data using systems biology approaches. *Metabolomics* **14**, 37 (2018).
23. Camacho, D., de la Fuente, A. & Mendes, P. The origin of correlations in metabolomics data. *Metabolomics* **1**, 53–63, <https://doi.org/10.1007/s11306-005-1107-3> (2005).
24. Werner, M., Brooks, S. H. & Knott, L. B. Additive, multiplicative, and mixed analytical errors. *Clinical chemistry* **24**, 1895–1898 (1978).
25. Balwiercz, P. J. *et al.* Methods for analyzing deep sequencing expression data: constructing the human and mouse promoterome with deepcage data. *Genome biology* **10**, R79 (2009).
26. Mehlkopf, A., Korbee, D., Tiggelman, T. & Freeman, R. Sources of t1 noise in two-dimensional nmr. *Journal of Magnetic Resonance* (1969) **58**, 315–323 (1984).
27. Van Batenburg, M. F., Coulier, L., van Eeuwijk, F., Smilde, A. K. & Westerhuis, J. A. New figures of merit for comprehensive functional genomics data: the metabolomics case. *Analytical chemistry* **83**, 3267–3274 (2011).
28. Rocke, D. M. & Lorenzato, S. A two-component model for measurement error in analytical chemistry. *Technometrics* **37**, 176–184 (1995).
29. Karakach, T. K., Wentzell, P. D. & Walter, J. A. Characterization of the measurement error structure in 1D ^1H NMR data for metabolomics studies. *Analytica Chimica Acta* **636**, 163–174 (2009).
30. Pearson, K. & Lee, A. On the laws of inheritance in man: I. Inheritance of physical characters. *Biometrika* **2**, 357–462 (1903).
31. Winne, P. H. & Belfry, M. J. Interpretive problems when correcting for attenuation. *Journal of Educational Measurement* 125–134 (1982).
32. Liu, K., Stamler, J., Dyer, A., McKeever, J. & McKeever, P. Statistical methods to assess and minimize the role of intra-individual variability in obscuring the relationship between dietary lipids and serum cholesterol. *Journal of chronic diseases* **31**, 399–418 (1978).
33. McCulloch, C. E. & Neuhaus, J. M. Generalized linear mixed models. *Encyclopedia of biostatistics* **4** (2005).
34. Verbeke, G. & Molenberghs, G. *Linear mixed models for longitudinal data* (Springer Science & Business Media, 2009).
35. Leger, M. N., Vega-Montoto, L. & Wentzell, P. D. Methods for systematic investigation of measurement error covariance matrices. *Chemometrics and Intelligent Laboratory Systems* **77**, 181–205 (2005).
36. Wentzell, P. D., Cleary, C. S. & Kompany-Zareh, M. Improved modeling of multivariate measurement errors based on the wishart distribution. *Analytica chimica acta* **959**, 1–14 (2017).
37. Comrey, A. L. & Lee, H. B. *A first course in factor analysis* (Psychology press, 2013).
38. Day, N. *et al.* Correlated measurement error—implications for nutritional epidemiology. *International Journal of Epidemiology* **33**, 1373–1381 (2004).
39. Pereira, V., Waxman, D. & Eyre-Walker, A. A problem with the correlation coefficient as a measure of gene expression divergence. *Genetics* **183**, 1597–1600 (2009).

40. Reynier, F. *et al.* Importance of correlation between gene expression levels: application to the type i interferon signature in rheumatoid arthritis. *PLoS one* **6**, e24828 (2011).
41. Springer, M. D. *The algebra of random variables* (Wiley and Sons, 1979).
42. Bishara, A. J. & Hittner, J. B. Reducing bias and error in the correlation coefficient due to nonnormality. *Educational and psychological measurement* **75**, 785–804 (2015).
43. Kowalski, C. J. On the effects of non-normality on the distribution of the sample product-moment correlation coefficient. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **21**, 1–12 (1972).
44. Magnusson, P. K. *et al.* The Swedish twin registry: establishment of a biobank and other recent developments. *Twin Research and Human Genetics* **16**, 317–329 (2013).
45. Haug, K. *et al.* Metabolights—an open-access general-purpose repository for metabolomics studies and associated meta-data. *Nucleic acids research* **41**, D781–D786 (2012).
46. Ganna, A. *et al.* Large-scale non-targeted metabolomic profiling in three human population-based studies. *Metabolomics* **12**, 4 (2016).

Acknowledgements

This work has been partially funded by The Netherlands Organization for Health Research and Development (ZonMW) through the PERMIT project (Personalized Medicine in Infections: from Systems Biomedicine and Immunometabolism to Precision Diagnosis and Stratification Permitting Individualized Therapies, project contract number 456008002) under the PerMed Joint Transnational call JTC 2018 (Research projects on personalised medicine - smart combination of pre-clinical and clinical research with data and ICT solutions). The authors acknowledge Peter Wentzell (Halifax, Canada) for kindly making available the NMR data set.

Author contributions

E.S. and A.S. conceived the study and performed theoretical calculations. E.S., M.H. and A.S. analysed and interpreted the results. E.S. and M.H. performed simulations. E.S., M.H. and A.S. wrote, reviewed and approved the manuscript in its final form.

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to E.S.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020