

NIASRA

NATIONAL INSTITUTE FOR APPLIED
STATISTICS RESEARCH AUSTRALIA



***National Institute for Applied Statistics Research
Australia***

University of Wollongong, Australia

Working Paper

01-21

A Few Statistical Principles for Data Science

Noel Cressie

*Copyright © 2021 by the National Institute for Applied Statistics Research Australia, UOW.
Work in progress, no part of this paper may be reproduced without permission from the Institute.*

National Institute for Applied Statistics Research Australia, University of Wollongong,
Wollongong NSW 2522, Australia Phone +61 2 4221 5076, Fax +61 2 4221 4998.

Email: karink@uow.edu.au

A few statistical principles for data science

Noel Cressie^{1*}

University of Wollongong

Summary

In any other circumstance, it might make sense to define the extent of the terrain (Data Science) first, and then locate and describe the landmarks (Principles). But this data revolution we are experiencing defies a cadastral survey. Areas are continually being annexed into Data Science. For example, biometrics was traditionally statistics for agriculture in all its forms but now, in Data Science, it means the study of characteristics that can be used to identify an individual. Examples of non-intrusive measurements include height, weight, fingerprints, retina scan, voice, photograph/video (facial landmarks and facial expressions), and gait. A multivariate analysis of such data would be a complex project for a statistician, but a software engineer might appear to have no trouble with it at all. In any applied-statistics project, the statistician worries about uncertainty and quantifies it by modelling data as realisations generated from a probability space. Another approach to uncertainty quantification is to find similar data sets, and then use the variability of results between these data sets to capture the uncertainty. Both approaches allow ‘error bars’ to be put on estimates obtained from the original data set, although the interpretations are different. A third approach, that concentrates on giving a single answer and gives up on uncertainty quantification, could be considered as Data Engineering, although it has staked a claim in the Data Science terrain. This article presents a few (actually nine) statistical principles for data scientists that have helped me, and continue to help me, when I work on complex interdisciplinary projects.

Key words: Applied statistics; hierarchical statistical models; measurement error; regression; spatial statistics

1. Does your analysis have ‘error bars’?

Science and Engineering have long been held separate by our universities and learned academies. Broadly speaking, Science is concerned with ‘why,’ and Engineering is concerned

* Author to whom correspondence should be addressed.

† National Institute for Applied Statistics Research Australia (NIASRA)
University of Wollongong
Wollongong NSW 2522, Australia
Email: ncressie@uow.edu.au

Acknowledgment. This research was supported by the Australian Research Council under Discovery Project DP190100180. The Sydney Business School, University of Wollongong, generously provided an environment where most of this article was written. The referees have generously given comments that have helped me clarify the exposition of my nine principles. My thanks go to Dr B. Maloney for his skillful technical assistance. I would like to express my appreciation to Dr J. Wong for insightful discussions on this and other material.

9 with ‘how.’ Data Science seems to include both (Cressie 2020), and sometimes it may be
10 hard to find the science in a Data Science solution. For example, an idea is born on how
11 a particular data set could be analysed, perhaps with only a vague idea of the question
12 being answered. Then algorithms and software are developed and, in some cases, the analysis
13 results in superior performance as judged by figures-of-merit applied to that data set, which
14 leads to a publication. Are the results extendible and do they have ‘error bars’?

15 Usually figures-of-merit are concerned with the *efficiency* of a method. However, *validity*
16 should be established first, and error bars allow this to be done. No matter how efficient a
17 prediction method is, if its error bars define a nominal 90% prediction interval but the true
18 quantity (i.e., the predictand) is only contained in these intervals 60% of the time, the method
19 lacks validity. This important figure-of-merit is called *coverage* and is analogous to one in
20 the medical literature called specificity.

21 This article is aimed at data scientists who want their error bars to be based on
22 probabilities and their coverages to be (approximately) as stated. There are data scientists
23 who have error bars that are based on empirical distributions derived from ensembles of like
24 data sets or from partitioning a single data set into like subsets. While the principles I present
25 are for data scientists who quantify uncertainties using probabilities, they could act as a guide
26 for both types. If your analyses include uncertainty quantifications, then you may find some
27 of the principles in the following sections useful for hypothesis generation, for climbing the
28 knowledge pyramid from data to information to knowledge, and for scientific inference.

29 **2. My principal principles**

30 When I was planning this article, I hovered between the words ‘laws,’ ‘rules,’
31 ‘principles,’ and ‘guidelines’ to establish some landmarks in Data Science. Laws should not
32 be broken, exceptions prove the rule, and guidelines prevent accidents. But principles should
33 make the analysis go better without interdicting other approaches: Improvising on something
34 Groucho Marx once said, if you don’t like my principles, you may have others you like better
35 . . .

36 A data scientist might see variability between strata and within strata. A statistical
37 scientist uses probability distributions to capture some or all of this variability. For example,
38 probability densities would typically be used to model within-strata variabilities, and these
39 densities may be different across strata in order to capture the between-strata variability. In
40 many, many cases, those strata differences are modelled via a regression model, leaving
41 the errors around the regression line (i.e., the within-strata variability) to be described by
42 a probability distribution, often a Gaussian (or sometimes called normal) distribution.

43 In my opinion, the worst word in the Statistics lexicon, which happens to have four
 44 letters, is ‘mean.’ It is the first moment of a probability distribution, and it is also used
 45 to describe the average of a collection of numbers, leading to many misunderstandings in
 46 Data Science. To data scientists who capture variability through the empirical distributions of
 47 summary statistics from similar data sets, the average is their first moment, so *they* could
 48 legitimately call it a mean of their empirical distribution. To avoid confusion, it will be
 49 henceforth assumed in this article that variability is captured by *probability* distributions,
 50 and a mean *is* (is and only is) the first moment of that distribution.

51 In this section, five general statistical principles are presented. Different areas of study
 52 will have their own and, in Section 3, I present three principles that are important in Spatial
 53 Statistics. In an epilogical section (Section 4), I present one more, the principle of respecting
 54 units, which makes nine principles in total.

55 2.1. All models are wrong . . . but some are wrong-er than others

56 In my work on remote sensing of carbon dioxide, data \mathbf{z} from a single sounding is
 57 a more-than-2,000-dimensional vector of radiances observed at wavelengths in the near
 58 infra-red part of the spectrum. Cressie (2018) contains a literature review that features the
 59 relationship between \mathbf{z} and the state of the atmosphere \mathbf{y} (a 40-to-50 dimensional vector) via
 60 the statistical model,

$$\mathbf{z} = \mathbf{f}(\mathbf{y}) + \varepsilon, \quad (1)$$

61 where unobserved \mathbf{y} and ε are statistically independent; $\varepsilon \sim N(\mathbf{0}, \Sigma_\varepsilon)$, where $N(\boldsymbol{\mu}, \Sigma)$ refers
 62 to a multivariate Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix Σ ; \mathbf{f} is a
 63 known vector of nonlinear forward models obtained from atmospheric physics/chemistry; Σ_ε
 64 is a covariance matrix known from experiments on the remote sensing instrument, performed
 65 in the laboratory before the satellite is launched; $\mathbf{y} = \boldsymbol{\mu}_\alpha + \boldsymbol{\alpha}$, where $\boldsymbol{\alpha} \sim N(\mathbf{0}, \Sigma_\alpha)$; and $\boldsymbol{\mu}_\alpha$
 66 and Σ_α are the first two moments of the (hidden) atmospheric state \mathbf{y} . The goal is to predict
 67 \mathbf{y} from data \mathbf{z} . In my experience, it is the specification of Σ_α (the covariance matrix of the
 68 state of the atmosphere) where the statistical model can be seriously wrong (Cressie 2018).

69 **Principle 2.1:** *Establish a true model (TM), perhaps different from the scientist’s working*
 70 *model (WM). Critically, compute the TM-distributional properties of the WM estimators.*

71 At the very least, this principle gives the data scientist an idea of the sensitivity of the
 72 model to misspecification. In the remote sensing example, the WM-based statistical analysis
 73 provides a predicted atmospheric state, $\hat{\mathbf{y}}_{WM}$, called a *retrieval*. However, the atmosphere
 74 does not care what WM was chosen; it acts under the true model (TM), and the principal

75 figures-of-merit that are calculated are (using obvious notation):

$$\begin{aligned} \text{True bias: } E_{TM}(\hat{\mathbf{y}}_{WM} - \mathbf{y}), \\ \text{True uncertainty: } \text{var}_{TM}(\hat{\mathbf{y}}_{WM} - \mathbf{y}). \end{aligned} \tag{2}$$

76 While $\hat{\mathbf{y}}_{WM}$ may be optimal under WM, it is not under TM, and the properties given by (2)
77 will expose the potential weaknesses of a WM chosen for its simplicity or its computational
78 convenience.

79 It is tempting (but wrong) for scientists to calculate instead:

$$\begin{aligned} E_{WM}(\hat{\mathbf{y}}_{WM} - \mathbf{y}), \\ \text{var}_{WM}(\hat{\mathbf{y}}_{WM} - \mathbf{y}), \end{aligned} \tag{3}$$

80 which are often facile calculations compared to (2). This was the case for NASA's Orbiting
81 Carbon Observatory-2 (OCO-2) satellite. Over time, empirical distributions of retrieval-error
82 ensembles showed that bias and variance should have been calculated from (2), not (3). This is
83 discussed by Cressie (2018), and an illustration of applying Principal 2.1 is given in Nguyen,
84 Cressie & Hobbs (2019).

85 In the rest of this article, I shall assume that the WM and TM are one and the same.
86 However, this principle is potentially applicable in all the examples below.

87 2.2. What you see is not what you want to get

88 There is an equivalent way to write the model (1) that shows how statistical scientists
89 can naturally build complexity into their models:

$$\begin{aligned} \mathbf{z}|\mathbf{y} &\sim N(\mathbf{f}(\mathbf{y}), \Sigma_\varepsilon) \\ \mathbf{y} &\sim N(\boldsymbol{\mu}_\alpha, \Sigma_\alpha). \end{aligned} \tag{4}$$

90 This is a *hierarchical statistical model (HM)*, which is defined by layers of conditional
91 probabilities (here two layers).

92 The top layer is the *data model* and can be written in abbreviated notation as $[\mathbf{z}|\mathbf{y}]$. It
93 models what we see – the data (conditional on the process \mathbf{y}). The next layer is the *process*
94 *model* that can be written as $[\mathbf{y}]$. It models what we want to get – the latent process, or state,
95 hidden behind the data. In fact, in the remote sensing problem described in Section 2.1, the
96 distribution of \mathbf{y} is conditional on \mathbf{w} , the meteorology of the atmosphere. If need be, the
97 process model could be modified to now consist of $[\mathbf{y}|\mathbf{w}]$ and $[\mathbf{w}]$. Then the HM would be
98 made up of three levels, and the sequence of conditional probabilities would be $[\mathbf{z}|\mathbf{y}]$, $[\mathbf{y}|\mathbf{w}]$,

99 and $[\mathbf{w}]$. In the specific case of OCO-2 retrievals, the meteorological process, \mathbf{w} , is obtained
 100 from a numerical weather forecasting model and is considered known.

101 **Principle 2.2:** *Build statistical models conditionally, through a data model and a process*
 102 *model. Infer the unknown process from the predictive distribution.*

103 The predictive distribution is,

$$[\mathbf{y}|\mathbf{z}] = \frac{[\mathbf{z}|\mathbf{y}] \cdot [\mathbf{y}]}{[\mathbf{z}]}, \quad (5)$$

104 where for convenience any parameters θ in $[\mathbf{z}|\mathbf{y}]$ and $[\mathbf{y}]$ are dropped from the notation (but
 105 they are still there). The HM components are featured in the numerator, which is equal to
 106 the joint distribution, $[\mathbf{z}, \mathbf{y}]$; and the denominator, $[\mathbf{z}]$, is simply a ‘normaliser’ to ensure that
 107 $[\mathbf{y}|\mathbf{z}]$ is a density (or probability mass) function that integrates (or sums) to 1. Equation (5) is
 108 Bayes’ Theorem but, in this case, the unknowns are the state elements in \mathbf{y} , not the parameters
 109 θ . No prior distributions on θ have been assumed in the making of this HM! However, a prior
 110 could be assumed or, pragmatically, an estimate $\hat{\theta}$ could be obtained from \mathbf{z} .

111 Derivation of the predictive distribution in (5) can be problematic in many cases of
 112 practical interest, such as when $[\mathbf{y}]$ is a highly multivariate probability distribution of a
 113 complex scientific process. Computational methods, particularly Markov chain Monte Carlo
 114 (MCMC), are in constant development to allow realisations from $[\mathbf{y}|\mathbf{z}]$ to be simulated, from
 115 which distributional properties such as the predictive mean and the predictive quantiles can
 116 be used to predict the state \mathbf{y} and the error bars, respectively.

The simplest special case of (4) is when the data (here, z) and the process (here, y) are univariate, and $f(\cdot)$ is the identity function:

$$\begin{aligned} z|y &\sim \text{N}(y, \sigma_\varepsilon^2), \\ y &\sim \text{N}(\mu_\alpha, \sigma_\alpha^2). \end{aligned}$$

Then $[y|z]$ is Gaussian, characterised by its first two moments:

$$\begin{aligned} \text{E}(y|z) &= \{(1/\sigma_\alpha^2)\mu_\alpha + (1/\sigma_\varepsilon^2)z\}\{1/\sigma_\alpha^2 + 1/\sigma_\varepsilon^2\}^{-1}, \\ \text{var}(y|z) &= \{1/\sigma_\alpha^2 + 1/\sigma_\varepsilon^2\}^{-1}. \end{aligned}$$

117 The predictor, $\hat{y} = \text{E}(y|z)$, is a weighted combination of the data and the mean of the
 118 unknown state y , where the weights depend on the signal-to-noise ratio, $\sigma_\alpha^2/\sigma_\varepsilon^2$. These results
 119 can be generalised to the multivariate case with linear retrieval equations in (4), namely
 120 $\mathbf{f}(\mathbf{y}) = \mathbf{c} + \mathbf{K}\mathbf{y}$. In the case of OCO-2, \mathbf{K} has more rows than columns, and the first two

121 moments of the predictive distribution, which is Gaussian, are:

$$\begin{aligned} E(\mathbf{y}|\mathbf{z}) &= \mathbf{y}_\alpha + \mathbf{G}(\mathbf{z} - \mathbf{c} - \mathbf{K}\mathbf{y}_\alpha), \\ \text{var}(\mathbf{y}|\mathbf{z}) &= \{\boldsymbol{\Sigma}_\alpha^{-1} + \mathbf{K}^\top \boldsymbol{\Sigma}_\varepsilon^{-1} \mathbf{K}\}^{-1}, \end{aligned} \tag{6}$$

122 where $\mathbf{G} = \{\boldsymbol{\Sigma}_\alpha^{-1} + \mathbf{K}^\top \boldsymbol{\Sigma}_\varepsilon^{-1} \mathbf{K}\}^{-1} \mathbf{K}^\top \boldsymbol{\Sigma}_\varepsilon^{-1}$ is sometimes called the gain matrix.

From (6), it can be seen that the precision matrix (i.e., the inverse of the variance matrix), written $\text{prec}(\cdot)$, of each component of the predictive distribution satisfies,

$$\text{prec}(\mathbf{y}|\mathbf{z}) = \text{prec}(\mathbf{y}) + \mathbf{K}^\top \text{prec}(\mathbf{z}|\mathbf{y}) \mathbf{K},$$

123 a result that holds when the matrix \mathbf{K} is any size. This decomposition of precision
124 demonstrates that, when going from $[\mathbf{y}]$ to $[\mathbf{y}|\mathbf{z}]$, the precision increases (i.e., the variance
125 decreases). Moreover, the predictive mean is unbiased; that is, $E(E(\mathbf{y}|\mathbf{z})) = E(\mathbf{y}) = \boldsymbol{\mu}_\alpha$.

126 In geostatistics, the importance of (5) is deeply misunderstood, because Matheron
127 (1963) originally formulated kriging in terms of what he called a regionalised variable,
128 $\{z(\mathbf{s}) : \mathbf{s} \in D \subset \mathbb{R}^d\}$, where \mathbb{R}^d is d -dimensional Euclidean space and D is a subset of \mathbb{R}^d
129 with volume $|D| > 0$. In this formulation, the data model and the process model are collapsed
130 into the single probability distribution, $\{z(\mathbf{s}) : \mathbf{s} \in D\}$. Then the goal is to predict $z(\mathbf{s}_0)$
131 from data $\mathbf{z} = (z(\mathbf{s}_1), \dots, z(\mathbf{s}_n))^\top$, for which the generic kriging predictor is $E(z(\mathbf{s}_0)|\mathbf{z})$,
132 the mean of $[z(\mathbf{s}_0)|\mathbf{z}]$.

133 However, when measurement of the process is taken into account, there is a HM that
134 differentiates between the observations \mathbf{z} and the underlying latent process $\{y(\mathbf{s}) : \mathbf{s} \in D\}$
135 and that is observed imperfectly through \mathbf{z} . The spatial trend and the spatial covariance
136 function are defined in the process model, although they are estimated from the noisy data \mathbf{z} .
137 Once the measurement error is included in the model, it is clear that geostatistics should do
138 kriging using $E(y(\mathbf{s}_0)|\mathbf{z})$ and not $E(z(\mathbf{s}_0)|\mathbf{z})$; see Cressie & Wikle (2011, Section 4.1).

139 Consequently, much of the earlier geostatistics software did not make inference on $y(\mathbf{s}_0)$
140 but chose to make inference on $z(\mathbf{s}_0)$ where the measurement error is confounded with the
141 process error. User beware: some still do, but those written for environmental applications
142 (e.g., `geoR`, `gstat`, `FRK`) give the correct kriging equations. The difference is most apparent
143 when the kriging variance is computed as $\text{var}(z(\mathbf{s}_0)|\mathbf{z})$ at location \mathbf{s}_0 , but then interpreted
144 incorrectly as the predictive variance $\text{var}(y(\mathbf{s}_0)|\mathbf{z})$ of the true process at \mathbf{s}_0 .

145 **2.3. Geophysicists conserve energy but what do data scientists conserve?**

146 Building physical models usually involves ensuring that mass or energy is conserved. If
 147 the system is leaking energy, then it needs to be plugged or the energy needs to be followed
 148 as it moves into another system.

149 Now consider a designed experiment where data \mathbf{z} are obtained and the statistical model
 150 fitted is

$$\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\xi}. \quad (7)$$

151 The model in (7) is a linear regression with covariate (or design) matrix \mathbf{X} , $\boldsymbol{\beta}$ is an unknown
 152 p -dimensional vector of regression coefficients, and $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top$ consists of random
 153 variables that are independent and identically distributed (iid) Gaussian random variables
 154 with mean 0 and variance σ_ξ^2 .

Suppose that the measuring instrument was carefully calibrated and, in the study
 protocol, its sample variance from repeated measurements was reported as a number that
 I shall write as $s_\varepsilon^2 > 0$. Note that the uncertainty in the measurements usually needs to
 be quantified outside the experiment in order to identify the linear-model-error variance.
 Scientific interest is primarily in $\boldsymbol{\beta}$, but σ_ξ^2 is by no means a nuisance parameter. Its restricted
 maximum likelihood (REML) estimate is:

$$s_\xi^2 = (\mathbf{z} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{z} - \mathbf{X}\hat{\boldsymbol{\beta}}) / (n - p),$$

155 where $\hat{\boldsymbol{\beta}}$ is the maximum likelihood (and also the ordinary least squares) estimate of $\boldsymbol{\beta}$.

156 It could almost be a rule that in any study of this sort, you will see $s_\xi^2 > s_\varepsilon^2$. Where did
 157 the component of variability, $(s_\xi^2 - s_\varepsilon^2)$, go?

158 **Principle 2.3:** *In any well defined statistical model, there is conservation of variability.*

Indeed, the model given by (7) is not defined well enough: The error term ξ_i is, in fact, made
 up of two components of variance:

$$\xi_i = \delta_i + \varepsilon_i; \quad i = 1, \dots, n,$$

159 where $\{\varepsilon_i : i = 1, \dots, n\}$ represent measurement errors. Often forgotten are $\{\delta_i : i =$
 160 $1, \dots, n\}$, which represent model errors resulting from using the linear-regression model,
 161 $\{\mathbf{x}_i^\top \boldsymbol{\beta} : i = 1, \dots, n\}$, and which are key in accounting for the inexactness of using any
 162 model (linear or nonlinear).

163 The scientific process $\{y_i : i = 1, \dots, n\}$ is given by $y_i = \mathbf{x}_i \boldsymbol{\beta} + \delta_i$, and an observation
 164 of it is $z_i = y_i + \varepsilon_i$, for $i = 1, \dots, n$. The HM captures the variability of \mathbf{z} and \mathbf{y} beautifully,

165 as follows:

$$\begin{aligned} \mathbf{z}|\mathbf{y} &\sim \mathbf{N}(\mathbf{y}, \sigma_\varepsilon^2 \mathbf{I}) \text{ [or, } \mathbf{z} = \mathbf{y} + \varepsilon], \\ \mathbf{y} &\sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \sigma_\delta^2 \mathbf{I}) \text{ [or, } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \delta]. \end{aligned} \quad (8)$$

166 Hence, from (8), we obtain the earlier result in vector form:

$$\mathbf{z} = (\mathbf{X}\boldsymbol{\beta} + \delta) + \varepsilon = \mathbf{X}\boldsymbol{\beta} + (\delta + \varepsilon) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\xi}, \quad (9)$$

167 where δ and ε are independent mean-zero random vectors.

168 Comparing (7) and (8), we see that the model given by (7) should be augmented
169 with the equation, $\boldsymbol{\xi} = \delta + \varepsilon$, which results in the conservation-of-variability equation,
170 $\sigma_\xi^2 = \sigma_\delta^2 + \sigma_\varepsilon^2$. Consequently, knowing s_ξ^2 (an estimate of σ_ξ^2 that is obtained outside the
171 experiment) and, having computed the REML estimate of σ_ε^2 , an estimate of the model error,
172 σ_δ^2 , can be obtained as,

$$s_\delta^2 = s_\xi^2 - s_\varepsilon^2 \quad (10)$$

173 (provided the right-hand side is non-negative).

174 In its simplest form, conservation of variability says that the *total variability* is equal to
175 the variability due to *model uncertainty plus* the variability due to *measurement uncertainty*.
176 When variability is captured using a probability space defined by a HM, this principle can be
177 expressed as:

$$\text{var}(\mathbf{z}) = \text{var}(\mathbf{E}(\mathbf{z}|\mathbf{y})) + \mathbf{E}(\text{var}(\mathbf{z}|\mathbf{y})). \quad (11)$$

For example, consider the HM defined by (8). Since $\text{var}(\mathbf{z}) = \text{var}(\boldsymbol{\xi})$ and $\text{var}(\mathbf{y}) = \text{var}(\delta)$, we have

$$\sigma_\xi^2 \mathbf{I} = \sigma_\delta^2 \mathbf{I} + \sigma_\varepsilon^2 \mathbf{I},$$

178 and variability is conserved.

179 This might seem obvious to you, or perhaps even trivial. Again consider (7), and suppose
180 you want to predict an unknown value, y_{n+1} , outside the data set, but only the p -dimensional
181 estimate $\hat{\boldsymbol{\beta}}$ and the covariate, \mathbf{x}_{n+1} , are at your disposal. Most regression textbooks would
182 say that you should use as predictor, $\mathbf{x}_{n+1}^\top \hat{\boldsymbol{\beta}}$. However, a scientist wants to predict the
183 value of the scientific process, $y_{n+1} = \mathbf{x}_{n+1}^\top \boldsymbol{\beta} + \delta_{n+1}$. Using the well defined HM (8), its
184 predictor is $\hat{y}_{n+1} = \mathbf{E}(y_{n+1}|\mathbf{z}) = \mathbf{x}_{n+1}^\top \hat{\boldsymbol{\beta}} + \hat{\delta}_{n+1}$, where $\hat{\delta}_{n+1} = \mathbf{E}(\delta_{n+1}|\mathbf{z})$. Since $\{\delta_i\}$ are
185 iid $\mathbf{N}(0, \sigma_\delta^2)$, $\hat{\delta}_{n+1} = \mathbf{E}(\delta_{n+1}|\mathbf{z}) = 0$; however, $\text{var}(\delta_{n+1}|\mathbf{z})$ is *not zero*, and that has to be
186 recognised in order to perform valid predictive inference, as given below (e.g., Cressie 2020).

The prediction error is $(\hat{y}_{n+1} - y_{n+1})$, and its first two moments are:

$$\begin{aligned} E(\hat{y}_{n+1} - y_{n+1}) &= 0, \\ E(\hat{y}_{n+1} - y_{n+1})^2 &= E(\mathbf{x}_{n+1}^\top (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}))^2 + E(\text{var}(\delta_{n+1}|\mathbf{z})), \end{aligned}$$

187 since the expectation of the cross-product is zero. Thus, if $\text{var}(\delta_{n+1}|\mathbf{z})$ were forgotten by
188 the data scientist, such as would happen if the presence of $\{\delta_i\}$ were not recognised in the
189 statistical model, a forecasting decision about future values of the process $\{y_i\}$ would be
190 overly optimistic and potentially harmful.

191 Principle 2.3 also shows up in the analysis of variance (ANOVA) method, where
192 the ‘between sum of squares’ corresponds to the variance of the conditional expectation
193 in (11), the ‘within sum of squares’ corresponds to the expectation of the conditional
194 variance in (11), and the ‘total sum of squares’ corresponds to the left-hand side of (11).
195 Conservation of variability implicitly includes variability and *covariability*. For example, if
196 two random variables ε_1 and ε_2 have correlation $\rho \neq 0$, then it should be recognised that
197 $\text{var}(\varepsilon_1 + \varepsilon_2) \neq \text{var}(\varepsilon_1) + \text{var}(\varepsilon_2)$. This has often been ignored by scientists doing an error
198 budget (e.g., in remote sensing retrievals; see Connor et al. 2016). Depending on the sign
199 of ρ , $\text{var}(\varepsilon_1) + \text{var}(\varepsilon_2)$, may be either larger than or smaller than the total variability, since
200 $\text{var}(\varepsilon_1 + \varepsilon_2) = \text{var}(\varepsilon_1) + \text{var}(\varepsilon_2) + 2\rho\text{var}(\varepsilon_1)^{1/2}\text{var}(\varepsilon_2)^{1/2}$.

201 The rules of probability theory can explain easily how the variability can seem to
202 disappear, and then they can show us where to find it. It is less natural to do so simply with
203 empirical distributions, because the pairs $\{(z_i, y_i)\}$ involve the unavailable, hidden variable
204 $\{y_i\}$. Further, if $\{z_i\}$ and $\{x_j\}$ are two sets of observations, there is no guarantee that they
205 will occur in pairs, and hence the empirical correlation, r , might be difficult to obtain.

The famous bias–variance trade-off is another manifestation of Principle 2.3. In the
context of estimation of a fixed but unknown parameter θ with an estimator $\hat{\theta}(\mathbf{z})$, where
recall \mathbf{z} is the data vector, the mean-squared error is the sum of the estimator’s squared bias
and its variance:

$$E(\hat{\theta}(\mathbf{z}) - \theta)^2 = (E(\hat{\theta}(\mathbf{z})) - \theta)^2 + \text{var}(\hat{\theta}(\mathbf{z})).$$

206 Generally speaking, an estimator that decreases the bias will increase the variance, and *vice*
207 *versa*. The mean-squared error might be decreased using different estimators, but there is a
208 trade-off to be made when trying to decrease both bias and variance at the same time.

209 **2.4. The holy grail: all scales of variation are additive**

210 A statistical analysis cannot get very much simpler than fitting a simple linear regression
 211 (a special case of (7)) to $\{(z_i, x_i) : i = 1, \dots, n\}$, as follows:

$$z_i = \beta_1 + \beta_2 x_i + \xi_i; i = 1, \dots, n, \quad (12)$$

212 where $\{\xi_i : i = 1, \dots, n\}$ are iid $N(0, \sigma_\xi^2)$. However, is (12) appropriate if $\{z_i : i = 1, \dots, n\}$
 213 are photon counts at wavelengths in a given band of the electro-magnetic spectrum, or if they
 214 are percentages of trace elements in soil?

215 Exploratory data analysis (EDA) might reveal a highly skewed histogram of $\{z_i\}$. After
 216 plotting the histogram of $\{\log z_i\}$, we might achieve a more symmetric appearance; assume
 217 for the moment that we do. This would prompt a plot of not only $\{z_i\}$ versus $\{x_i\}$, but
 218 also of $\{\log z_i\}$ versus $\{x_i\}$, and the latter might look more linear. Further EDA, in the
 219 form of residual plots, after fitting a simple linear regression to the data $\{z_i\}$ and one to the
 220 transformed $\{\log z_i\}$, would be carried out. And, for example, it might be that the residual
 221 variability of $\{z_i\}$ appears to increase with x (i.e., heteroskedasticity), but the residual
 222 variability of $\{\log z_i\}$ shows no dependence on x (i.e., homoskedasticity).

223 Admittedly, this is a story, not methodology, but in my experience with analysing data,
 224 it happens enough to propose the next principle (Cressie 1985).

225 **Principle 2.4:** *Seek a transformation of the scientific process where all components of*
 226 *variation act and interact additively.*

227 This principle looks more like a doctrine but, as I explain below, a fitting of (12) to the data
 228 $\{z_i\}$ is following this principle, where the transformation is simply the identity.

229 Suppose the quantity y has variabilities that can be expressed as ‘large-scale’ and ‘small-
 230 scale.’ Scientists usually put their theories in the large scale and their errors (both in the theory
 231 and in the measurements) in the small scale. Statistical scientists know that to make inference
 232 on the large scale of the scientific process, its small scale has to be modelled well, usually
 233 by a random error term, δ . However, the model I now give for the scientific process is more
 234 basic, given in terms of those large and small scales of variability: It is additive and of the
 235 following form,

$$y = (\mu^{(1)} + \dots + \mu^{(p)}) + (\delta^{(1)} + \dots + \delta^{(N)}), \quad (13)$$

236 where p and N are positive integers and, for convenience, the subscript ‘ i ’ has been dropped.
 237 For example, the large-scale variation in simple linear regression has $p = 2$, $\mu^{(1)} = \beta_1$, and
 238 $\mu^{(2)} = \beta_2 x$. For the small-scale variation in (13), N is large, and $\{\delta^{(l)} : l = 1, \dots, N\}$ are
 239 physically interpretable components with variabilities that are very small but such that their
 240 total variability is the variability of the difference, $y - \mu^{(1)} - \dots - \mu^{(p)}$; see Principle 2.3.

The model (13) is a physical model where the large scales act additively with each other, the small scales act additively with each other, and the large scales and the small scales interact additively. It can be made statistical by assuming that these many small-scale effects, $\delta^{(1)}, \dots, \delta^{(N)}$, are random, statistically independent, have mean zero and variances that are $O(1/N)$, with the same leading coefficient, denoted here as σ_δ^2 . Now let N be large and define the small-scale variations collectively as δ ; then (13) can be written as,

$$y = \mu^{(1)} + \dots + \mu^{(p)} + \delta,$$

241 where

- 242 1. $E(y) = \mu^{(1)} + \dots + \mu^{(p)}$ (additivity),
- 243 2. $\text{var}(y) = \sigma_\delta^2$ (homoskedasticity),
- 244 3. $\delta \sim N(0, \sigma_\delta^2)$ (central limit theorem).

245 The regression (7) is obtained when the large-scale effects $\{\mu^{(k)} : k = 1, \dots, p\}$ are
 246 identified with $\{\beta_k x_k : k = 1, \dots, p\}$. It is the default model used in much of science, but
 247 this discussion shows that it originates from the imposition of additivity within and between
 248 scales of variability. From Principle 2.2, the measurement of a phenomenon should be
 249 separated from how the phenomenon behaves in nature and, indeed, the generalised linear
 250 model does this through a link function (McCullagh & Nelder 1989). Hence, Principle 2.4
 251 covers many statistical models and offers a plausible explanation of why a linear relation,
 252 homoskedasticity, and Gaussianity are often found to occur together after a transformation of
 253 the data (e.g., Cressie 1978). Of course, it is easy to construct probability models where the
 254 principle does not hold, but one should keep in mind that probability theory is used to model
 255 nature's variability, not the other way around.

256 When nonlinearity is inherent to the physical system, such as would occur when there
 257 are barriers or thresholds, the quest looks to be futile. However, after following Principle 2.2,
 258 it may just be that the grail, Principle 2.4, is hidden deep in the process model (Berliner,
 259 Wikle & Cressie 2000).

260 2.5. Could swans be red?

261 Taleb (2007) published a book about 'Black Swan' events that, when they happen, are
 262 considered to have been unpredictable. Such events could have a major impact on Earth's
 263 environment or, as has happened in 2020, on global-population health and economies around
 264 the globe.

265 This black-swan meme has its roots in seventeenth-century Europe. Up to then, no
 266 swans had ever been observed that were black but, in 1697, Dutch explorers became the

267 first Europeans to observe black swans during a voyage that took them to Western Australia.
 268 The implication in Taleb's book is that if the scientific world in the 1600s could not predict
 269 black swans, how can scientists predict catastrophic environmental events in the twenty-first
 270 century?

271 My response is that 300+ years ago the best scientific minds in Europe were too
 272 certain about their science; that is, if they had been asked to put probabilities (according
 273 to abundance) on the colours that swans could be: red, orange, . . . , violet, white, black, they
 274 would have put $0, 0, \dots, 1, 0$, which is a well defined probability mass function. Certainly,
 275 black-swan events are *not* predictable if the scientific model is 100% certain that they do
 276 not exist. Because their probability model gave black swans and indeed coloured (including
 277 red) swans, zero probability, this unimagined event did not emerge out of the 'ether' in
 278 subsequent inferences, until one was observed. Do we have to wait until a highly unusual
 279 event occurs before we are forced to change the probability model? The real lesson from the
 280 black-swan meme is that scientific knowledge is never perfect, that modellers need to explore
 281 the parameter space thoroughly, and that they need to 'spend' some probability on highly
 282 unusual events.

283 At the time of writing this article, our species is under attack from a virus that was
 284 originally called 'novel coronavirus,' so new that it took several weeks before the infection it
 285 caused had a name: COVID-19. Virologists certainly had assigned a small positive probability
 286 that each new decade would have its own 'novel' pandemic. However, politicians (and most
 287 economists) appear to have put zero probability on a severe, worldwide economic disruption.

288 Given a swine-flu-like pandemic has occurred, the conditional probability of some
 289 economic disruption is *not zero*. But, given a severe pandemic has occurred, the conditional
 290 probability of severe economic disruption is substantial. This conditional probability is then
 291 multiplied by the probability of a severe pandemic, which is by no means zero given the
 292 ability of viruses to mutate and occasionally jump from animals to humans. The product
 293 of these two is the joint probability of a severe pandemic followed by severe economic
 294 disruption, which is not negligible. This has happened twice in the last hundred years, and
 295 it will likely happen more often with humans and animals sharing in an ever-more-crowded
 296 environment.

297 **Principle 2.5:** *When building probability models, look carefully where zero probabilities*
 298 *are assumed (perhaps implicitly) and, with the same care, move appropriate probabilities*
 299 *away from zero. Calculate joint probabilities from products of conditional (not marginal)*
 300 *probabilities, unless entropy is maximal.*

301 Unfortunately, uncertainty quantification through joint probabilities all too often comes
 302 from multiplying marginal probabilities as if each event in the causal chain were independent.
 303 It does not 'hurt' for small probabilities to be assigned to $\Pr(\text{bird is colour } c \mid \text{bird is a swan})$

304 where c also covers the colour ‘black.’ Then $\Pr(\text{bird is a swan and it is black}) = \Pr(\text{bird is}$
 305 $\text{black} \mid \text{bird is a swan}) \times \Pr(\text{bird is a swan})$. The point being made here is philosophical, and
 306 it could be applied to $\Pr(\text{pandemic followed by global economic disruption})$, for different
 307 severities of both events. The worst thing would be to make $\Pr(\text{global economic disruption}$
 308 $\mid \text{pandemic})$ equal to zero, since then there would be a lack of planning for the health and
 309 economic crises brought on by a pandemic (Mackenzie 2020).

310 If nothing were known about the conditional probability, a fall-back is the maximum-
 311 entropy model where the marginal probability, $\Pr(\text{bird is colour } c)$, is used in place of
 312 $\Pr(\text{bird is colour } c \mid \text{bird is a swan})$. The marginal probability is not zero, since it is based
 313 on abundance and, for example, there are birds that are predominantly coloured red (e.g.,
 314 Australia’s king parrot). The maximum-entropy principle is discussed in Cressie, Richardson
 315 & Jaussent (2004).

316 Experiences over the last 100 years mean that $\Pr(\text{global economic disruption})$, whether
 317 caused by a pandemic or not, should not be zero, yet governments described the events
 318 of March–April 2020 as unimaginable. Events of small (but not zero) probabilities with
 319 consequences expressed through a loss function, allow a non-zero expected loss to be
 320 calculated, which can then be used to make optimal decisions that mitigate the consequences
 321 (e.g., see Berger 1985 for a discussion of decision theory).

322 I have just presented five statistical principles that should be useful in Data Science.
 323 Data often come with location information, in which case the data scientist will likely be
 324 using spatial-analysis methods. In the next section, I present three principles that are specific
 325 to Spatial Statistics.

326 **3. A few spatial statistical principles**

327 Those of us who work in Spatial Statistics will know of Tobler’s First Law of Geography
 328 (Tobler 1970). In spatial statistical science, it really is a ‘principle’ rather than a ‘law’ and, in
 329 the following subsections, I shall present it and two other principles that I have found useful
 330 in this area of research.

331 **3.1. Patches in close proximity are commonly more alike . . .**

332 In his famous 1935 book on experimental design, Fisher (1935, p. 66), wrote: ‘After
 333 choosing the area we usually have no guidance beyond the widely verified fact that patches
 334 in close proximity are commonly more alike, as judged by the yield of crops, than those
 335 which are further apart.’ A spatial statistician sees this as the making of a principle but, at that
 336 time, Fisher made a sharp right turn. In his analyses of field trials he applied randomisation
 337 to eradicate the pest: spatial correlation! Cressie & Wikle (2011, Ch.1 and Ch. 4) give

338 some historical perspective to the work of researchers who took roads less travelled and
 339 developed this vibrant area we now call Spatial Statistics. Some of this development comes
 340 from Geography, and so it is fitting that the first and most important principle in this section
 341 has become known as the First Law of Geography. Originally articulated by Tobler (1970), it
 342 is given here in exactly his words.

343 **Principle 3.1** *Everything is related to everything else, but near things are more related than*
 344 *distant things (Tobler 1970).*

345 This principle is at the core of what we do in spatial and spatio-temporal statistics.
 346 For example, in remote sensing of Earth's surface, the scene of interest D contains *many*
 347 retrievals, $\{z(\mathbf{s}_i) : i = 1, \dots, n\}$, where $\{\mathbf{s}_i : i = 1, \dots, n\}$ are the (lon, lat) locations of
 348 the n data inside D (retrieved over a short time period). There is a hidden spatial process
 349 $\{y(\mathbf{s}) : \mathbf{s} \in D\}$ that the geophysicist would like to infer and, in a spatial HM, a spatial
 350 statistical model for $\{y(\mathbf{s}) : \mathbf{s} \in D\}$ is built around Principle 3.1.

For example, consider a simple process model, appropriate for a small scene D :

$$y(\mathbf{s}) = \beta_1 + \beta_2 \text{lat}(\mathbf{s}) + \delta(\mathbf{s}), \text{ for } \mathbf{s} \in D,$$

where the mean of the process is a linear function of latitude, and $\{\delta(\mathbf{s}) : \mathbf{s} \in D\}$ is a spatial error process with mean zero and stationary covariance function, $C_\delta(\mathbf{h}) = \text{cov}(\delta(\mathbf{s} + \mathbf{h}), \delta(\mathbf{s})) = \sigma_\delta^2 \cdot \exp(-\|\mathbf{h}\|/\phi)$. Notice that $C_y(\mathbf{h}) = \text{cov}(y(\mathbf{s} + \mathbf{h}), y(\mathbf{s})) = C_\delta(\mathbf{h})$; $C_y(\mathbf{0}) = C_\delta(\mathbf{0}) = \sigma_\delta^2$; and the scale parameter ϕ controls how 'more related' things are. With this parameterisation, $\phi = 0$ is the degenerate case of no spatial correlation and, as ϕ increases, the spatial correlation increases for a given distance $\|\mathbf{h}\|$. The data model in this example is also simple, that of additive independent measurement error. The data vector is $\mathbf{z} = (z(\mathbf{s}_1), \dots, z(\mathbf{s}_n))^\top$ and

$$z(\mathbf{s}_i) = y(\mathbf{s}_i) + \varepsilon_i; i = 1, \dots, n,$$

351 where $\{\varepsilon_i : i = 1, \dots, n\}$ are independent mean-zero errors.

Assuming for the moment that all the parameters $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma_\delta^2, \phi, \sigma_\varepsilon^2)^\top$ are known, inference on $\{y(\mathbf{s}) : \mathbf{s} \in D\}$ comes from summaries of the predictive distribution, $[\{y(\mathbf{s}) : \mathbf{s} \in D\} | \mathbf{z}]$. For a Gaussian process model $y(\cdot)$ and Gaussian measurement errors $\{\varepsilon_i\}$, the predictive distribution is also a Gaussian process whose first two moments, $E(y(\mathbf{s}) | \mathbf{z})$ and $\text{cov}(y(\mathbf{s}), y(\mathbf{u}) | \mathbf{z})$, for $\mathbf{s}, \mathbf{u} \in D$, can be obtained analytically. From Section 2.2, a common predictor of $y(\mathbf{s})$ is $\hat{y}(\mathbf{s}) = E(y(\mathbf{s}) | \mathbf{z})$, whose statistical properties are needed for inference. It is straightforward to see that $E(\hat{y}(\mathbf{s})) = E(y(\mathbf{s}))$ (i.e., the predictor is

unbiased) and the mean-squared predictor error is:

$$E(\hat{y}(\mathbf{s}) - y(\mathbf{s}))^2 = E(\text{var}(y(\mathbf{s})|\mathbf{z})) = \text{var}(y(\mathbf{s})|\mathbf{z}),$$

352 which is the predictive variance (the last equality being due to the Gaussian assumptions
353 made). Further summaries might come from well chosen percentiles (e.g., 2.5%, 25%, 50%,
354 75%, 97.5%) of $[y(\mathbf{s})|\mathbf{z}]$.

355 3.2. What is one person's mean function could be another person's spatial error

In spatial statistics, it might appear that the spatial-prediction problem is quite difficult, because the number of predictions to be made is often greater than the number of data. Even for the problem of parameter estimation, the number of 'degrees of freedom' in n spatially dependent data will not be n , as I now show. Applying Principle 3.1 with a stationary spatial covariance model $C_z(\mathbf{h}) = \text{cov}(z(\mathbf{s} + \mathbf{h}), z(\mathbf{s}))$, for $\mathbf{h} \in \mathbb{R}^d$, that exhibits only positive correlations, it is easy to see that for $\bar{z} = (1/n) \sum_{i=1}^n z(\mathbf{s}_i)$,

$$\text{var}(\bar{z}) = (1/n)\{\sigma_z^2 + 2 \sum_{i < j} C_z(\mathbf{s}_i - \mathbf{s}_j)/n\} > \frac{\sigma_z^2}{n},$$

where $\sigma_z^2 = C_z(\mathbf{0})$. Hence one can define the *effective degrees of freedom*, n_{eff} , as

$$n_{\text{eff}} = \sigma_z^2 / \text{var}(\bar{z}) < n;$$

356 that is, under the almost ubiquitous spatial model of positive spatial correlation, $n_{\text{eff}} < n$. In
357 the remote sensing context where there are many observations taken over a short period of
358 time, Zhang, Cressie & Wunch (2017) gave an example where $n = 2961$ and $n_{\text{eff}} = 202.3$,
359 less than one-tenth of the sample size!

360 These calculations are a warning that, in the spatial (and spatio-temporal) setting,
361 intuition learned from the 'iid errors' model has to be modified, sometimes substantially.
362 Critically, the probability model that captures the spatial variability has to be well specified,
363 and we now discuss how this can be done.

364 The classical spatial statistical model consists of a deterministic mean process, $\{\mu(\mathbf{s}) :$
365 $\mathbf{s} \in D\}$, and a random spatial error process, $\delta(\mathbf{s}) = y(\mathbf{s}) - \mu(\mathbf{s})$, so that

$$y(\mathbf{s}) = \mu(\mathbf{s}) + \delta(\mathbf{s}), \text{ for } \mathbf{s} \in D; \quad (14)$$

366 see Cressie (1993, Ch. 2). It is often a personal choice what components go into $\mu(\mathbf{s})$.
367 For linear regression, $\mu(\mathbf{s}) = \mathbf{x}(\mathbf{s})^\top \beta$, for $\mathbf{s} \in D$, but the set of possible covariates, $\mathbf{x}(\mathbf{s})$,

368 typically comes from the modeller’s subject-matter knowledge, augmented by *spatial trend*
 369 terms that are linear and higher-order functions of $\mathbf{s} = (s_1, \dots, s_d)^\top$.

370 What is the effect of not having included an important spatially varying covariate,
 371 $x_{p+1}(\mathbf{s})$ say, in $\mu(\mathbf{s})$? From Principle 2.3, $\{\delta(\mathbf{s})\}$ has to absorb the contribution to spatial
 372 variability that $\{x_{p+1}(\mathbf{s})\}$ would have made had it been included in the regression. This
 373 means that a fixed effect that has been inadvertently left out will be picked up in the spatial
 374 statistical model (14), as a random effect.

375 **Principle 3.2:** *Assume a decomposition of a spatial process into fixed effects plus random*
 376 *effects. While it is not unique, the decomposition must be chosen to conserve variability.*

This principle is a refinement of Principle 2.3 that is adapted here for Spatial Statistics. The variability of the deterministic-mean component is measured differently from that of a random-error component. Define the *regional variance* of $\mu(\cdot)$ as,

$$s_\mu^2 = \left(\frac{1}{|D|} \right) \int_D (\mu(\mathbf{s}) - \bar{\mu})^2 \mathbf{d}\mathbf{s},$$

377 where $\bar{\mu} = (1/|D|) \int_D \mu(\mathbf{s}) \mathbf{d}\mathbf{s}$ is the regional average of $\mu(\cdot)$ (Lahiri et al. 1999). Suppose
 378 that a stationary covariance function, $C_\delta(\cdot)$, describes the covariability in $\{\delta(\mathbf{s}) : \mathbf{s} \in$
 379 $D\}$, with $\sigma_\delta^2 = C_\delta(\mathbf{0})$; recall that $C_\delta(\mathbf{h}) = \text{cov}(\delta(\mathbf{s} + \mathbf{h}), \delta(\mathbf{s}))$, for $\mathbf{h} \in \mathbb{R}^d$. Let $\hat{C}_\delta(\cdot)$
 380 denote the empirical covariance function, and $s_\delta^2 = \hat{C}_\delta(\mathbf{0})$. Then, according to Principle
 381 3.2, *approximately* speaking, $s_\mu^2 + s_\delta^2$ should not change as different decompositions
 382 given by (14) are fitted. Now suppose that the additional covariate vector, $\mathbf{x}_{p+1} =$
 383 $(x_{p+1}(\mathbf{s}_1), \dots, x_{p+1}(\mathbf{s}_n))^\top$, is included in the regression, and it is not in the column space
 384 of \mathbf{X} . Then generally, the new empirical covariance function $\hat{C}_\delta(\cdot)$, yields a new estimate
 385 $s_\delta^2 = \hat{C}_\delta(\mathbf{0})$. The principle says that this new s_δ^2 should decrease to allow for the additional
 386 spatial variability in $\{\mu(\mathbf{s}) : \mathbf{s} \in D\}$. (Understandably, the general shape of the new empirical
 387 covariance function will change as well.)

388 How will we ever know which model is better after each is fitted to the same spatial
 389 data $\{z(\mathbf{s}_i)\}$? Model-selection criteria such as ‘Akaike Information Criterion,’ ‘Deviance
 390 Information Criterion,’ and cross-validation will remove bad models but, in the ‘difficult
 391 middle,’ there is no way to know whether a graph of $\{z(\mathbf{s}_i)\}$ versus $\{x_{p+1}(\mathbf{s}_i)\}$ is showing
 392 the behaviour of a component of the mean function or the behaviour of a component of the
 393 error process. Germane to Principle 3.2, it is well known to spatial statisticians that a single
 394 simulation of a strongly spatially dependent random effect, $\{\delta(\mathbf{s}) : \mathbf{s} \in D\}$, can look like
 395 deterministic spatial trend. If a replicate of the data were available, and if the analogous
 396 plot showed the same behaviour as seen in the original plot, then $\{x_{p+1}(\mathbf{s}) : \mathbf{s} \in D\}$ would
 397 probably belong in the mean function $\mu(\cdot)$. From just one realisation, the decomposition (14)

398 is not unique, and hence what is one person's mean function could be another person's spatial
 399 error.

400 3.3. COS is the DNA of Spatial Statistics

401 Most of us have heard of, or encountered, *Simpson's Paradox*, which basically says that
 402 two variables x and y could be positively correlated but, *conditional* on a third variable w ,
 403 it is perfectly acceptable that x and y be negatively correlated (or uncorrelated, or positively
 404 correlated)! Usually, Simpson's Paradox is expressed in terms of categorical data, which we
 405 could do here by defining ordinal categories, or bins, from the ranges of w , x , and y .

406 Let a sample $\{(w_i, x_i, y_i) : i = 1, \dots, n\}$ be assigned to the pre-defined bins, creating a
 407 three-way contingency table with counts in each cell of the table. A two-way table showing
 408 (marginal) dependence between binned x and binned y is obtained by aggregating the counts
 409 across the third variable, namely binned w . Then a measure of dependence in the two-way
 410 table (e.g., the statistic *gamma*, due to Goodman & Kruskal 1954) could show positive
 411 dependence, but the same measure applied to the two-way tables of binned x and binned
 412 y *conditional* on each of the values of binned w , could all show negative (or no, or positive)
 413 dependence.

414 What can you do about it? First, understand why it happens, and then spend the rest of
 415 your data-science career looking for those lurking variables like w that you may have missed!
 416 It manifests in other settings as well, as I now discuss.

Let x and y be jointly Gaussian random variables and related through the simple linear regression model given by (12); the probability model of how y varies with x is as follows: Conditional on x , the random variable y is Gaussian with mean and variance, respectively,

$$\begin{aligned} E(y|x) &= E(y) + \{\rho_{xy}\text{var}(y)^{1/2}/\text{var}(x)^{1/2}\}\{x - E(x)\}, \\ \text{var}(y|x) &= \text{var}(y)\{1 - \rho_{xy}^2\}, \end{aligned}$$

417 where $\rho_{xy} = \text{corr}(x, y)$.

418 Now consider regression of y on both x and w , where again joint Gaussianity of
 419 random variables w , x , and y is assumed. To make the calculations easier to follow, in
 420 the rest of this subsection consider the simple case where $E(w) = E(x) = E(y) = 0$, and
 421 $\text{var}(w) = \text{var}(x) = \text{var}(y) = 1$. Then the covariances, $\text{cov}(x, y)$, $\text{cov}(w, y)$, and $\text{cov}(x, w)$
 422 are identical to correlations, which are denoted here as ρ_{xy} , ρ_{wy} , and ρ_{xw} , respectively. The
 423 regression lines to be compared are,

$$E(y|x, w) = \left\{ \frac{\rho_{wy} - \rho_{xy}\rho_{xw}}{1 - \rho_{xw}^2} \right\} w + \left\{ \frac{\rho_{xy} - \rho_{wy}\rho_{xw}}{1 - \rho_{xw}^2} \right\} x, \quad (15)$$

and

$$E(y|x) = 0 + \rho_{xy} \cdot x.$$

Then, conditional on w ,

$$\begin{aligned} \text{corr}(x, y|w) &= \left\{ \frac{\rho_{xy} - \rho_{wy}\rho_{xw}}{1 - \rho_{xw}^2} \right\} \frac{\text{var}(x|w)^{1/2}}{\text{var}(y|w)^{1/2}} \\ &= \left\{ \frac{\rho_{xy} - \rho_{wy}\rho_{xw}}{(1 - \rho_{xw}^2)^{1/2}(1 - \rho_{wy}^2)^{1/2}} \right\}, \end{aligned} \quad (16)$$

424 which is to be compared to $\text{corr}(x, y) = \rho_{xy}$.

425 If w has zero correlation with both x and y , then from (16), $\text{corr}(x, y|w) = \text{corr}(x, y) =$
 426 ρ_{xy} , which is an intuitively reasonable result. In general, $\rho_{wy} \neq 0$ and $\rho_{xw} \neq 0$, so from (15)
 427 it is clear that a lurking variable w can wreak havoc on any honest attempt to interpret a
 428 simple linear regression of y on x . But, what does Simpson's Paradox have to do with spatial
 429 statistics?

430 The answer is 'plenty,' if you think about w as being a variable that describes a range
 431 of geographical strata. For example, the Australian Bureau of Statistics divides Australia
 432 up into *small areas*, at different levels of aggregation. A study of y = mean weekly income
 433 in NSW, Australia, regressed on selected demographic variables x at the finest level of
 434 aggregation, was given in Burden, Cressie & Steel (2015). Given the great divide in Australia
 435 between 'city' and 'country,' the most obvious variable w to choose is: $w = 1$ if the small
 436 area is in Greater Sydney (city), and otherwise $w = 0$ (country). In this set-up, it is easy to
 437 imagine that the two regressions, $E(y|x, w = 1)$ and $E(y|x, w = 0)$, would result in more
 438 interpretable results than the single regression, $E(y|x)$. (In fact, Burden, Cressie & Steel
 439 2015, used Principle 3.2 and captured the 'geography' by using a spatial error term rather
 440 than a geographic covariate w .)

441 Simpson's Paradox is potentially everywhere in the spatial context, because regressions
 442 of y on x can be done at different levels of spatial aggregation. A regression of y on x
 443 at the finest level of aggregation may show positive dependence, but when the variables
 444 are aggregated to a coarser level, the regression may show negative (or no, or positive)
 445 dependence! In Sociology, this phenomenon has been called the *ecological effect* (Robinson
 446 1950), an unfortunate and misleading name that has no direct connection to Ecology. In
 447 Geography, the combination of Simpson's Paradox and the ecological effect has been called
 448 the Modifiable Areal Unit Problem, or MAUP; a spatial statistical discussion of MAUP is
 449 given in Cressie (1996).

Because it is so natural to aggregate processes over space, meta-data in the spatial data
 set should include the *spatial support* of each datum. Define $y(B) = (1/|B|) \int_B y(s) ds$;

then $y(B)$ has *spatial support* B with volume, $|B| = \int_B ds > 0$. In geostatistics, B is called a *block*, and we say that $y(B)$ is an aggregation (or block average) of the original process. Then a *change-of-support* (COS) is said to have occurred from point support to spatial support B , with a corresponding change of statistical properties. Data on mutually exclusive supports, $\{B_1, \dots, B_n\}$, are typically represented as:

$$z(B_i) = y(B_i) + \varepsilon_i; i = 1, \dots, n.$$

450 The simplest model would be $\{\varepsilon_i : i = 1, \dots, n\}$ iid $N(0, \sigma_\varepsilon^2)$, but modifications to account
 451 for the protocols of measurement and the possible overlap of the supports $\{B_i\}$ have been
 452 developed (e.g., Wikle & Berliner 2005).

453 The two earlier principles of Spatial Statistics (nearby things are more alike; and
 454 decompose spatial variability into fixed effects plus random effects) are joined here by a
 455 change-of-support principle.

456 **Principle 3.3:** *The variance of the aggregation, $y(B)$, is a decreasing function of the volume,*
 457 *$|B|$.*

458 This is the most important of a number of COS properties (e.g., Gotway & Young 2002),
 459 whose discussion I have left for a future time.

460 Let $y(\mathbf{s}) = y_0 + \delta(\mathbf{s})$, for $\mathbf{s} \in D$, where y_0 is a non-degenerate random variable with
 461 possibly non-zero mean and independent of $\{\delta(\mathbf{s}) : \mathbf{s} \in D\}$. Then for $B \subset D$,

$$\text{var}(y(B)) = \sigma_0^2 + \text{var}(\delta(B)), \quad (17)$$

462 where $\sigma_0^2 = \text{var}(y_0) > 0$. Therefore, if $\text{var}(\delta(B))$ decreases to 0 as the volume $|B|$ increases,
 463 $\text{var}(y(B))$ decreases to $\sigma_0^2 > 0$. It is this type of behaviour that is of interest to geoscientists
 464 analysing remote sensing data. In that literature (reviewed in Cressie 2018), they distinguish
 465 between two types of error, ‘systematic error’ and ‘random error,’ as follows.

466 *Random error:* An error is a random error if the average of a collection of n of them has
 467 variability that decreases to zero like $1/n$, as n becomes large.

468 *Systematic error:* An error is a systematic error if the average of a collection of n of them has
 469 variability that does not decrease to zero, as n becomes large.

470 These might be considered ‘verbal working definitions,’ but a statistical scientist looks
 471 at these and tries to find an appropriate probability model. The ones I give below are building
 472 blocks for parts of a HM. In practice, the most difficult aspect is to know which errors are of
 473 which type and how to group them together for inference.

474 For ‘random error,’ the obvious probability model is: Let $\delta(\mathbf{s}_1), \dots, \delta(\mathbf{s}_n)$ be iid
 475 random variables with mean 0 and variance σ_δ^2 . Now the average, $\bar{\delta} = (\sum_{i=1}^n \delta(\mathbf{s}_i))/n$, has

476 $E(\bar{\delta}) = 0$ and $\text{var}(\bar{\delta}) = \sigma_{\delta}^2/n$, which decreases to zero like $1/n$. In a ‘data-rich’ situation
 477 and with enough averaging of the data, random errors can be shown to be annihilated by the
 478 averaging (using the law of large numbers). In a spatial setting, Principle 3.1 suggests that the
 479 independence assumption between the $\{\delta(\mathbf{s}_i)\}$ is not appropriate. Under increasing-domain
 480 asymptotics, it can be shown that $\text{var}(\bar{\delta})$ still decreases like $1/n$ (Cressie 1993), however
 481 strong spatial dependence can make the errors look more systematic than random (Morris &
 482 Ebey 1984).

For ‘systematic error,’ an obvious probability model is a random-effects model: Let $e(\mathbf{s}_1), \dots, e(\mathbf{s}_n)$ be written as

$$e(\mathbf{s}_i) = y_0 + \delta(\mathbf{s}_i); \quad i = 1, \dots, n,$$

where y_0 has mean 0 and variance $\sigma_0^2 > 0$, and y_0 and $\{\delta(\mathbf{s}_i)\}$ are independent. Now, if $\{\delta(\mathbf{s}_i)\}$ are iid mean 0 and variance σ_{δ}^2 , then the average error $\bar{e} = \sum_{i=1}^n e(\mathbf{s}_i)/n$ has variance,

$$\text{var}(\bar{e}) = \sigma_0^2 + \sigma_{\delta}^2/n.$$

483 This does not decrease to zero as n becomes large, and hence $e(\cdot)$ is a form of systematic
 484 error. Zhang, Cressie & Wunch (2019) used spatial dependence in the $\{\delta(\mathbf{s}_i)\}$ and an additive
 485 random effect y_0 with mean 0, as part of the OCO-2 calibration of remote sensing data to
 486 ‘ground-truth’ data.

487 Consider spatial prediction errors, $\{\hat{y}(\mathbf{s}_i) - y(\mathbf{s}_i) : i = 1, \dots, n\}$, which are located
 488 at $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$ in the spatial domain D . Typically, these prediction errors have individual
 489 means that are not zero, and it is often the case in Spatial Statistics that there is no
 490 replication to estimate them. A way out of the conundrum is to assume exchangeability
 491 (e.g., Spiegelhalter, Abrams & Myles 2004, Section 3.17), which results in a spatial statistical
 492 model that exhibits both systematic error and random error; see Cressie (2018, Rejoinder).

493

4. Epilogue

494 In the previous sections, I presented five principles for Statistical Science and three
 495 special ones for Spatial Statistics. In what follows, I add a ninth principle that speaks to
 496 all of Science, not just Data Science, and I give some concluding remarks.

497 4.1. You can’t add apples and oranges

498 There is one further principle that has been my ‘rock’ in the diverse applied-statistics
 499 projects in which I have participated. I once saw a cartoon that showed a small town’s
 500 ‘Welcome’ sign, and it looked something like this:

Centerville

501

| | | |
|------------|---|----------|
| Population | : | 2,390 |
| Elevation | : | 862 feet |
| Total | : | 3,252 |

502 Keeping track of units is a fundamental part of all of science, including Data Science.

503 **Principle 4.1:** *Only add quantities that have the same units. Multiply quantities that have*
 504 *different units, and cancel units from the product whenever possible. Take logs, exponents,*
 505 *and other special functions of unit-free quantities.*

506 The first part of the principle is illustrated with the cartoon I referred to above. You might
 507 say you would never add apples and oranges but, if in (7), z is measured in petagrams (Pg) of
 508 carbon in Earth's atmosphere and $x_1 = 1$, $x_2 = t$ (in years), and $x_3 = t^2$, then a unit analysis
 509 using the second part of the principle would reveal that regression coefficient β_1 is in units of
 510 Pg, β_2 is in Pg/yr, and β_3 is in Pg/yr². Principle 4.1 is respected, since it is the regression
 511 components $\{\beta_k x_k\}$ that are added. However, it is the regression coefficients $\{\beta_k\}$ that
 512 are often interpreted, so I often standardise the covariates by, respectively, subtracting their
 513 averages and dividing by their sample standard deviations. Then, after this standardisation,
 514 $\beta_1, \beta_2, \dots, \beta_k$ all have the same units as those of z .

515 Statistical scientists know that probabilities have no units. Using Principle 4.1, a unit
 516 analysis of the fundamental probability equation, $\int f(y) dy = 1$, where $f(\cdot)$ is the probability
 517 density function of the random variable y , whose units are, say, Pg of carbon, reveals that
 518 $f(y)$ has units of (Pg)⁻¹. Using the same principle, $E(y)$ has units of Pg, and so forth. While
 519 probability density functions have units, cumulative distribution functions and probability
 520 mass functions do not.

The last part of Principle 4.1 applies to any special function that can be expressed as a Taylor series. The most common ones in science are logs and exponents, whose Taylor series are

$$\log_e(1+x) = x - x^2/2 + x^3/3 - \dots$$

and

$$e^x = 1 + x + x^2/2! + x^3/3! + \dots,$$

521 which only make sense when x has no units; see the first part of Principle 4.1. Euler's number,
 522 e , has no units, because $e = \lim_{n \rightarrow \infty} (1 + 1/n)^n$.

Counts, percentages, and correlations have no units. Also, the Box-Cox transformation (Box & Cox 1964),

$$g_{\lambda}(x) = \begin{cases} (x^{\lambda} - 1)/\lambda, & \text{for } \lambda \in (-\infty, 0) \cup (0, \infty) \\ \log x, & \text{for } \lambda = 0 \end{cases},$$

523 only makes sense if x has been rendered unit-free. One way to accomplish this is to specify
524 x relative to a given standard.

525 Every term in the process model (and data model) in a HM, should be assigned their
526 rightful units. The scientific models that make up parts of the process model are sometimes
527 derived theoretically, sometimes empirically. Beware of a ‘scientific constant.’ It may be an
528 (estimated) regression coefficient, in which case its units (the ratio of the dependent variable’s
529 units to the corresponding covariate’s units) are key, since that ‘constant’ might change if the
530 units change.

531 In conclusion, it is a critical part of every collaboration to do a ‘unit analysis,’ which
532 avoids obvious mistakes made by confusing different parts of different systems of units (e.g.,
533 metric and imperial), as well as the more subtle ones discussed above.

534 4.2. These are my nine principles

535 **Principle 2.1:** *Establish a true model (TM), perhaps different from the scientist’s working*
536 *model (WM). Critically, compute the TM-distributional properties of the WM estimators. [All*
537 *models are wrong . . . but some are wrong-er than others.]*

538 **Principle 2.2:** *Build statistical models conditionally, through a data model and a process*
539 *model. Infer the unknown process from the predictive distribution. [What you see is not what*
540 *you want to get.]*

541 **Principle 2.3:** *In any well defined statistical model, there is conservation of variability.*
542 *[Geophysicists conserve energy but what do data scientists conserve?]*

543 **Principle 2.4:** *Seek a transformation of the scientific process where all components of*
544 *variation act and interact additively. [The holy grail: all scales of variation are additive.]*

545 **Principle 2.5:** *When building probability models, look carefully where zero probabilities*
546 *are assumed (perhaps implicitly) and, with the same care, move appropriate probabilities*
547 *away from zero. Calculate joint probabilities from products of conditional (not marginal)*
548 *probabilities, unless entropy is maximal. [Could swans be red?]*

549 **Principle 3.1** *Everything is related to everything else, but near things are more related than*
550 *distant things (Tobler 1970). [Patches in close proximity are commonly more alike . . .]*

551 **Principle 3.2:** *Assume a decomposition of a spatial process into fixed effects plus random*
552 *effects. While it is not unique, the decomposition must be chosen to conserve variability.*
553 [What is one person's mean function could be another person's spatial error.]

554 **Principle 3.3:** *The variance of the aggregation, $y(B)$, is a decreasing function of the volume,*
555 *$|B|$. [COS is the DNA of Spatial Statistics.]*

556 **Principle 4.1:** *Only add quantities that have the same units. Multiply quantities that have*
557 *different units, and cancel units from the product whenever possible. Take logs, exponents,*
558 *and other special functions of unit-free quantities. [You can't add apples and oranges.]*

559 **4.3. A disclaimer**

560 These nine principles of Statistical Science are personal, leading to a certain amount
561 of self-referencing, but I hope others find them useful. There are no theorem-proof results,
562 but there are back-of-the-envelope calculations that I use to justify the principles in simple
563 settings. At the very least, they should give data scientists boundaries in their analytics that
564 should be respected, and criteria by which supervised and unsupervised machine-learning
565 methods could be assessed.

566 **4.4. Happy birthday!**

567 I once introduced Adrian Baddeley at a conference session I was chairing (ASC2010,
568 Fremantle, Western Australia) as a national treasure. I repeat it here, and I wish him a very
569 happy 65th birthday!

570

References

- 571 BERGER, J.O. (1985). *Statistical Decision Theory*. Springer, New York, NY.
- 572 BERLINER, L.M., WIKLE, C.K. & CRESSIE, N. (2000). Long-lead prediction of Pacific SSTs via Bayesian
573 dynamic modeling. *Journal of Climate* **13**, 3953–3968.
- 574 BOX, G.E.P. & COX, D.R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society,*
575 *Series B* **26**, 211–252.
- 576 BURDEN, S., CRESSIE, N. & STEEL, D. (2015). The SAR model for very large datasets: A reduced-rank
577 approach. *Econometrics* **3**, 317–338.
- 578 CONNOR, B., BÖSCH, H., MCDUFFIE, J., TAYLOR, T., FU, D., FRANKENBERG, C., O’DELL, C., PAYNE,
579 V.H., GUNSON, M., POLLOCK, R., HOBBS, J., OYAFUSO, F. & JIANG, Y. (2016). Quantification of
580 uncertainties in OCO-2 measurements of XCO₂: Simulations and linear error analysis. *Atmospheric*
581 *Measurement Techniques* **9**, 5227–5238.
- 582 CRESSIE, N. (1978). Removing nonadditivity from two-way tables with one observation per cell. *Biometrics*
583 **34**, 505–513.
- 584 CRESSIE, N. (1985). When are relative variograms useful in geostatistics? *Journal of the International*
585 *Association for Mathematical Geology (now Mathematical Geosciences)* **17**, 693–702.
- 586 CRESSIE, N. (1993). *Statistics for Spatial Data*, rev. edn. Wiley, New York, NY.
- 587 CRESSIE, N. (1996). Change of support and the modifiable areal unit problem. *Geographical Systems* **3**,
588 159–180.
- 589 CRESSIE, N. (2018). Mission CO₂ntrol: A statistical scientist’s role in remote sensing of atmospheric carbon
590 dioxide (with discussion and a rejoinder). *Journal of the American Statistical Association* **113**, 152–181.
- 591 CRESSIE, N. (2020). When is it Data Science and when is it Data Engineering? Comment on “Prediction,
592 Estimation, and Attribution” by B. Efron. *Journal of the American Statistical Association*, **115**, 660–
593 662.
- 594 CRESSIE, N., RICHARDSON, S. & JAUSSENT, I. (2004). Ecological bias: Use of maximum-entropy
595 approximations. *Australian and New Zealand Journal of Statistics* **46**, 233–255.
- 596 CRESSIE, N. & WIKLE, C.K. (2011). *Statistics for Spatio-Temporal Data*. John Wiley & Sons, Hoboken,
597 NJ.
- 598 FISHER, R.A. (1935). *The Design of Experiments*. Oliver and Boyd, Edinburgh, UK.
- 599 GOODMAN, L.A. & KRUSKAL, W.H. (1954). Measures of association for cross classifications. *Journal of*
600 *the American Statistical Association* **49**, 732–764.
- 601 GOTWAY, C. & YOUNG, L.J. (2002). Combining incompatible spatial data. *Journal of the American*
602 *Statistical Association* **97**, 632–648.
- 603 LAHIRI, S.N., KAISER, M.S., CRESSIE, N. & HSU, N.J. (1999). Prediction of spatial cumulative
604 distribution functions using subsampling (with discussion and a rejoinder). *Journal of the American*
605 *Statistical Society* **94**, 86–110.
- 606 MACKENZIE, D. (2020). *Covid-19: The Pandemic that Should Never Have Happened, and How to Stop the*
607 *Next One*. The Bridge Street Press, London, UK.
- 608 MATHERON, G. (1963). *Traité de Geostatistique Appliquée, Tome II: Le Krigeage*. Memoires du Bureau de
609 Recherche Géologique et Minière, No. **24**. Paris, France.
- 610 MCCULLAGH, P. & NELDER, J.A. (1989). *Generalized Linear Models, 2nd edn*. John Wiley & Sons, New
611 York, NY.
- 612 MORRIS, M.D. & EBEL, S.F. (1984). An interesting property of the sample mean under a first-order
613 autoregressive model. *The American Statistician* **38**, 127–129.
- 614 NGUYEN, H., CRESSIE, N. & HOBBS, J. (2019). Sensitivity of Optimal Estimation satellite retrievals
615 to misspecification of the prior mean and covariance, with application to OCO-2 retrievals. *Remote*
616 *Sensing* **11**, 2770.

- 617 ROBINSON, W.S. (1950). Ecological correlations and the behavior of individuals. *American Sociological*
618 *Review* **15**, 351–357.
- 619 SPIEGELHALTER, D.J., ABRAMS, K.R. & MYLES, J.P. (2004). *Bayesian Approaches to Clinical Trials and*
620 *Health-Care Evaluation*. John Wiley & Sons, Chichester, UK.
- 621 TALEB, N. (2007). *The Black Swan: The Impact of the Highly Improbable*. Random House, New York, NY.
- 622 TOBLER, W.R. (1970). A computer movie simulating urban growth in the Detroit region. *Economic*
623 *Geography* **46**, 234–240.
- 624 WIKLE, C.K. & BERLINER, L.M. (2005). Combining information across spatial scales. *Technometrics* **47**,
625 80–91.
- 626 ZHANG, B., CRESSIE, N. & WUNCH, D. (2017). Statistical properties of atmospheric greenhouse gas
627 measurements looking down from space and looking up from the ground. *Chemometrics and Intelligent*
628 *Laboratory Systems* **162**, 214–222.
- 629 ZHANG, B., CRESSIE, N. & WUNCH, D. (2019). Inference for errors-in-variables models in the presence of
630 systematic errors with an application to a remote sensing campaign. *Technometrics* **61**, 187–201.