

The reconciliation of multiple conflicting estimates: Entropy-based and axiomatic approaches

Rutgers University has made this article freely available. Please share how this access benefits you.
Your story matters. [\[https://rucore.libraries.rutgers.edu/rutgers-lib/58264/story/\]](https://rucore.libraries.rutgers.edu/rutgers-lib/58264/story/)

This work is the **VERSION OF RECORD (VoR)**

This is the fixed version of an article made available by an organization that acts as a publisher by formally and exclusively declaring the article "published". If it is an "early release" article (formally identified as being published even before the compilation of a volume issue and assignment of associated metadata), it is citable via some permanent identifier(s), and final copy-editing, proof corrections, layout, and typesetting have been applied.

Citation to Publisher Rodrigues, João F.D. & Lahr, Michael L. (2018). The reconciliation of multiple conflicting estimates: Entropy-based and axiomatic approaches. *Entropy* 20. <https://doi.org/10.3390/e20110815>.

Citation to *this* Version: Rodrigues, João F.D. & Lahr, Michael L. (2018). The reconciliation of multiple conflicting estimates: Entropy-based and axiomatic approaches. *Entropy* 20. Retrieved from <http://dx.doi.org/doi:10.7282/T37P92XH>.

Terms of Use: Copyright for scholarly resources published in RUcore is retained by the copyright holder. By virtue of its appearance in this open access medium, you are free to use this resource, with proper attribution, in educational and other non-commercial settings. Other uses, such as reproduction or republication, may require the permission of the copyright holder.

Article begins on next page

Article

The reconciliation of multiple conflicting estimates: entropy-based and axiomatic approaches

João F. D. Rodrigues ^{1*} and Michael L. Lahr ²

¹ Institute of Environmental Sciences CML, Leiden University, Einsteinweg 2, 2333 CC Leiden, the Netherlands

² Edward J. Bloustein School of Planning & Public Policy, Rutgers, The State University of New Jersey, 33 Livingston Avenue, New Brunswick, NJ 08901-1982, USA

* Correspondence: E-mail: j.rodrigues@cml.leidenuniv.nl; Tel.: +31 645 062 746

Academic Editor: name

Version October 18, 2018 submitted to Entropy

Abstract: When working with economic accounts it may occur that multiple estimates of a single datum exist, with different degrees of uncertainty or data quality. This paper addresses the problem of defining a method that can reconcile conflicting estimates, given best guess and uncertainty values. We proceeded from first principles, using two different routes. First, under an entropy-based approach, the data reconciliation problem is addressed as a particular case of a wider data balancing problem, and an alternative setting is found in which the multiple estimates are replaced by a single one. Afterwards, under an axiomatic approach, a set of properties is defined, which characterizes the ideal data reconciliation method. Under both approaches, the conclusion is that the formula for the reconciliation of best guesses is a weighted arithmetic average, with the inverse of uncertainties as weights, and that the formula for the reconciliation of uncertainties is a harmonic average.

Keywords: Uncertainty modelling; economic accounts; conflicting estimates; entropy-based approach; axiomatic approach.

1. Introduction

With improvements in information technology, the world has become more unified and interconnected. Information is now typically shared quickly and easily from all over the globe, such that barriers formed by linguistic and geographic boundaries essentially have been torn down. This has enabled people from disparate cultures and backgrounds to share ideas and information. One outcome of this regime change has been a boosting of the perceived benefits of statistical information. While some benefits of such statistical information have been known since at least Quetelet's (1835) tome on so-called "social physics" was published, today's massive socio-economic statistical repositories in Europe, North America, and East Asia are enabling a data revolution of sorts. Indeed, the fields of data mining and data analytics are fast becoming important fields of academic study. Mirroring the rise of data availability and the nature of some of the data itself, the term "big data" has been coined [1] to refer to the extremely voluminous and complex data sets that require specialized processing application software to deal with them.

The most prominent stewards of socio-economic data are government statistical agencies, which focus on producing and disseminating data products secured via surveys (for example the American Community Survey), censuses (such as Japan's 2015 Population Census), and administrative procedures (like information needed to get an academic promotion in Spain). As a result, data storage is now ubiquitously electronic, replicated offsite to guard against storage failure, and measured in petabytes. Electronic storage enables low-cost dissemination of data. It also facilitates the integration

32 of records across disparate databases – for example, into a system of national accounts, which is what
33 countries use to generate their estimates of gross domestic product, as suggested by the United Nations.
34 Both lead to concerns about confidentiality of data and how it can be protected [2].

35 Our point in broaching the above is that data producers, disseminators, and users alike can
36 run into the problem of having access to multiple estimates for a single quantity of interest. In the
37 particular experience of the authors, which motivated the present study, these multiple estimates are a
38 consequence of non-disclosure by a statistical office in order to ensure confidentiality. Hence we act as
39 what Duncan *et al.* [2] called “data snooper”. For concreteness, in such instances we are interested in
40 obtaining number of employees at county level from the U.S. Bureau of Statistics’s Quarterly Census
41 of Employment and Wages (QCEW) data and U.S. Bureau of the Census’s County Business Patterns.
42 In both datasets these figures are suppressed for selected sectors in some counties and even states. But
43 information is provided for larger spatial and sectoral units, so it is possible to use this higher-level
44 information to obtain multiple estimates of the quantities of interest. It is common for official statistical
45 data to have a hierarchical structure so this problem is quite general. Garfinkel *et al.* [3] note that the
46 increasing ability of data snoopers is making ever more data stewards reluctant to provide certain data
47 products because they are finding it increasingly difficult to ensure confidentiality to the agents from
48 whom they obtain the data. This is despite some use of noise as a disclosure limitation [4].

49 In this paper, we focus on the problem of combining such multiple estimates into a single value.
50 In the case of economic accounts Miller and Blair [5, pp. 384-6] have called it “the reconciliation issue”.
51 The reconciliation issue considered here should not be confused with the more general problem of
52 data balancing, in which a set of multiple data points need to satisfy a set of constraints: that problem
53 is addressed in other studies, such as Kruihof [6], Stone *et al.* [7], Byron [8], Van der Ploeg [9], Lahr
54 and Mesnard [10], Chen [11]. General solutions to confidentiality disclosure or data censoring issues
55 are provided by [12] and [13]. Herein we set out to assist current and future data snoopers and miners,
56 by identifying what a data reconciliation method should be from first principles when a fairly general
57 formulation of the reconciliation constraints is possible.

58 In particular, we consider that the multiple estimates for a particular datum can be characterized
59 by a best guess and uncertainty. If we interpret each estimate as a random variable with an underlying
60 probability distribution, the best guess is the expected value and the uncertainty is standard deviation.
61 In the case of multiple data sources, the conflict enabling the multiple estimates is self-evident. When
62 numbers are published with some data censored and for which estimates can be obtained using partial
63 information [14], the conflict can arise from a higher (or lower) hierarchical spatial or sectoral level (e.g.,
64 average employee number if the number of establishments is available). To the best of our knowledge
65 no first-principle approach to this problem has yet been published, although more heuristic approaches
66 can be found in Bourque *et al.* [15], Miernyk *et al.* [16], Jensen and McGaurr [17], Gerking [1819], Weale
67 [20], Boomsma and Oosterhaven [21], Rassier *et al.* [22]. We tackle the same problem from two different
68 angles.

69 Using concepts and techniques from Bayesian inference [23] and in particular the minimum
70 cross-entropy method [24], we first address the problem of data reconciliation as a particular case of
71 more general data balancing [25]. That is, we consider there are two or more initial estimates for a
72 particular datum, but this datum is itself embedded in a set of constraints connecting it to other data
73 that are potentially unbalanced. We look for simplifications of the general setting under which this
74 original problem can be transformed into another balancing problem where the multiple estimates are
75 replaced by a single one. We prove that, if the initial uncertainty estimates are close to one another, the
76 data reconciliation method of best guesses is a weighted arithmetic average and the data reconciliation
77 method of uncertainties is a harmonic average.

78 Afterwards we address the same problem from an axiomatic perspective, laying out the desirable
79 properties of a data reconciliation method. Such an approach has roots in different fields, from table
80 deflation [26] and supply-use transformations [27] to environmental responsibility [28]. It turns out
81 that the canonical data reconciliation method, i.e., the one that satisfies all required properties, is none

other than a suitable generalization of the entropy-based method derived earlier. That generalization centers on the introduction of the number of previously combined priors and a ranking of estimates by their relative quality.

2. Entropy-based approach

2.1. Basic concepts

Bayesian inference was first developed by Laplace [29] and later expanded by others, such as Jeffreys [30], Jaynes [31] and Jaynes [23]. According to the Bayesian paradigm, a probability is a degree of belief about the likelihood of an event, and should reflect all relevant available information about that event. According to Weise and Woger [32], if an empirical quantity is subject to measurement errors, it must be described by a random variable, whose expectation is the best guess and whose standard-deviation is the uncertainty estimate.

More formally, a *prior* datum θ_i is characterized by a probability distribution $\pi(q_i)$, which expresses the degree of belief that the datum takes realization q_i . The *best guess* is $\mu_i = E[\theta_i]$ and the *uncertainty* is $\sigma_i = \sqrt{\text{Var}[\theta_i]}$. When multiple data are considered, e.g., θ_i and θ_j , it is necessary to introduce the *correlation* between them, $\rho_{ij} = \text{Cov}[\theta_i, \theta_j] / \sigma_i \sigma_j$. Rodrigues [33] further provides a series of rules to determine the properties of a strictly positive prior datum, using the maximum-entropy principle [34].

The type of data we are interested in are connected to one another through *accounting identities* of the form:

$$\theta_0 = \sum_{i=1}^n \theta_i,$$

where θ_0 is an *aggregate* datum and the θ_i 's are *disaggregate* data. If the set of data is arranged in a vector θ of length n_T , the set of n_K accounting identities can be defined through a concordance matrix \mathbf{G} , where, for a given accounting identity i , $G_{ij} = 1$ if θ_j is a disaggregate datum, $G_{ij} = -1$ if θ_j is an aggregate datum and $G_{ij} = 0$ otherwise.

If the prior configuration is unbalanced, then $\mathbf{G}\theta \neq \mathbf{0}$, where $\mathbf{0}$ is a vector of zeros. Rodrigues [25] derives an analytical solution and a series of approximations that, given a concordance matrix and prior configuration, provide a *posterior* configuration, \mathbf{t} , such that $\mathbf{G}\mathbf{t} = \mathbf{0}$. **The notational convention used here is that Greek letters refer to priors while Latin cognates will refer to posteriors, i.e., m_i , s_i and r_{ij} are, respectively, the best guess and uncertainty of t_i and correlation between t_i and t_j .**

2.2. Problem formulation

We are now in position to formulate the data reconciliation problem. Given *initial priors* θ' and θ'' and a system with $n_T + 1$ numerical data $\{\theta_1, \dots, \theta_{n_T-1}, \theta', \theta''\}$, and $n_K + 1$ accounting identities, where accounting identity $n_K + 1$ takes the form $\theta' = \theta''$, our goal is to determine the *final prior* θ , in a new system with n_T numerical data, $\{\theta_1, \dots, \theta_{n_T-1}, \theta\}$ and the n_K first accounting identities of the original system, in which the posteriors $\{t_1, \dots, t_{n_T-1}\}$ are identical in both data balancing problems, and $t = t' = t''$. **Conceptually, we are approaching data reconciliation as a form of preliminary data balancing, as illustrated in Figure 1. The conflicting estimates are initial priors of the same datum, and the reconciled value is a final prior. Note the following notational convention: while other variables (and their properties) are denoted with subscripts, initial priors/posteriors (and their properties) are denoted with one (') or two (") primes, and the final prior/posterior is denoted with neither subscripts nor primes.**

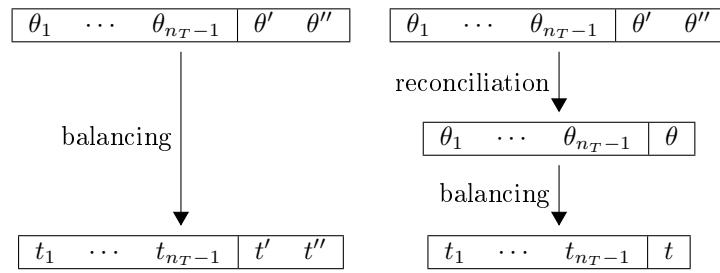


Figure 1. On the left-hand side balancing in a single step, with multiple initial estimates (priors) of the same datum, θ' and θ'' , balanced to the same quantity (posterior), $t' = t''$. On the right-hand side balancing in two steps: first the reconciliation procedure combines the multiple initial estimates (initial priors), θ' and θ'' , into a final prior, θ ; afterwards the full system is balanced, leading to posterior t . We impose that the result from both procedures is the same, $t = t' = t''$.

122 Three situations emerge: either the datum to be reconciled is only a disaggregate datum; it is
 123 only an aggregate datum; or it is both a disaggregate and an aggregate datum, in different accounting
 124 identities. We will deal with the three cases separately.

125 We now present simple systems to illustrate the three possible cases. As a benchmark consider a
 126 tabular system (i.e., with data organized in rows and columns) with no multiple estimates consisting
 127 of a 2×3 table \mathbf{A} with row sums \mathbf{b} and columns sums \mathbf{c} . Furthermore, consider that the sum of both
 128 \mathbf{b} and \mathbf{c} is known as d . If \mathbf{i} is a vector of ones of appropriate length, all vectors are in column format
 129 by default, and prime (') adjoined to a matrix or vector denotes transpose, then the previous set of
 130 constraints means that:

$$\begin{aligned} \mathbf{A}\mathbf{i} &= \mathbf{b}; \\ \mathbf{A}'\mathbf{i} &= \mathbf{c}; \\ \mathbf{b}'\mathbf{i} &= d; \\ \mathbf{c}'\mathbf{i} &= d. \end{aligned}$$

131
 132 The vectorized form of this system and the concordance table is presented in Table 1. In the
 133 baseline system there is a total of twelve variables (columns of the concordance matrix \mathbf{G}) and seven
 134 constraints (rows thereof). The first six variables are disaggregate values (corresponding to the initial \mathbf{A}
 135 matrix), the following five are mixed (row and column sums \mathbf{b} and \mathbf{c}), and the last one is an aggregate
 136 datum (d). The first two constraints (rows of \mathbf{G}) are the row sums of \mathbf{A} , the following three are its
 137 columns sums, and the last two are the sums of \mathbf{b} and \mathbf{c} . To understand how \mathbf{G} is constructed let us
 138 consider the first constraint, which is the row sum of \mathbf{A} . Formally, this is:

$$A_{11} + A_{12} + A_{13} - b_1 = 0,$$

139 hence in the first row of \mathbf{G} the entries corresponding to the columns of A_{11} , A_{12} and A_{13} have 1s, the
 140 entry corresponding to the column of b_1 has -1 and all entries are zero.

141 We are now in position to formalize the three situations of multiple estimates of a single datum as
 142 variants of Table 1 in which an additional row and column has been added to \mathbf{G} .

143 The case of disaggregate datum occurs if the datum for which multiple estimates exist is an
 144 interior point, which for concreteness we consider to be element A_{23} : the set of constraints is shown in
 145 Table 2. As an illustration of the case of there being two estimates of an aggregate datum consider it to

Table 1. Prior vector and concordance matrix, with no multiple estimates.

θ'	A_{11}	A_{12}	A_{13}	A_{21}	A_{22}	A_{23}	b_1	b_2	c_1	c_2	c_3	d
G	1	1	1	0	0	0	-1	0	0	0	0	0
	0	0	0	1	1	1	0	-1	0	0	0	0
	1	0	0	1	0	0	0	0	-1	0	0	0
	0	1	0	0	1	0	0	0	0	-1	0	0
	0	0	1	0	0	1	0	0	0	0	-1	0
	0	0	0	0	0	0	1	1	0	0	0	-1
	0	0	0	0	0	0	0	0	1	1	1	-1

146 be d : the set of constraints is shown in Table 3. Finally, consider as example of an element that is both
 147 aggregate and disaggregate that of b_1 : the set of constraints is shown in Table 4.

Table 2. Prior vector and concordance matrix, with multiple estimates of A_{23} .

θ'	A_{11}	A_{12}	A_{13}	A_{21}	A_{22}	A'_{23}	b_1	b_2	c_1	c_2	c_3	d	A''_{23}
G	1	1	1	0	0	0	-1	0	0	0	0	0	0
	0	0	0	1	1	1	0	-1	0	0	0	0	0
	1	0	0	1	0	0	0	0	-1	0	0	0	0
	0	1	0	0	1	0	0	0	0	-1	0	0	0
	0	0	1	0	0	0	0	0	0	0	-1	0	1
	0	0	0	0	0	0	1	1	0	0	0	-1	0
	0	0	0	0	0	0	0	0	1	1	1	-1	0
	0	0	0	0	0	1	0	0	0	0	0	0	-1

Table 3. Prior vector and concordance matrix, with multiple estimates of d .

θ'	A_{11}	A_{12}	A_{13}	A_{21}	A_{22}	A_{23}	b_1	b_2	c_1	c_2	c_3	d'	d''
G	1	1	1	0	0	0	-1	0	0	0	0	0	0
	0	0	0	1	1	1	0	-1	0	0	0	0	0
	1	0	0	1	0	0	0	0	-1	0	0	0	0
	0	1	0	0	1	0	0	0	0	-1	0	0	0
	0	0	1	0	0	1	0	0	0	0	-1	0	0
	0	0	0	0	0	0	1	1	0	0	0	-1	0
	0	0	0	0	0	0	0	0	1	1	1	0	-1
	0	0	0	0	0	0	0	0	0	0	0	1	-1

Table 4. Prior vector and concordance matrix, with multiple estimates of b_1 .

θ'	A_{11}	A_{12}	A_{13}	A_{21}	A_{22}	A_{23}	b'_1	b_2	c_1	c_2	c_3	d	b''_1
G	1	1	1	0	0	0	-1	0	0	0	0	0	0
	0	0	0	1	1	1	0	-1	0	0	0	0	0
	1	0	0	1	0	0	0	0	-1	0	0	0	0
	0	1	0	0	1	0	0	0	0	-1	0	0	0
	0	0	1	0	0	1	0	0	0	0	-1	0	0
	0	0	0	0	0	0	0	1	0	0	0	-1	1
	0	0	0	0	0	0	0	0	1	1	1	-1	0
	0	0	0	0	0	0	1	0	0	0	0	0	-1

148 It is perhaps instructive to describe how the reconciliation problems differ from the features of
 149 the baseline system. The three variants of the baseline are constructed by adding a single variable,
 150 the conflicting estimate, which by convenience is always appended to the original system. It is also
 151 necessary to add an extra constraint, connecting the two conflicting estimates. Finally, the baseline
 152 system is also changed so that in one of the original occurrences of the datum to be reconciled is the
 153 first conflicting estimate and the second occurrence is the other conflicting estimate.

154 Note that in this simple example there are only two constraints affecting each datum, but that
 155 naturally is not generally the case. The number of constraints per datum is arbitrary and can be either
 156 one or larger than two. An example of what this system might represent is employment count by
 157 region and sector, with an extra dimension being type of ownership (private or local, state, or federal
 158 government), as reported in the QCEW database.

159 2.3. From balancing to reconciliation

160 Rodrigues [25] shows that if the posterior configuration is balanced, then its first- and
 161 second-moment constraints are:

$$\mathbf{0} = \mathbf{G}\mathbf{m}; \quad (1)$$

$$\mathbf{0} = \text{diag}(\mathbf{G}\mathbf{S}|\mathbf{G}'|'), \quad (2)$$

162 where \mathbf{m} and \mathbf{S} are the posterior best-guess vector and covariance matrix, and the latter is defined as
 163 $\mathbf{S} = \hat{\mathbf{s}}\mathbf{R}\hat{\mathbf{s}}$, where \mathbf{s} is the vector of posterior uncertainties and \mathbf{R} is the vector of posterior correlations,
 164 and $\hat{\cdot}$ denotes diagonal matrix. Likewise $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the prior best guess vector and covariance matrix,
 165 and the latter is defined as $\boldsymbol{\Sigma} = \hat{\boldsymbol{\sigma}}\mathbf{P}\hat{\boldsymbol{\sigma}}$.

166 The analytical solution of the data-balancing problem is:

$$\tilde{\mathbf{S}}^{-1} = \tilde{\boldsymbol{\Sigma}}^{-1} + \mathbf{G}'\hat{\boldsymbol{\beta}}|\mathbf{G}|; \quad (3)$$

$$\tilde{\mathbf{S}}^{-1}\tilde{\mathbf{m}} = \tilde{\boldsymbol{\Sigma}}^{-1}\tilde{\boldsymbol{\mu}} + \mathbf{G}'\boldsymbol{\alpha}. \quad (4)$$

167 Notice that Equations 3-4 contain symbols adjoined with \sim (which we refer to as Gaussian
 168 parameters) while Equations 1-2 do not. The connection between the Gaussian parameters and
 169 the corresponding observable quantities is described in Rodrigues [25]: when relative uncertainty,
 170 σ_j/μ_j or s_j/m_j , is low, then the Gaussian parameter and the observable are identical. When relative
 171 uncertainty is high, the best guess Gaussian parameter tends to $-\infty$ and the uncertainty Gaussian
 172 parameter tends to ∞ , in such a way that if relative uncertainty is unitary, $-\tilde{\mu}_j/\tilde{\sigma}_j^2 = 1/\mu_j = 1/\sigma_j$
 173 and $-\tilde{m}_j/\tilde{s}_j^2 = 1/m_j = 1/s_j$. There is no closed-form expression between observables and Gaussian
 174 parameters in the multivariate case.

175 If both the prior uncertainty of aggregate data and initial prior correlations are high, we obtain a
 176 simplified weighted least-squares (WLS) method in which the weights are prior uncertainties:

$$\mathbf{m} = \boldsymbol{\mu} + \hat{\boldsymbol{\sigma}}\mathbf{G}'\boldsymbol{\alpha}, \quad (5)$$

177 and posterior correlations are set by considering that relative uncertainty is constant, $\mathbf{s} = \mathbf{m} \odot \boldsymbol{\sigma} \oslash \boldsymbol{\mu}$,
 178 where \odot and \oslash are Hadamard (or entrywise) product and division, and the update takes place in
 179 small steps.

180 This WLS method is a generalization of the [standard biproportional balancing method \(RAS\)](#)
 181 for arbitrary structure and uncertainty data [25]. However, it is in a way too simple for the data
 182 reconciliation problem, because it keeps relative uncertainty constant. In the data reconciliation
 183 problem this assumption is untenable, whenever the relative uncertainty of the initial priors differs.

184 Thus, we now look for a simplification of the general solution (Equations 3-4) that is still feasible
 185 and that allows both for best guess and uncertainty reconciliation. Let us consider that correlations
 186 change little from prior to posterior, so that only uncertainties are adjusted. Equations 3-4 become:

$$\begin{aligned}\hat{\mathbf{s}}^{-1}\mathbf{R}^{-1}\hat{\mathbf{s}}^{-1} &= \hat{\sigma}^{-1}\mathbf{P}^{-1}\sigma^{-1} + \mathbf{G}'\boldsymbol{\beta}; \\ \hat{\mathbf{s}}^{-1}\mathbf{R}^{-1}\hat{\mathbf{s}}^{-1}\mathbf{m} &= \hat{\sigma}^{-1}\mathbf{P}^{-1}\hat{\sigma}^{-1}\boldsymbol{\mu} + \mathbf{G}'\boldsymbol{\alpha},\end{aligned}$$

187 where we dropped the \sim , meaning that all variables are observables. If correlations are not adjusted,
188 then $\mathbf{R} = \mathbf{P}$, and if variances change little $\mathbf{s} \simeq \sigma$. The previous expressions become:

$$\begin{aligned}\mathbf{s}^{-1} &= \sigma^{-1} + \mathbf{P}\hat{\sigma}\mathbf{G}'\boldsymbol{\beta}; \\ \hat{\mathbf{s}}^{-1}\mathbf{m} &= \hat{\sigma}^{-1}\boldsymbol{\mu} + \mathbf{P}\hat{\sigma}\mathbf{G}'\boldsymbol{\alpha}.\end{aligned}$$

189 For convenience, consider now that a datum corresponding to entry (i, j) in the tabular matrix is
190 t_{ij} , while the sums of row or column i is t_i , and the Lagrange parameters of a row sum or column sum
191 are adjoined with superscript R or C . For a particular entry, the previous matrix equation reads:

$$\begin{aligned}\frac{1}{s_{ij}} &= \frac{1}{\sigma_{ij}} + \sigma_i^R \beta_i^R + \sigma_j^C \beta_j^C + \sum_{k \neq i, j} (\sigma_{ik} \beta_k^C + \sigma_{kj} \beta_k^R); \\ \frac{1}{s_i} &= \frac{1}{\sigma_i} + \sigma_i (\beta_i^R + \beta_i^C),\end{aligned}$$

192 where $\sigma_i^R = \sum_j \sigma_{ij}$ and $\sigma_i^C = \sum_j \sigma_{ji}$. If the adjustment from prior to posterior is small, then $\sigma_i^R \simeq \sigma_i^C \simeq \sigma_i$.
193 If $\beta_i^* = \sigma_i \beta_i^R$, $\sigma_i \gg \sigma_{ij}$ and $\sigma_i \gg \sigma_{ji}$, then the previous expression matrix expressions simplify to:

$$\mathbf{s}^{-1} = \sigma^{-1} + \mathbf{G}'\boldsymbol{\beta}^*; \quad (6)$$

$$\hat{\mathbf{s}}^{-1}\mathbf{m} = \hat{\sigma}^{-1}\boldsymbol{\mu} + \mathbf{G}'\boldsymbol{\alpha}^*, \quad (7)$$

194 where the derivation of Equation 7 follows along identical lines to that Equation 6. We now use
195 these expressions to obtain a tentative solution of the data reconciliation problem, even though they
196 were derived under rather strict assumptions.

197 2.4. A tentative solution

198 We now examine the implications of applying Equations 6-7 to different to the different data
199 reconciliation configurations described in Section 2.2: multiple estimates of (a) an aggregate datum; (b)
200 a disaggregate datum; and (c) a datum that is both aggregate and disaggregate. We shall see that the
201 same expression applies to all these problems.

202 For clarity, the analysis is carried out using scalar expressions, and, for brevity, only to the case
203 of two constraints per datum. The strategy of the proof is the same for all configurations: to derive
204 constraints connecting prior and posterior in the original problem and in a modified problem in which
205 there is only a single datum where originally there were the conflicting estimates.

206 2.4.1. Aggregate datum

207 Consider that there are two initial priors of a datum, θ'_0 and θ''_0 and that the datum is involved
208 in two accounting identities, the first summing over elements 1 to n' and the second summing over
209 $n' + 1$ to n'' :

$$\begin{aligned}
 t'_0 &= \sum_{i=1}^{n'} t_i; \\
 t''_0 &= \sum_{i=n'+1}^{n'+n''} t_i; \\
 t'_0 &= t''_0,
 \end{aligned}$$

210 where each t_i , for $i > 0$, can be affected by other accounting identities. The Lagrange parameters
 211 associated with these three expressions in Equation 6 are denoted, respectively, by β'_0 , β''_0 and β_0 . We
 212 wish to determine a final prior θ_0 , such that:

$$\begin{aligned}
 t_0 &= \sum_{i=1}^{n'} t_i; \\
 t_0 &= \sum_{i=n'+1}^{n'+n''} t_i.
 \end{aligned}$$

213 Equation 6 reads, for the original problem:

$$\begin{aligned}
 \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta'_0 + \dots \quad \text{if } 1 \leq i \leq n'; \\
 \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta''_0 + \dots \quad \text{if } n' + 1 \leq i \leq n''; \\
 \frac{1}{s'_0} &= \frac{1}{\sigma'_0} - \beta'_0 + \beta_0; \\
 \frac{1}{s''_0} &= \frac{1}{\sigma''_0} - \beta''_0 - \beta_0,
 \end{aligned}$$

214 where \dots refers to other Lagrange parameters. And in the modified problem:

$$\begin{aligned}
 \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta'_0 + \dots \quad \text{if } 1 \leq i \leq n'; \\
 \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta''_0 + \dots \quad \text{if } n' + 1 \leq i \leq n''; \\
 \frac{1}{s_0} &= \frac{1}{\sigma_0} - \beta'_0 - \beta''_0.
 \end{aligned}$$

215 Notice that for every datum $i > 0$ the original and modified problem are identical. Because the
 216 posteriors of the aggregate datum are all identical, $s'_0 = s''_0 = s_0$, we can write:

$$\begin{aligned}
 2\frac{1}{s_0} &= 2\left(\frac{1}{\sigma_0} - \beta'_0 - \beta''_0\right) \\
 &= \frac{1}{\sigma'_0} + \frac{1}{\sigma''_0} - \beta'_0 - \beta''_0 + \beta_0 - \beta_0.
 \end{aligned}$$

217 A similar expression can be obtained from Equation 7 for the final prior best guess, leading to the
 218 solution:

$$\frac{1}{\sigma_0} = \frac{1}{2} \left(\frac{1}{\sigma_0'} + \frac{1}{\sigma_0''} \right);$$

$$\frac{\mu_0}{\sigma_0} = \frac{1}{2} \left(\frac{\mu_0'}{\sigma_0'} + \frac{\mu_0''}{\sigma_0''} \right).$$

219 Thus, both the final prior of the absolute uncertainty, σ , and the relative uncertainty, σ/μ , are
 220 obtained as the harmonic average of the initial prior absolute and relative uncertainties.

221 2.4.2. Disaggregate datum

222 Consider now that there are two initial priors of an interior point, θ_1' and θ_1'' , which is affected by
 223 two accounting identities, such that the posteriors satisfy:

$$t_0' = t_1' + \sum_{i=2}^{n'} t_i;$$

$$t_0'' = t_1'' + \sum_{i=n'+1}^{n'+n''} t_i;$$

$$t_1' = t_1''.$$

224 The Lagrange parameters associated with these three expressions are, as before, β_0' , β_0'' and β_0 .
 225 We wish to determine a final prior θ_1 , such that:

$$t_0' = t_1 + \sum_{i=2}^{n'} t_i;$$

$$t_0'' = t_1 + \sum_{i=n'+1}^{n'+n''} t_i.$$

226 Equation 6 reads, for the original problem:

$$\frac{1}{s_i} = \frac{1}{\sigma_i} + \beta_0' + \dots \quad \text{if } 2 \leq i \leq n';$$

$$\frac{1}{s_i} = \frac{1}{\sigma_i} + \beta_0'' + \dots \quad \text{if } n' + 1 \leq i \leq n'';$$

$$\frac{1}{s_0'} = \frac{1}{\sigma_0'} - \beta_0' + \dots;$$

$$\frac{1}{s_0''} = \frac{1}{\sigma_0''} - \beta_0'' + \dots;$$

$$\frac{1}{s_1'} = \frac{1}{\sigma_1'} + \beta_0' + \beta_1;$$

$$\frac{1}{s_1''} = \frac{1}{\sigma_1''} + \beta_0'' - \beta_1,$$

227 and in the modified problem:

$$\begin{aligned}
\frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta'_0 + \dots & \text{if } 2 \leq i \leq n'; \\
\frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta''_0 + \dots & \text{if } n' + 1 \leq i \leq n''; \\
\frac{1}{s'_0} &= \frac{1}{\sigma'_0} - \beta'_0 + \dots; \\
\frac{1}{s''_0} &= \frac{1}{\sigma''_0} - \beta''_0 + \dots; \\
\frac{1}{s_1} &= \frac{1}{\sigma_1} + \beta'_0 + \beta''_0.
\end{aligned}$$

228 As before, the data for which there are no conflicting estimates (t'_0, t''_0 and t_i with $i > 1$) are subject
229 to the same set of constraints in the original and in the modified problem. Because the posteriors of the
230 disaggregate datum are all identical, $s'_1 = s''_1 = s_1$, we can write:

$$\begin{aligned}
2\frac{1}{s_1} &= 2\left(\frac{1}{\sigma_1} + \beta'_0 + \beta''_0\right) \\
&= \frac{1}{\sigma'_1} + \frac{1}{\sigma''_1} + \beta'_0 + \beta''_0 + \beta_1 - \beta_1.
\end{aligned}$$

231 At this stage it becomes clear that we will encounter exactly the same solution as in the case of an
232 aggregate datum:

$$\begin{aligned}
\frac{1}{\sigma_1} &= \frac{1}{2} \left(\frac{1}{\sigma'_1} + \frac{1}{\sigma''_1} \right); \\
\frac{\mu_1}{\sigma_1} &= \frac{1}{2} \left(\frac{\mu'_1}{\sigma'_1} + \frac{\mu''_1}{\sigma''_1} \right).
\end{aligned}$$

233 2.4.3. Mixed datum

234 Consider now that there are two initial priors, θ'_1 and θ''_1 , of a datum that is both aggregate and
235 disaggregate, in different accounting identities, and whose posteriors satisfy:

$$\begin{aligned}
t_0 &= t'_1 + \sum_{i=2}^{n'} t_i; \\
t''_1 &= \sum_{i=n'+1}^{n''} t_i; \\
t'_1 &= t''_1.
\end{aligned}$$

236 As before the Lagrange parameters are denoted as β'_0, β''_0 and β_1 . We wish to determine a final
237 prior θ_1 , such that:

$$\begin{aligned}
t_0 &= t_1 + \sum_{i=2}^{n'} t_i; \\
t_1 &= \sum_{i=n'+1}^{n''} t_i.
\end{aligned}$$

238 Equation 6 reads, for the original problem:

$$\begin{aligned}\frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta'_0 + \dots & \text{if } 2 \leq i \leq n'; \\ \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta''_0 + \dots & \text{if } n' + 1 \leq i \leq n''; \\ \frac{1}{s_0} &= \frac{1}{\sigma_0} - \beta'_0 + \dots; \\ \frac{1}{s'_1} &= \frac{1}{\sigma'_1} + \beta'_0 + \beta_1; \\ \frac{1}{s''_1} &= \frac{1}{\sigma''_1} - \beta''_0 - \beta_1,\end{aligned}$$

239 and in the modified problem:

$$\begin{aligned}\frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta'_0 + \dots & \text{if } 2 \leq i \leq n'; \\ \frac{1}{s_i} &= \frac{1}{\sigma_i} + \beta''_0 + \dots & \text{if } n' + 1 \leq i \leq n''; \\ \frac{1}{s_0} &= \frac{1}{\sigma_0} - \beta'_0 + \dots; \\ \frac{1}{s_1} &= \frac{1}{\sigma_1} + \beta'_0 - \beta''_0.\end{aligned}$$

240 As has become routine, for datum 0 and for every datum $i > 1$ the original and modified problem
241 are identical. Because $s'_1 = s''_1 = s_1$, we can write:

$$\begin{aligned}2\frac{1}{s_1} &= 2\left(\frac{1}{\sigma_1} + \beta'_0 - \beta''_0\right) \\ &= \frac{1}{\sigma'_1} + \frac{1}{\sigma''_1} + \beta'_0 - \beta''_0 + \beta_1 - \beta_1.\end{aligned}$$

242 Thus, it is clear that the solution is again identical.

243 3. Axiomatic approach

244 3.1. Axiomatic formulation

245 In Section 2 we obtained a data reconciliation algorithm from first principles, as an operation of
246 data balancing under a particular structure. However, we can also reason about the data reconciliation
247 algorithm in terms of its properties, i.e., we will not determine what it is, but what it ought to be.

248 If θ' and θ'' are two initial priors, the data reconciliation algorithm is a function $f(\cdot)$ that generates
249 a final prior $\theta = f(\theta', \theta'')$, where each prior θ is characterized by a best guess, μ , an absolute uncertainty,
250 σ , and a relative uncertainty, $u = \sigma/\mu$, which can take values in the range $0 \leq u \leq 1$. Let $x_{\min} =$
251 $\min\{x', x''\}$ and $x_{\max} = \max\{x', x''\}$, where x can be μ , σ or u .

252 We now propose a series of properties that define the data reconciliation method.

253 **Property 1** (Lower and upper bounds). *The parameters of the final prior lie within the range set by the*
254 *parameters of the initial priors, $\mu_{\min} \leq \mu \leq \mu_{\max}$, $\sigma_{\min} \leq \sigma \leq \sigma_{\max}$ and $u_{\min} \leq u \leq u_{\max}$.*

255 **Property 2** (Commutativity). *The order in which the initial priors are combined does not matter, $f(\theta', \theta'') =$*
 256 *$f(\theta'', \theta')$.*

257 **Property 3** (Associativity). *Several initial priors can be combined and the resulting final prior is invariant to*
 258 *the order of reconciliation, $f(\theta', f(\theta'', \theta''')) = f(f(\theta', \theta''), \theta''')$.*

259 **Property 4** (Identity). *If the initial prior best guesses are identical, $\mu' = \mu''$ then the final prior best guess is*
 260 *identical, $\mu = \mu' = \mu''$. If the initial prior uncertainties are identical, $\sigma' = \sigma''$ then the final prior uncertainty*
 261 *is identical, $\sigma = \sigma' = \sigma''$.*

262 **Property 5** (Monotonicity). *The relative adjustment from initial to final prior increases with the relative*
 263 *magnitude of initial uncertainty:*

$$\frac{\mu - \mu'}{\mu'' - \mu} = g\left(\frac{\sigma'}{\sigma''}\right); \quad (8)$$

$$\frac{\sigma - \sigma'}{\sigma'' - \sigma} = h\left(\frac{\sigma'}{\sigma''}\right). \quad (9)$$

264 where $dg(x)/dx > 0$ and $dh(x)/dx > 0$.

265 **Property 6** (Absorption). *If initial prior θ' is known with minimal uncertainty, $u' = 0$, and θ'' is not,*
 266 *$u'' > 0$, then the final prior is identical to the first initial prior, $f(\theta', \theta'') = \theta'$. If initial prior θ' is known with*
 267 *maximal uncertainty, $u' = 1$, and θ'' is not, $u'' < 1$, then the final prior is identical to the second initial prior,*
 268 *$f(\theta', \theta'') = \theta''$.*

269 We believe that these six properties are uncontroversial and self-explanatory. **However, it turns**
 270 **out that the problem as formulated here has no solution, i.e., no formula can satisfy all of the above**
 271 **properties. We later overcome this hurdle by generalizing the problem formulation, to include two**
 272 **additional concepts: a hierarchy of data quality and the number of combined priors.**

273 3.2. The canonical data reconciliation method

274 The properties outlined in Section 3.1 constrain the range of data reconciliation algorithms but
 275 do not define a unique solution. However, Equations 8-9 suggests how it may be possible to obtain
 276 a solution. Let us consider that $g(x)$ and $h(x)$ take the simple yet flexible form of $g(x) = ax^b$ and
 277 $h(x) = cx^d$.

278 The condition of identity (Property 4), in the case of $\mu' = \mu''$ and $\sigma' = \sigma''$ leads to the
 279 indeterminacy:

$$\frac{\mu - \mu'}{\mu'' - \mu} = \frac{0}{0}.$$

280 But if the limit is approached as $\mu' = \mu - \delta$ and $\mu'' = \mu + \delta$, when $\delta \rightarrow 0$, then:

$$\frac{\mu - \mu'}{\mu'' - \mu} = \frac{\delta}{\delta} = 1.$$

281 Thus, under the condition of identity, Equations 8-9 imply that:

$$1 = a1^b;$$

$$1 = c1^d,$$

282 so $a = c = 1$. Let us further consider the simplest possible case $b = d = 1$, so that $g(\cdot)$ and $h(\cdot)$ are the
283 identity lines. Applying $g(x) = x$ and $h(x) = x$ to Equations 8-9 leads to:

$$\frac{\mu - \mu'}{\sigma'} = \frac{\mu'' - \mu}{\sigma''};$$

$$\frac{\sigma - \sigma'}{\sigma'} = \frac{\sigma'' - \sigma}{\sigma''}.$$

284 Rearranging terms:

$$\mu \left(\frac{1}{\sigma'} + \frac{1}{\sigma''} \right) = \frac{\mu'}{\sigma'} + \frac{\mu''}{\sigma''};$$

$$\sigma \left(\frac{1}{\sigma'} + \frac{1}{\sigma''} \right) = \frac{\sigma'}{\sigma'} + \frac{\sigma''}{\sigma''}.$$

285 Recalling that $u = \sigma/\mu$ we obtain the canonical data reconciliation method as:

$$\mu = \left(\frac{1/\sigma'}{1/\sigma' + 1/\sigma''} \right) \mu' + \left(\frac{1/\sigma''}{1/\sigma' + 1/\sigma''} \right) \mu''; \quad (10)$$

$$\frac{1}{\sigma} = \frac{1}{2} \left(\frac{1}{\sigma'} + \frac{1}{\sigma''} \right). \quad (11)$$

286 Equation 10 can be expressed in two other ways:

$$\mu = \frac{\sigma}{2} \left(\frac{\mu'}{\sigma'} + \frac{\mu''}{\sigma''} \right); \quad (12)$$

$$\frac{1}{u} = \frac{1}{2} \left(\frac{1}{u'} + \frac{1}{u''} \right). \quad (13)$$

287 Thus, if the ratio of relative adjustment of best guesses and uncertainties is identical to the ratio
288 of absolute uncertainties of the initial priors, the best-guess data reconciliation method is a weighted
289 average, where the weights are proportional to the inverse of absolute uncertainty, and the absolute
290 and relative uncertainty data reconciliation methods are harmonic averages.

291 Does this data reconciliation method satisfy the properties of Section 3.1? It is trivial to check that
292 Properties 1, 2, 4 and 5 are satisfied. But this is not the case for Properties 3 and 6. In the following
293 subsections we present suitable extensions of the canonical data reconciliation method to address these
294 problems.

295 3.3. The number of combined priors

296 The canonical data reconciliation method is not associative. The properties of $f(\theta', f(\theta'', \theta'''))$ are:

$$\frac{1}{u} = \frac{1}{2} \left(\frac{1}{u'} + \frac{1}{2} \left(\frac{1}{u''} + \frac{1}{u'''} \right) \right) = \frac{1}{2u'} + \frac{1}{4u''} + \frac{1}{4u'''} \\ \frac{1}{\sigma} = \frac{1}{2} \left(\frac{1}{\sigma'} + \frac{1}{2} \left(\frac{1}{\sigma''} + \frac{1}{\sigma'''} \right) \right) = \frac{1}{2\sigma'} + \frac{1}{4\sigma''} + \frac{1}{4\sigma'''}.$$

297 While the properties of $f((\theta', \theta''), \theta''')$ are:

$$\frac{1}{u} = \frac{1}{2} \left(\frac{1}{2} \left(\frac{1}{u'} + \frac{1}{u''} \right) + \frac{1}{u'''} \right) = \frac{1}{4u'} + \frac{1}{4u''} + \frac{1}{2u'''} \\ \frac{1}{\sigma} = \frac{1}{2} \left(\frac{1}{2} \left(\frac{1}{\sigma'} + \frac{1}{\sigma''} \right) + \frac{1}{\sigma'''} \right) = \frac{1}{4\sigma'} + \frac{1}{4\sigma''} + \frac{1}{2\sigma'''}.$$

298 Thus, $f(\theta', f(\theta'', \theta''')) \neq f((\theta', \theta''), \theta''')$. But upon some reflection, this result is in fact reasonable.
299 The final prior is the combination of two initial priors with equal weights. If some of these initial priors
300 is itself a combination of other initial priors, this information has to be considered explicitly.

301 Let us introduce a new quantity, n , as the *number of combined priors*, so that now a prior θ is defined
302 by a best guess, μ , an absolute uncertainty, σ , and n . Consider the following data reconciliation rule:

$$\mu = \left(\frac{n'/\sigma'}{n'/\sigma' + n''/\sigma''} \right) \mu' + \left(\frac{n''/\sigma''}{n'/\sigma' + n''/\sigma''} \right) \mu''; \quad (14)$$

$$\frac{1}{\sigma} = \frac{1}{n} \left(\frac{n'}{\sigma'} + \frac{n''}{\sigma''} \right); \quad (15)$$

$$n = n' + n''. \quad (16)$$

303 As before, Equation 14 can be expressed in two other ways:

$$\mu = \frac{\sigma}{n} \left(\frac{n'}{\sigma'} \mu' + \frac{n''}{\sigma''} \mu'' \right); \quad (17)$$

$$\frac{1}{u} = \frac{1}{n} \left(\frac{n'}{u'} + \frac{n''}{u''} \right). \quad (18)$$

304 This data reconciliation rule satisfies the first five properties of Section 3.1.

305 3.4. Ranking of data quality

306 The canonical data reconciliation method satisfies the absorption property of minimal uncertainty.
307 If $\sigma' = 0$ and $\sigma'' > 0$, then:

$$\frac{1}{\sigma'} + \frac{1}{\sigma''} \simeq \frac{1}{\sigma'},$$

308 and Equations 10-11 become:

$$\mu \simeq \left(\frac{1/\sigma'}{1/\sigma'} \right) \mu' + \left(\frac{1/\sigma''}{1/\sigma'} \right) \mu'' = (1)\mu' + (0)\mu''; \\ \sigma \simeq 2\sigma' = 0,$$

309 so $\mu = \mu'$ and $\sigma = \sigma'$. However, it does not satisfy the absorption property of maximal uncertainty. If
 310 $\sigma' = \mu'$ and $\sigma'' < \mu''$, then $u' = 0$ and Equations 10, 11 and 13 become:

$$\begin{aligned}\mu &= (\mu'' + \sigma'') \frac{\mu'}{\mu' + \sigma''}; \\ \sigma &= 2\sigma'' \frac{\mu'}{\mu' + \sigma''}; \\ u &= 2 \frac{u''}{1 + u''},\end{aligned}$$

311 and thus $\mu \neq \mu''$ and $\sigma \neq \sigma''$.

312 In order to ensure that the absorption of maximal uncertainty is satisfied, we use the concept of
 313 *data quality*, introduced in Rodrigues [25]. The idea is that, besides an uncertainty estimate, which
 314 formalizes quantitatively a degree of confidence in the accuracy of the best guess of a datum, it is also
 315 possible to formalize qualitatively a degree of confidence in the accuracy of a datum relative to others.

316 For the purpose of data balancing, Rodrigues [25] suggests that a datum that is considered to be
 317 of higher quality should be kept fixed while lower quality data are adjusted. The natural corollary, in
 318 the problem of data reconciliation, is to consider that when one wishes to combine two initial priors of
 319 differing levels of data quality, the prior of lower quality should be disregarded.

320 If a datum has unitary relative uncertainty, then it is maximally uninformative, and it is reasonable
 321 to disregard it. After all, a maximally uninformative prior should only be used if no better alternative
 322 is available. We therefore suggest that, if $\sigma' = \mu'$ and $\sigma'' < \mu''$, then $\theta = \theta''$ directly, without using
 323 Equations 10-11.

324 3.5. Summary

325 We now present the expressions for the combination of n initial priors, θ_i , with $i = 1, \dots, n$ into a
 326 single final prior θ . Addressing this problem requires the specification, for each prior, θ_i , of its best
 327 guess, μ_i , its absolute uncertainty, σ_i , and the number of previously combined priors, n_i .

328 If all relative uncertainties, $u_i = \sigma_i / \mu_i$, are in the range $0 < u_i < 1$, then the final prior properties
 329 are defined as:

$$\mu = \sum_{i=1}^n \left(\frac{n_i / \sigma_i}{n / \sigma} \right) \mu_i; \quad (19)$$

$$\frac{1}{\sigma} = \frac{1}{n} \left(\sum_{i=1}^n \frac{n_i}{\sigma_i} \right); \quad (20)$$

$$n = \sum_{i=1}^n n_i. \quad (21)$$

330 Equation 19 can be expressed as:

$$\frac{1}{u} = \frac{1}{n} \sum_{i=1}^n \frac{n_i}{u_i}. \quad (22)$$

331 If some initial priors have zero relative uncertainty, $u_i = 0$, then all other initial priors should
 332 be disregarded. If some initial priors have unitary relative uncertainty, $u_i = 1$, then it is they which
 333 should be disregarded.

334 In Figure 2 we illustrate the behaviour of Equation 19, when $n = 2$ and $n_1 = n_2 = \mu_2 = 1$. The
 335 plot shows different curves of the combined posterior best guess μ as a function of the prior best guess

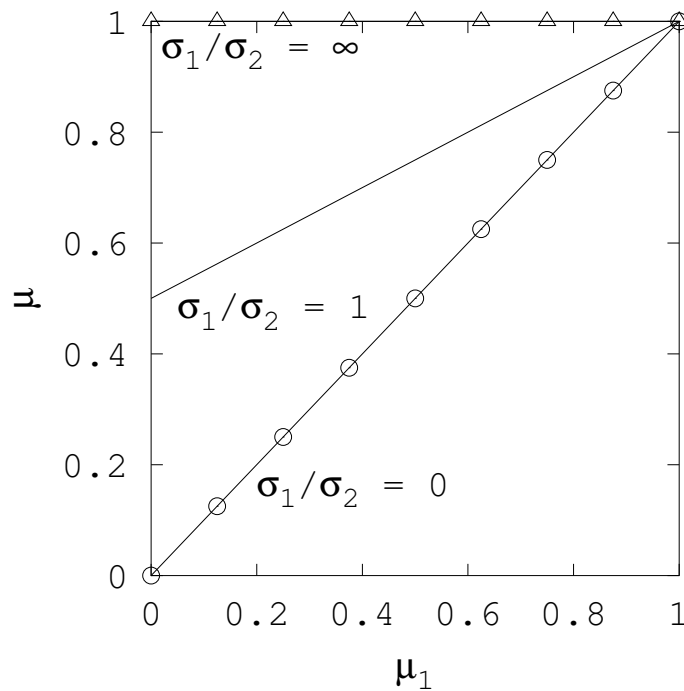


Figure 2. Best guess of the combined prior as a function of initial prior best guess, μ_1 , when $\mu_2 = 1$, for different values of σ_1/σ_2 : Solid line (—) when $\sigma_1/\sigma_2 = 1$, solid line with circle markers (—○—) when $\sigma_1/\sigma_2 = 0$ and solid line with triangle markers (—△—) when $\sigma_1/\sigma_2 = \infty$.

336 μ_1 , where each curve corresponds to a different ratio of uncertainties, σ_1/σ_2 . When both uncertainties
 337 are identical, the posterior best guess is the arithmetic average of the two prior best guesses. When the
 338 uncertainties differ, the best guess prior with the largest uncertainty contributes the least to combined
 339 best guess: in the limit case in which $\sigma_1 \gg \sigma_2$ the prior μ_1 is ignored and the posterior is similar to μ_2 ;
 340 when $\sigma_1 \ll \sigma_2$ the reverse occurs and $\mu \simeq \mu_1$.

341 In turn, Figure 3 we illustrate the behaviour of Equation 20. We still consider that $n = 2$ and
 342 $n_1 = n_2$ but now no explicit assumption about best guesses is necessary. Instead, the uncertainty of the
 343 second variable is fixed, $\sigma_2 = 1$, and the curve shows the value of the combined posterior as a function
 344 of prior uncertainty σ_1 . The figure describes an arc slightly above the diagonal line. When both prior
 345 uncertainties are identical ($\sigma_1 = 1$), then the posterior equals the priors, as expected. As σ_1 becomes
 346 smaller than σ_2 the combined prior becomes closer to σ_1 than to σ_2 , but always larger, $\sigma > \sigma_1$, except in
 347 the limit case $\sigma_1 \rightarrow 0$, in which case $\sigma \rightarrow \sigma_1$.

348 4. Conclusions and discussion

349 Herein we investigated using two distinct pathways the problem of reconciling multiple
 350 conflicting estimates in the course of database development. We assume that the developer (data
 351 snooper) is tooled with a best guess and uncertainty for each of those conflicting estimates.

352 First, we apply a maximum-entropy Bayesian inference method, under the limiting condition
 353 that the adjustment from prior to posterior uncertainties is small. Second, we obtain a canonical
 354 data reconciliation method through an axiomatic approach that is as simple as possible but satisfied
 355 important qualitative properties. Each approach verifies the other.

356 The resulting formula for the best guess, Equation 19, is a weighted average showing that, as
 357 the count of conflicting priors underlying a particular prior rises, the value of that prior increases in
 358 importance in terms of obtaining a solution. We get a similar result with the inverse of the uncertainty,

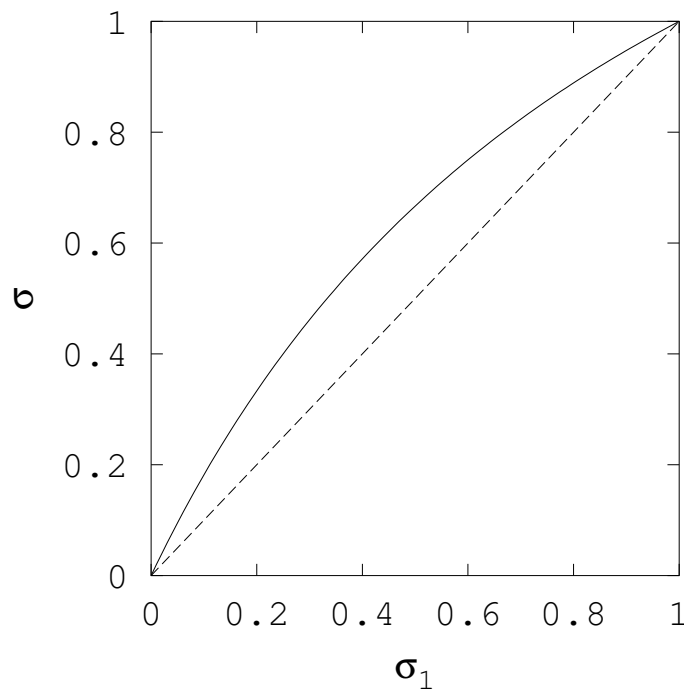


Figure 3. Solid line (–) and relative uncertainty, uncertainty of the final prior, σ , as a function of initial prior uncertainty, σ_1 , when $\sigma_2 = 1$. Dashed line (--) is the identity line.

359 that is, the narrower the uncertainty of an estimate the more it contributes to the final solution. The
 360 resulting formula for the uncertainty, Equation 20, is a harmonic average where the same factors are
 361 present: as the count of conflicting priors underlying a particular prior rises, the value of that prior
 362 increases in importance; and the narrower the uncertainty of a prior, the more it contributes to the final
 363 solution.

364 Of course, limitations to our approach must be mentioned. And the key one is certainly that, in
 365 some practical applications, the data snooper will lack information on either or both best guess and
 366 uncertainty. It may be that instead, one only has upper and lower bounds for the datum of interest
 367 to inform its best guess and uncertainty. This is certainly the case in some instances when data are
 368 censored, e.g., the anti-suppression problem of Gerking *et al.* [13] and Isserman and Westervelt [12].
 369 Future work using variable ranges with externally informed priors would be a natural extension of
 370 what is presented here. Indeed, some initial forays into this line of investigation are already underway,
 371 see, e.g., Makarkina and Lahr [35].

372 It should be mentioned that although the focus of attention here was on conflicting estimates
 373 arising from economic accounts there are other circumstances in which a formally identical problem
 374 arises, for example in expert elicitation [36].

375 **Author Contributions:** J.R. performed the conceptualization and formal analysis, M.L. provided the motivation
 376 and both J.R. and M.L wrote the manuscript.

377 **Funding:** Funding to be specified when accepted.

378 **Acknowledgments:** Any errors the paper may contain are the sole responsibility of the authors.

379 **Conflicts of Interest:** The authors declare no conflict of interest.

380

- 381 1. Lohr, S. The Origins of 'Big Data': An Etymological Detective
382 Story. *New York Times* **2013**. Available online in January 2018 at
383 <https://bits.blogs.nytimes.com/2013/02/01/the-origins-of-big-data-an-etymological-detective-story/>.
- 384 2. Duncan, G.T.; Elliot, M.; Salazar-Gonzalez, J.J. *Statistical Confidentiality: Principles and Practice*; Springer:
385 New York, 2011.
- 386 3. Garfinkel, R.; Gopal, R.; Goes, P. Privacy Protection of Binary Confidential Data against Deterministic,
387 Stochastic, and Insider Attack. *Management Science* **2002**, *48*, 749–764.
- 388 4. Evans, T.; Zayatz, L.; Slanta, J. Using Noise for Disclosure Limitation of Establishment Tabular Data.
389 *Journal of Official Statistics* **1998**, *14*, 537–551.
- 390 5. Miller, R.E.; Blair, P.D. *Input-Output Analysis: Foundations and Extensions*; Cambridge University Press:
391 Cambridge, UK, 2009. Second Edition.
- 392 6. Kruihof, R. Telefoonverkeersrekening. *De Ingenieur* **1937**, *52*, E15–E25.
- 393 7. Stone, R.; Meade, J.E.; Champernowne, D.G. The Precision of National Income Estimates. *The Review of*
394 *Economic Studies* **1942**, *9*, 111–125.
- 395 8. Byron, R. The Estimation of Large Social Account Matrices. *Journal of the Royal Statistical Society, Series A*
396 **1978**, *141*, 359–367.
- 397 9. Van der Ploeg, F. Reliability and the Adjustment of Sequences of Large Economic Accounting Matrices.
398 *Journal of the Royal Statistical Society, Series A* **1982**, *145*, 169–194.
- 399 10. Lahr, M.L.; Mesnard, L.d. Biproportional Techniques in Input-Output Analysis: Table Updating and
400 Structural Analysis. *Economic Systems Research* **2004**, *16* (2), 115–134.
- 401 11. Chen, B. A Balanced System of U.S. Industry Accounts and Distribution of the Aggregate Statistical
402 Discrepancy by Industry. *Journal of Business and Economic Statistics* **2012**, *30*, 202–211.
- 403 12. Isserman, A.; Westervelt, J. 1.5 Million Missing Numbers: Overcoming Employment Suppression in
404 County Business Patterns Data. *International Regional Science Review* **2006**, *29*, 311–335.
- 405 13. Gerking, S.; Isserman, A.; Hamilton, W.; Pickton, T.; Smirnov, O.; Sorenson, D. Anti-suppressants and
406 the Creation and Use of Non-Survey Regional Input-Output Models. *Regional Science Perspectives in*
407 *Economic Analysis: A Festschrift in Memory of Benjamin H. Stevens*; Lahr, M.; Miller, R., Eds.; Elsevier:
408 New York, 2001; pp. 379–406.
- 409 14. Rodrigues, J.; Marques, A.; Wood, R.; Tukker, A. A network approach for assembling and linking
410 input-output models. *Economic Systems Research* **2016**, *28*, 518–538.
- 411 15. Bourque, P.J.; Chambers, E.J.; Chiu, J.S.Y.; Denman, F.L.; Dowdle, B.; Gordon, G.; Thomas, M.; Tiebout, C.;
412 Weeks, E.E. *The Washington Economy: An Input-Output Study*; University of Washington: Seattle, 1967.
- 413 16. Miernyk, W.H.; Shellhammer, K.L.; Brown, D.M.; Coccari, R.L.; Gallagher, C.J.; Wineman, W.H. *Simulating*
414 *Regional Economic Development: An Interindustry Analysis of the West Virginia Economy*; D.C. Heath and Co.:
415 Lexington, 1970.
- 416 17. Jensen, R.C.; McGaurr, D. Reconciliation of purchases and sales estimates in an Input-Output table. *Urban*
417 *Studies* **1976**, *13*, 59–65.
- 418 18. Gerking, S. Reconciling 'rows only' and 'columns only' coefficients in an Input-Output model. *International*
419 *Regional Science Review* **1976**, *1*, 623–626.
- 420 19. Gerking, S. Reconciling reconciliation procedures in regional Input-Output Analysis. *International Regional*
421 *Science Review* **1979**, *4*, 23–36.
- 422 20. Weale, M. The reconciliation of values, volumes and prices in national accounts. *Journal of the Royal*
423 *Statistical Society. Series A* **1988**, *151* (1), 211–221.
- 424 21. Boomsma, P.; Oosterhaven, J. A double-entry method for the construction of bi-regional Input-Output
425 tables. *Journal of Regional Science* **1992**, *32*, 269–284.
- 426 22. Rassier, D.; Howells, T.; Morgan, E.; Empey, N.; Roesch, C. Implementing a reconciliation and balancing
427 model in the U.S. industry accounts. Working Paper WP2007-4, Bureau of Economic Analysis, 2007.
- 428 23. Jaynes, E.T. *Probability Theory: The Logic of Science*; Cambridge University Press: Cambridge, UK, 2003.
- 429 24. Kullback, S.; Leibler, R.A. On Information and Sufficiency. *Annals of Mathematical Statistics* **1951**, *22*
430 (1), 79–86.

- 431 25. Rodrigues, J.F.D. A Bayesian Approach to the Balancing of Statistical Economic Data. *Entropy* **2014**, *16*
432 (3), 1243–1271.
- 433 26. Persons, W.M. Fisher's Formula for Index Numbers. *The Review of Economics and Statistics* **1921**, *3*, 103–113.
- 434 27. Kop Jansen, P.; ten Raa, T. The Choice of Model in the Construction of Input-Output Coefficients Matrices.
435 *International Economic Review* **1990**, *31*, 213–227.
- 436 28. Rodrigues, J.; Domingos, T.; Giljum, S.; Schneider, F. Designing an indicator of environmental responsibility.
437 *Ecological Economics* **2006**, *59*, 256–266.
- 438 29. Laplace, P.S. *Essai Philosophique sur les Probabilités*; Courcier Imprimeur: Paris, France, 1814.
- 439 30. Jeffreys, H. *Theory of Probability*; Clarendon Press: Oxford, UK, 1939.
- 440 31. Jaynes, E.T. *Papers on Probability, Statistics and Statistical Physics*; Kluwer Academic Publishers: Dordrecht,
441 Netherlands, 1983.
- 442 32. Weise, K.; Woger, W. A Bayesian theory of measurement uncertainty. *Measurement Science and Technology*
443 **1992**, *4* (1), 1–11.
- 444 33. Rodrigues, J.F.D. Maximum-Entropy Prior Uncertainty and Correlation of Statistical Economic Data.
445 *Journal of Business & Economic Statistics* **2016**, *34*, 357–367.
- 446 34. Jaynes, E.T. Information Theory and Statistical Mechanics I. *Physical Review* **1957**, *106*, 620–630.
- 447 35. Makarkina, G.V.; Lahr, M.L. Estimating Nationwide Impacts using an Input-Output Model
448 with Fuzzy Parameters. 25th International Input-Output conference. Atlantic City, NJ, 2017.
449 <https://rucore.libraries.rutgers.edu/rutgers-lib/54916/>.
- 450 36. A., E.; Mikhailov, L.; Keane, J. Inconsistency reduction in decision making via multi-objective optimisation.
451 *European Journal of Operational Research* **2018**, *267*, 212–226.

452 © 2018 by the authors. Submitted to *Entropy* for possible open access publication under the terms and conditions
453 of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).