

U.S. Department of Commerce
National Oceanic and Atmospheric Administration
National Weather Service
National Centers for Environmental Prediction
5830 University Research Court
College Park, MD 20740

Office Note 481

doi:10.7289/V5736NVN

**A THEORETICAL EXAMINATION OF THE CONSTRUCTION AND
CHARACTERIZATION OF SUPER-OBSERVATIONS OBTAINED BY
OPTIMALITY PRINCIPLES GUIDED BY INFORMATION THEORY**

R. James Purser*
IM Systems Group, Rockville, Maryland

May 18, 2015

THIS IS AN UNREVIEWED MANUSCRIPT, PRIMARILY INTENDED FOR INFORMAL
EXCHANGE OF INFORMATION AMONG THE NCEP STAFF MEMBERS

* email: jim.purser@noaa.gov

Abstract

This note characterizes the optimal construction of (possibly) multi-component super-observation (or ‘super-obs’) based upon the criterion of minimizing the information lost in the super-obsing process. It is asserted that, by an artificial intervention that adjusts the weights given to the super-ob, it is possible to ‘structurally precondition’ the assimilation problem to speed up the convergence of a Krylov-based iterative minimization (such as the conjugate gradient method, for example) without significantly changing the convergent limit of the process. By an examination of this optimal formulation in the context of a compact cluster of point data it is shown that, to the leading approximation in an asymptotic scaling parameter describing the cluster’s size, the optimal multi-component super-ob is essentially identical to the multipole characterization of generalized super-obs suggested (on an intuitive basis) in an earlier note by Purser, Parrish and Masutani.

1. INTRODUCTION

For dense satellite sounding data and for fields of radar reflectivity the cost of assimilating each individual measurement outweighs the marginal benefit that the given measurement, amongst so many others like it, confers on the resulting assimilation. Discarding a large proportion of such data is an economical solution which entails some loss of information. In some circumstances, a better resolution of this problem is to combine the data in sets of appropriate sizes and to characterize the collective measurements within each set by an appropriate average value, or a very much reduced set of summarizing values, based on an averaging procedure designed to minimize the information loss in a more intelligent way than is implied by blindly thinning. The compression of information from a cluster of raw data into a smaller set of representative quantities and their associated precision weights gives rise to what is known as a “super-observation” (Lorenz 1981), usually abbreviated to “super-ob”. More generally in a hierarchy of nested clusters, the contributing data of a higher level super-ob can themselves include super-obs.

We propose that an examination of the independent precision-weight values of each super-ob, followed by an intervention that artificially caps these values to a judiciously determined upper limit, will have the effect of ‘structurally preconditioning’ the assimilation system – essentially reducing the largest eigenvalue of the Hessian of the minimization process by slightly altering the problem being solved. Formally, the limit of convergence of the iterative minimization will also be slightly altered in this process, but the perturbation induced by this structural preconditioning will be, in practice, so small as to be insignificant. If the preconditioning is effective, then we can expect that the very large condition number of the original problem will be reduced (very substantially in some cases) so that, after any given *finite* number of conjugate gradient, or Krylov, descent iterations (van der Vorst 2003), the overall fit of the assimilation to the truth will be, on average, improved.

In general, a super-ob can be a multi-component object. In this note we state the conditions under which a super-ob of this type can be said to be ‘optimal’ subject to constraints on the

number of components it possesses. The optimality we invoke is interpreted as minimizing the loss of information in a formal sense. In terms of this definition we also show how super-obs can be given an asymptotic characterization. The leading terms in each component gives rise to an approximately optimal super-ob of exactly the multi-pole form proposed in Purser et al. (2000).

2. INFORMATION-THEORETIC CHARACTERIZATION OF OPTIMAL SUPER-OBSERVATIONS

Consider a compact cluster of m observations which we wish to replace by a smaller number of virtual measurements constructed to emulate the dominant effects of the original set. We assume, for simplicity, that the statistics of errors are Gaussian and that the observation operators are linear. We shall be concerned, initially, only with the marginal effect of the data compression of this one cluster, assuming that all other data are already assimilated (with or without super-obbing). Also, since the Gaussianity, and linear operator assumptions allow us to treat the assimilated data sequentially, we can assume that the effective background for the assimilation of our cluster is the optimally assimilated field of values that has already combined the original background with all data except those belonging to our cluster. The effective background error is therefore not the original background error but rather the error of this partial assimilation. However, we shall idealize the covariance statistics of this partially assimilated field by assuming that, at the location of our cluster at least, the error covariance of what has become the effective background is spatially homogeneous and smooth so that it might be expanded in a power series in spatial coordinates. The values attributed to the composite (multi-component) super-ob that stands in for the chosen cluster are also assumed to be linear in the values of the contributing raw observations.

The criteria we shall examine to guide our construction of the super-obs will be based upon such measures of quality as are exemplified by Shannon’s information entropy (Shannon 1948, Shannon and Weaver 1949, Cover and Thomas 1991) and by the ‘degrees of freedom for signal’, (DFS). Entropy has been used as a criterion for quantifying the information content of observing systems, e.g., Peckham (1974), Eyre (1990), Rodgers (1998). Criteria of these kinds have more generally served to focus studies in the field of statistical ‘optimal design’, where they are referred to as ‘trace criteria’ (Kiefer 1974, Pázman 1986, Purser and Huang 1993). There exists a whole spectrum of such trace criteria, information entropy and DFS being just two special cases. Conveniently, the same choice of super-obbing definition and composition weights is obtained for any particular choice from this selection of trace criteria. The reason for this will be made clear below.

Let the m individual data belonging to the compact cluster be denoted by the vector, \mathbf{y} , the field of partially assimilated data (not accounting for this cluster), by the ‘effective background’ vector in state space, \mathbf{x}^b , and the corresponding state space field of fully assimilated data by \mathbf{x}^a . Denote the covariance operator for the errors of \mathbf{y} by \mathbf{R} , covariances for the errors of \mathbf{x}^b by \mathbf{B} and the linearized measurement operator, \mathbf{H} . As is well known, the optimal analysis can be written:

$$\mathbf{x}^a = \mathbf{x}^b + (\mathbf{B}^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H} \mathbf{x}^b). \quad (2.1)$$

Unfortunately, since \mathbf{B}^{-1} , even when it is the inverse of a homogeneous background covariance, is not practically expressible in explicit terms. It’s therefore customary to express the analysis

increment, $\mathbf{x}^a - \mathbf{x}^b$, in terms of a new control variable, \mathbf{v} , and a (generally rectangular) matrix, \mathbf{C} :

$$\mathbf{x}^a - \mathbf{x}^b = \mathbf{C}\mathbf{v}, \quad (2.2)$$

where,

$$\mathbf{C}\mathbf{C}^\top = \mathbf{B}. \quad (2.3)$$

We can also conveniently write \mathbf{W} to denote the symmetric operator of ‘precision weights’ for the cluster of measurements,

$$\mathbf{W} = \mathbf{R}^{-1}, \quad (2.4)$$

and

$$\mathbf{d} = \mathbf{y} - \mathbf{H}\mathbf{x}^b, \quad (2.5)$$

to denote the vector of data ‘innovations’. Then the assimilation formula becomes:

$$\mathbf{v} = (\mathbf{I} + \mathbf{C}^\top \mathbf{H}^\top \mathbf{W} \mathbf{H} \mathbf{C})^{-1} \mathbf{C}^\top \mathbf{H}^\top \mathbf{W} \mathbf{d}. \quad (2.6)$$

The covariance, \mathbf{A} , of error of the assimilated field, \mathbf{x}^a , is given by,

$$\mathbf{A} = (\mathbf{B}^{-1} + \mathbf{H}^\top \mathbf{W} \mathbf{H})^{-1}, \quad (2.7)$$

so that the matrix operator,

$$\mathbf{A}^{-1} \mathbf{B} = \mathbf{I} + \mathbf{H}^\top \mathbf{W} \mathbf{H} \mathbf{B}, \quad (2.8)$$

describes quantitatively how the information of the background is augmented via the assimilation. For example, the trace of the logarithm of this quantity is directly proportional to the Shannon entropy:

$$\begin{aligned} S_0 &= \text{trace} \log(\mathbf{I} + \mathbf{H}^\top \mathbf{W} \mathbf{H} \mathbf{C} \mathbf{C}^\top) \\ &= \text{trace} \log(\mathbf{I} + \mathbf{M} \mathbf{M}^\top) \\ &= \text{trace} \log(\mathbf{I} + \mathbf{Q}) \end{aligned} \quad (2.9)$$

where, again for algebraic convenience,

$$\mathbf{M} = \mathbf{D} \mathbf{H} \mathbf{C}, \quad (2.10)$$

and

$$\mathbf{Q} = \mathbf{M} \mathbf{M}^\top, \quad (2.11)$$

with

$$\mathbf{D}^\top \mathbf{D} = \mathbf{W}. \quad (2.12)$$

A detailed derivation of (2.9) is provided in Appendix A.

The meaning of entropy in this context is the logarithm of the ratio of the squares of the effective state-space volumes of the background probability relative to the analysis probability; by choosing the logarithm here to have a base of four, the measure S_0 is therefore exactly in the conventional communication-theory units of ‘bits’. But in terms of natural logarithms

(more convenient for our purposes) the formula might be equivalently written, using one of the definitions of the natural logarithm function:

$$S_0 = \lim_{p \rightarrow 0} S_p, \quad (2.13)$$

where

$$S_p = \text{trace} \frac{\mathbf{I} - (\mathbf{I} + \mathbf{Q})^{-p}}{p}. \quad (2.14)$$

In this case, entropy is revealed to be just one of a continuum of the so called ‘trace-class’ of optimality criteria, as discussed by Kiefer (1974), Pázman (1986) and by Purser and Huang (1993). Other choices of $p > 0$ saturate at a finite measure of the information contributed by each ‘degree of freedom’, or eigen-component, of \mathbf{Q} . The special choice, $p = 1$, in this formula leads to a well-used formula for the ‘degrees of freedom for signal’ (Wahba, 1985, Purser and Huang 1993), because for this choice only, the saturation limit per degree of freedom is exactly one.

The trace-class of optimality criteria we have just described can be used to guide the construction of super-obs. The cluster of raw data will yield a value of S_p that is a measure of the full amount of information available (relative to what we have called the background). Any simplified summary of these data into some form of super-ob, being suboptimal, must generally yield a smaller S_p since some information is clearly being discarded. But we can use our freedom to choose the manner in which this super-ob is constructed to maximize the associated S_p within the limitations imposed by the constraints (such as the number, \bar{m} , of independent components possessed by this super-ob) under which this construction takes place.

We restrict ourselves to forms of super-obs whose component values are linear in the original data. We use an over-bar to denote the generally vector-valued super-ob and its associated operators. Thus, the values of the composite super-ob are components of vector, $\bar{\mathbf{y}}$, defined by:

$$\bar{\mathbf{y}} = \mathbf{L}^\top \mathbf{y}, \quad (2.15)$$

for some rectangular matrix, \mathbf{L} , whose $\bar{m} \leq m$ columns are the vectors of ‘combination weights’ defining each component of the super-ob. The corresponding measurement operator, $\bar{\mathbf{H}}$, is therefore:

$$\bar{\mathbf{H}} = \mathbf{L}^\top \mathbf{H}. \quad (2.16)$$

The effective weight for this super-ob is:

$$\bar{\mathbf{W}} = (\mathbf{L}^\top \mathbf{R} \mathbf{L})^{-1}, \quad (2.17)$$

and the new information measure is therefore:

$$\bar{S}_p = \text{trace} \frac{\mathbf{I} - (\mathbf{I} + \bar{\mathbf{Q}})^{-p}}{p}, \quad (2.18)$$

or its limiting value in the case of Shannon information, where,

$$\bar{\mathbf{Q}} = \bar{\mathbf{M}} \bar{\mathbf{M}}^\top, \quad (2.19)$$

with,

$$\overline{\mathbf{M}} = \overline{\mathbf{D}}\overline{\mathbf{H}}\overline{\mathbf{C}}, \quad (2.20)$$

and with

$$\overline{\mathbf{D}}^\top \overline{\mathbf{D}} = \overline{\mathbf{W}}. \quad (2.21)$$

Now \overline{S}_p can also be written:

$$\overline{S}_p = \sum_{a=1}^{\overline{m}} \frac{1 - (1 + \overline{\mu}_a^2)^{-p}}{p}, \quad (2.22)$$

where $\overline{\mu}_a^2$ are the eigenvalues of $\overline{\mathbf{Q}}$ or, equivalently, of what we can call the ‘matrix Rayleigh quotient’,

$$\overline{\mathbf{Q}} = \mathbf{L}^\top \mathbf{H} \mathbf{B} \mathbf{H}^\top \mathbf{L} (\mathbf{L}^\top \mathbf{R} \mathbf{L})^{-1}. \quad (2.23)$$

[We can also think of the $\overline{\mu}_a$ as being the positive singular values in the singular value decomposition of any factor, $\overline{\mathbf{M}}$, that satisfies (2.19).] By replacing the columns of \mathbf{L} by any invertible linear combination of them, say $\mathbf{L}' = \mathbf{L}\mathbf{G}$, for a square matrix \mathbf{G} for which $|\mathbf{G}| \neq 0$, we find that $\overline{\mathbf{Q}}$ is unchanged; the Rayleigh quotient (as in the classical scalar version of it) depends only upon the subspace spanned by the columns of \mathbf{L} and not at all upon the particular values of these vectors themselves. By writing,

$$\mathbf{L} = \mathbf{D}^\top \mathbf{V}, \quad (2.24)$$

the matrix Rayleigh quotient becomes:

$$\overline{\mathbf{Q}} = \mathbf{V}^\top \mathbf{Q} \mathbf{V} (\mathbf{V}^\top \mathbf{V})^{-1}. \quad (2.25)$$

Geometrically, the operator, $\overline{\mathbf{Q}}$, is a ‘section’ through the operator \mathbf{Q} , confining it to a subspace (defined by the span of \mathbf{L}) of dimension \overline{m} . The eigenvalues of this $\overline{\mathbf{Q}}$ are therefore subject to certain minimum-maximum properties of eigenvalues, which are described in Chapter 1, §4, of Courant and Hilbert (1989). Of these properties, the one of importance here is that the h -th eigenvalue (ordered from largest to smallest) of $\overline{\mathbf{Q}}$ cannot exceed the h -th eigenvalue of \mathbf{Q} . (This result is often referred to as “Courant’s mini-max principle” although it has been argued that the same result was known earlier independently by Rayleigh and Ritz.) Equality of the respective h -th eigenvalues is attained when the subspace spanned by the first h columns of \mathbf{L} (and \mathbf{V}) is exactly that spanned by the first h eigenvectors of \mathbf{Q} . For example, this result is obtained when the successive columns of \mathbf{V} are just the eigenvectors of \mathbf{Q} ordered according to eigenvalue dominance. Therefore, since the family of trace-class criteria (including DFS and Shannon entropy) are monotonic additive functions of these eigenvalues, we have shown:

Theorem 1

The ‘best’ choices of \overline{m} super-ob combination weights, \mathbf{L} , by the criteria either of maximizing the Shannon entropy or of maximizing the DFS, are those \mathbf{L} whose columns span the subspace formed by the \overline{m} most dominant eigenvectors of \mathbf{Q} . \square .

Remarks

When the columns of \mathbf{V} are chosen to be the dominant eigenvectors of \mathbf{Q} , each normalized to unit length, then the effective weight matrix, $\overline{\mathbf{W}}$, of the resulting composite super-ob is

simply the identity operator. But if, alternatively, we choose to normalize these eigenvectors such that,

$$\mathbf{V}^T \mathbf{V} = \boldsymbol{\mu}^{-2}, \quad (2.26)$$

then the magnitudes of the components of the corresponding diagonal weight matrix, obtained using (2.24) and (2.17):

$$\overline{\mathbf{W}} = \overline{\boldsymbol{\mu}}^2 \quad (2.27)$$

gives us an immediate quantitative assessment of the impact of these super-ob components on the assimilation system and, as we show in the next section, allows us to modify their weights (without changing the composition weights defined by \mathbf{L}) in a way that improves the condition number of the assimilation system as a whole without seriously risking the quality of the assimilation attained in the (presumably) more rapidly converging minimization of the cost function.

3. STRUCTURAL PRECONDITIONING BY SUPER-OB WEIGHT ADJUSTMENT

We now introduce a useful theorem on the eigenvalues of positive-definite or positive-semi-definite matrices:

Theorem 2

If \mathbf{A} and \mathbf{B} are two positive-definite square matrices of the same order then the ranked eigenvalues of their sum, $\mathbf{A} + \mathbf{B}$, are each greater than the corresponding ranked eigenvalues of \mathbf{A} ; if \mathbf{B} is merely positive semi-definite, then the ranked eigenvalues of the sum are not smaller than those of \mathbf{A} . \square .

Proof

Let \mathbf{V} be a normalized eigenvector with a the corresponding eigenvalue of a positive-definite matrix, \mathbf{A} , such that,

$$\mathbf{V}^T \mathbf{A} \mathbf{V} = a. \quad (3.1)$$

Let $d\mathbf{A}$ denote a positive-definite or semi-definite infinitesimal increment to \mathbf{A} , and consider the corresponding increments to the eigenvector and eigenvalue:

$$d\mathbf{V}^T \mathbf{A} \mathbf{V} + \mathbf{V}^T d\mathbf{A} \mathbf{V} + \mathbf{V}^T \mathbf{A} d\mathbf{V} = da. \quad (3.2)$$

But differentiating the normalization condition for the eigenvector we find that:

$$d\mathbf{V}^T \mathbf{V} = 0 \quad (3.3)$$

which, combined with the substitution of the eigenvector identity,

$$\mathbf{A} \mathbf{V} = \mathbf{V} a \quad (3.4)$$

implies that

$$\mathbf{V}^T d\mathbf{A} \mathbf{V} = da. \quad (3.5)$$

If the increment, $d\mathbf{A}$, is positive definite, then the evaluation of the term on the left of (3.5) is positive (though infinitesimal) and therefore, so is da . If $d\mathbf{A}$ is positive semi-definite, then the term on the left is non-negative, and therefore so is da . If we now integrate the infinitesimal

increments, $d\mathbf{A}$, to get a finite symmetric positive-definite, or positive-semi-definite, matrix, \mathbf{B} , the corresponding integrated, da , retains the same properties of positivity, and non-negativity, respectively. \square

One consequence of this theorem is that the largest \bar{m} ranked eigenvalues of the assimilation matrix, $\mathbf{I} + \mathbf{Q}$, corresponding to the inclusion of *all* data, are at least as great as the corresponding \bar{m} eigenvalues of any single multi-component optimal super-ob whose data contribute to this assimilation. In particular, under the almost universally true assumption that the *smallest* eigenvalue of the assimilation is essentially unity (Courtier 1997), the largest eigenvalue of any super-ob's $\bar{\mathbf{Q}}$ sets a practical lower bound to the condition number of the assimilation. If we find such eigenvalues that are very large, we have effectively identified a set of data which, while possibly innocuous individually, are collectively responsible for a severe conditioning difficulty.

The smallest $m - \bar{m}$ eigenvalues, μ_a^2 , of \mathbf{Q} , which are *not* represented amongst the eigenvalues of $\bar{\mathbf{Q}}$, contribute to the Shannon entropy deficit, $S_0 - \bar{S}_0$, and to the DFS deficit, $S_1 - \bar{S}_1$, representing information lost through this procedure. While there can be no hard rule, it is fair to say that:

(i) should the largest of these neglected μ_a^2 be comparable to or greater than unity, the adoption of super-obbing risks the loss of significant information potentially available in the raw data;

(ii) should an eigenvalue be *significantly less than unity*, then this degree of freedom of the measurement set can be neglected without detriment to the assimilation;

(iii) if an eigenvalue is *much larger than unity* (as can happen with dense clusters of very precise data), then by theorem 2, we can take this as a reliable diagnosis of severe ill-conditioning. The super-ob's contribution to this ill-conditioning can be mitigated, without serious damage to the assimilation, by artificially intervening to reduce this eigenvalue to a more reasonable 'capping value', without changing the associated eigenvector. Also by theorem 2, we can then infer that, by such an intervention, any change in the spectrum of the eigenvalues of the full assimilation and in the condition number itself (if it changes at all) must also be in the negative direction.

Taking the typical value of the capping limit – perhaps around $\mu_{\max}^2 = 10$, and recomputing the 'capped' effective eigenvalues, we can put them into a diagonal matrix, $\hat{\boldsymbol{\mu}}^2$, of order \bar{m} . Let the leading \bar{m} normalized eigenvectors of $\mathbf{Q} = \mathbf{M}\mathbf{M}^T$ be the columns of \mathbf{V} so that

$$\mathbf{V}^T \mathbf{Q} \mathbf{V} = \mathbf{V}^T \mathbf{D} \mathbf{H} \mathbf{B} \mathbf{H}^T \mathbf{D}^T \mathbf{V} = \hat{\boldsymbol{\mu}}^2, \quad (3.6)$$

Then, we might choose to make

$$\mathbf{L} = \mathbf{D}^T \mathbf{V} \hat{\boldsymbol{\mu}}^{-1}. \quad (3.7)$$

With this choice the projection of the background covariance into the space of the super-ob vector becomes the identity operator and the effective precision weight, without the 'capping' intervention, $\bar{\mathbf{W}}$ defined in (2.17), becomes

$$\bar{\mathbf{W}} = \hat{\boldsymbol{\mu}}^2. \quad (3.8)$$

Purser (1998) described algorithms for speeding up iterative approximations to the assimilation's descent problem that involved artificial manipulations of the spectrum of eigenvalues

similar to what is accomplished by these eigenvalue manipulations, and referred to such procedures as ‘structural preconditioning’. There it was shown how, in principle, the algorithms could be reformulated to ensure that the final convergent solution was still the exact minimum of the original (ill-conditioned) problem, although the eventual asymptotic rate of convergence could then never compete with that of the pure conjugate-gradient algorithm. Here we use ‘structural preconditioning’ to refer to just the simpler intervention that changes the descent algorithm’s eigenvalue spectrum *without* paying heed to the fact that, by intervening, we are also changing the convergent limit slightly (i.e., not strictly converging to the correct solution anymore). In the present context, then, ‘structural preconditioning’ of the super-ob involves simply replacing the precision matrix, $\overline{\mathbf{W}}$, of the super-ob by the capped version of it:

$$\hat{\mathbf{W}} = \hat{\boldsymbol{\mu}}^2, \quad (3.9)$$

where components of the diagonal matrix, $\hat{\boldsymbol{\mu}}^2$, of eigenvalues $\hat{\boldsymbol{\mu}}_a$ are now each defined:

$$\hat{\boldsymbol{\mu}}_a = \begin{cases} \mu_a & : \mu_a < \mu_{\max} \\ \mu_{\max} & : \mu_a \geq \mu_{\max} \end{cases}. \quad (3.10)$$

By doing so, not only do we guarantee that the computed solution to the optimal analysis with this substitution is typically indistinguishable from the true optimal solution, but we also guarantee that the conditioning of the assimilation as a whole has been improved. Since, in operational practice, the optimization involves a Krylov procedure (such as conjugate-gradients minimization) that is necessarily terminated after a finite number of iterations, the advantage of a sometimes dramatic reduction of the condition number of the problem effected by this eigenvalue capping intervention greatly outweighs the minor disadvantage of the (unattained) limit of the convergent iterations being very slightly sub-optimal. Noting first that the intervention in no way alters the definition of the super-ob combination weights, \mathbf{L} , we now turn to an examination of the form of these combination weights in the important asymptotic limit of a compact cluster of similar measurements and a smooth background error of much larger spatial characteristic scale.

4. ASYMPTOTIC CHARACTERIZATION OF COMPACT OPTIMAL MULTI-COMPONENT SUPER-OBS

The problem we examine in this section concerns the shape of the composition weights of an optimal composite super-ob formed from a compact cluster of similar measurements in the limit as the scale of the smooth covariance of background error goes to infinity. For simplicity, we assume that, through appropriate scaling, this scalar covariance is isotropic in spatial coordinates, \mathbf{X} , and can be expanded at locations \mathbf{X}_i and \mathbf{X}_j as a series in even degrees of separation: between two data indexed i and j where the scalar covariance is directly sampled:

$$(\mathbf{H}\mathbf{B}\mathbf{H}^\top)_{ij} = b_0 + \epsilon b_2 |\mathbf{X}_i - \mathbf{X}_j|^2 + \epsilon^2 b_4 |\mathbf{X}_i - \mathbf{X}_j|^4 + \dots, \quad (4.1)$$

the asymptotic parameter ϵ allows the other expansion coefficients to be kept order-unity (assuming that the coordinate origin is located close to the center of the cluster) while examining the asymptotic behavior as the covariance scale goes to infinity; ϵ is then essentially inversely proportional to the square of the covariance scale.

For each even p and $q \leq p$, and in any number of spatial dimensions, there exist coefficient matrices $\mathbf{a}_{q,p-q}$ that enable us to express each term of the expansion of $(\mathbf{H}\mathbf{B}\mathbf{H}^\top)_{ij}$ in terms of the row-vectors of the kind, $(\mathbf{v}_q)_i$, whose components consist of all products of degree q of the components of the position vector \mathbf{X}_i . For example, dropping the measurement index, i , we can consider the case of two spatial dimensions, $\mathbf{X}^\top = [X, Y]$, whereupon:

$$\mathbf{v}_0 = [1], \quad (4.2a)$$

$$\mathbf{v}_1 = [X, Y], \quad (4.2b)$$

$$\mathbf{v}_2 = [X^2, XY, Y^2], \quad (4.2c)$$

and so on, with the vector \mathbf{v}_q in d dimensions having $\binom{d+q-1}{q}$ components in the general case. Then for each even-valued p :

$$|\mathbf{X}_i - \mathbf{X}_j|^p = \sum_{q=0}^p (\mathbf{v}_q)_i \mathbf{a}_{q,p-q} (\mathbf{v}_{p-q})_j, \quad (4.3)$$

where the definitions of the matrices of coefficients \mathbf{a} are given recursively in Appendix B.

While we have defined $(\mathbf{v}_p)_i$ to be a row-vector of powers of the coordinates of a position vector, \mathbf{X}_i , of data point i , we can generalize \mathbf{v}_p (without the second subscript) to denote a *matrix* whose i -th row is that $(\mathbf{v}_p)_i$. Now, if we also use the Taylor series coefficients, b_p , of (4.1) to define new coefficient matrices that pertain specifically to this covariance model:

$$\mathbf{b}_{pq} = b_{p+q} \mathbf{a}_{pq}, \quad (4.4)$$

we can rewrite the expansion of the covariance \mathbf{B} as it appears projected as $\mathbf{H}\mathbf{B}\mathbf{H}^\top$ in the space of observations:

$$\begin{aligned} \mathbf{H}\mathbf{B}\mathbf{H}^\top &= \mathbf{v}_0 \mathbf{b}_{00} \mathbf{v}_0^\top + \epsilon (\mathbf{v}_0 \mathbf{b}_{02} \mathbf{v}_2^\top + \mathbf{v}_1 \mathbf{b}_{11} \mathbf{v}_1^\top + \mathbf{v}_2 \mathbf{b}_{20} \mathbf{v}_0^\top) \\ &\quad + \epsilon^2 (\mathbf{v}_0 \mathbf{b}_{04} \mathbf{v}_4^\top + \mathbf{v}_1 \mathbf{b}_{13} \mathbf{v}_3^\top + \mathbf{v}_2 \mathbf{b}_{22} \mathbf{v}_2^\top + \mathbf{v}_3 \mathbf{b}_{31} \mathbf{v}_1^\top + \mathbf{v}_4 \mathbf{b}_{40} \mathbf{v}_0^\top) + \dots \end{aligned} \quad (4.5)$$

However, we shall require that the coefficients, b_q , in the expansion of the covariance collectively conform to additional requirements which are sufficient to guarantee that \mathbf{B} has positive-definite projection, $\mathbf{H}\mathbf{B}\mathbf{H}^\top$, in *any* finite cluster of measurements. Introducing ‘even’ and ‘odd’ square block matrices:

$$\mathbf{b}_{[0:2s:2,0:2s:2]} = \begin{bmatrix} \mathbf{b}_{0,0} & \mathbf{b}_{0,2} & \cdots & \mathbf{b}_{0,2s} \\ \mathbf{b}_{2,0} & \mathbf{b}_{2,2} & \cdots & \mathbf{b}_{2,2s} \\ \dots & \dots & \dots & \dots \\ \mathbf{b}_{2s,0} & \mathbf{b}_{2s,2} & \cdots & \mathbf{b}_{2s,2s} \end{bmatrix}, \quad (4.6)$$

and

$$\mathbf{b}_{[1:2s+1:2,1:2s+1:2]} = \begin{bmatrix} \mathbf{b}_{1,1} & \mathbf{b}_{1,3} & \cdots & \mathbf{b}_{1,2s+1} \\ \mathbf{b}_{3,1} & \mathbf{b}_{3,3} & \cdots & \mathbf{b}_{3,2s+1} \\ \dots & \dots & \dots & \dots \\ \mathbf{b}_{2s+1,1} & \mathbf{b}_{2s+1,3} & \cdots & \mathbf{b}_{2s+1,2s+1} \end{bmatrix}, \quad (4.7)$$

we require that the determinants of such matrices are positive for all $s \geq 0$, i.e.,

$$|\mathbf{b}_{[0:2s:2,0:2s:2]}| > 0, \quad s \geq 0, \quad (4.8)$$

and

$$|\mathbf{b}_{[1:2s+1;2,1:2s+1;2]}| > 0, \quad s \geq 0. \quad (4.9)$$

These imply that each matrix is positive-definite and therefore non-singular.

We proceed by assuming that successive right-eigenvectors of $\mathbf{H}\mathbf{B}\mathbf{H}^\top\mathbf{W}$, are also smooth and each admit an asymptotic expansion in powers of the coordinates. It is convenient to generalize the eigen-problem so that we can deal simultaneously with the set of eigenvectors whose associated eigenvalues are of the same order of magnitude (as measured by the asymptotic parameter, ϵ). In other words, we attempt to identify distinct ‘eigen-spaces’ of the problem without necessarily striving to identify the spanning eigenvectors of each space individually. This modified eigen-problem is expressed by equations of the more lenient form than the classical eigenvector equation; we look for solutions to:

$$\mathbf{H}\mathbf{B}\mathbf{H}^\top\mathbf{W}\mathbf{U}_p = \mathbf{U}_p\mathbf{\Lambda}_p, \quad (4.10)$$

where \mathbf{U}_p is generally a *matrix* whose columns span an eigenspace, and where the quantity $\mathbf{\Lambda}_p$ is generally a square matrix (an ‘eigenmatrix’) instead of a single eigenvalue. The identified subspaces, indexed by p , will be segregated according to the asymptotic orders of magnitude attained by the contents of $\mathbf{\Lambda}$. The spanning array, \mathbf{U}_p , associated with each eigenspace is assumed to have a convergent Taylor expansion in spatial coordinates at each degree of an asymptotic covariance-scaling parameter, ϵ , but we do not require convergence of the asymptotic series itself. Thus, we can write an asymptotic expression true for any finite degree, \hat{r} , in the limit as $\epsilon \rightarrow 0$:

$$\mathbf{U}_p = \sum_{q=0}^{\infty} \sum_{r=0}^{\hat{r}} \epsilon^r \mathbf{v}_q \mathbf{U}_{p;qr} + \mathcal{O}(\epsilon^{\hat{r}+1}), \quad (4.11)$$

and we shall assume that the corresponding recombination matrix, $\mathbf{\Lambda}_p$, has an expansion:

$$\mathbf{\Lambda}_p = \sum_{r=p}^{\hat{r}} \epsilon^r \mathbf{\Lambda}_{p;r} + \mathcal{O}(\epsilon^{\hat{r}+1}). \quad (4.12)$$

We define the generalized Gram matrices,

$$\mathbf{A}_{[0;q,0;q]} = \begin{bmatrix} \mathbf{A}_{00}, & \cdots, & \mathbf{A}_{0q} \\ \cdots & \cdots & \cdots \\ \mathbf{A}_{q0}, & \cdots, & \mathbf{A}_{qq} \end{bmatrix}, \quad (4.13)$$

where,

$$\mathbf{A}_{pq} = \mathbf{v}_p^\top \mathbf{W} \mathbf{v}_q, \quad (4.14)$$

and define,

$$\mathbf{P}_{p;qr} = \sum_{s=0}^{\infty} \mathbf{A}_{qs} \mathbf{U}_{p;sr}. \quad (4.15)$$

Then we can examine the expansion of both sides of (4.10) to obtain, for each finite $\hat{r} \geq 0$:

$$\sum_{r=0}^{\hat{r}} \epsilon^r \left[\sum_{q=0}^{2r} \mathbf{v}_q \left(\sum_{s=0}^{[r-q/2]} \mathbf{b}_{q,2r-q-2s} \mathbf{P}_{p;2r-q-2s,s} \right) \right] = \sum_{r=p}^{\hat{r}} \epsilon^r \left[\sum_{q=0}^{2r} \left(\sum_{s=p}^r \mathbf{U}_{p;q,r-s} \mathbf{\Lambda}_{p;s} \right) \right] + \mathcal{O}(\epsilon^{\hat{r}+1}), \quad (4.16)$$

where the notation $[t]$ for a fractional index, t is always taken to mean “integer part of t ”.

First we note that no terms in ϵ^r for $r < p$ exist on the right-hand side of (4.16) to balance those we have put on the left-hand side. From this observation together with the asserted non-singularity of the positive-definite matrices of (4.6) and (4.7) we infer that:

$$\mathbf{P}_{p;qs} = \mathbf{0}, \quad 2s + q < p, \quad (4.17)$$

which, for each r in (4.16), restricts the upper limit of the q summation on the left-hand side of (4.16) to the finite maximum, $\hat{q} = 2r - p$. Now, since summed terms in $q > 2r - p$ must also vanish on right-hand side of (4.16) then, up to any p for which the leading eigen-matrix contribution, $\mathbf{\Lambda}_{p;p}$ continues to remain non-singular, we further deduce that:

$$\mathbf{U}_{p;qr} = 0, \quad \text{for } q - 2r > p. \quad (4.18)$$

This confirms not only that each Taylor series converges; it also proves that the Taylor series terminates at a finite q for each degree, r , of the asymptotic expansion. In each case where $r = p$ we find that the leading matrix term, $\mathbf{U}_{p;p0}$, in the expansion for \mathbf{U}_p has not more than as many columns as the number of columns in \mathbf{v}_p itself. We shall assume for simplicity that \bar{m} is large enough, up to some terminating p_{\max} , to ensure that there are enough remaining eigenvectors at each stage p to match the number of columns in the corresponding \mathbf{v}_p . The normalization of the summed terms for \mathbf{U}_p at each r is arbitrary; consequently, without any further loss of generality, we can conveniently choose these representative independent columns of \mathbf{U}_p such that,

$$\mathbf{U}_{p;pr} = \begin{cases} \mathbf{I} & : r = 0, \\ \mathbf{0} & : r > 0. \end{cases} \quad (4.19)$$

In the case where there are *not* enough remaining eigenvectors at the final stage $p = p_{\max}$, the normalization convention must be modified. For example, by replacing the identity operator on the right of (4.19) for $\mathbf{U}_{p;p0}$ by as many of the columns of \mathbf{I} as are needed.

We can consolidate all these results in a family of equations “ $\mathcal{E}(p; q, r)$ ” relating the left-hand side and right-hand side terms of each given q and r of (4.16):

$$\mathcal{E}(p; q, r) : \sum_{s=0}^{\lfloor r-q/2 \rfloor} \mathbf{b}_{q,2r-q-2s} \mathbf{P}_{p;2r-q-2s,s} = \sum_{s=p}^r \mathbf{U}_{p;q,r-s} \mathbf{\Lambda}_{p;s}, \quad q \leq 2r - p, \quad (4.20)$$

together with a rewriting of (4.15) as the appropriate *finite* summations in the equations that we now name “ $\mathcal{P}(p; q, r)$ ”:

$$\mathcal{P}(p; q, r) : \mathbf{P}_{p;qr} = \begin{cases} \sum_{s=0}^{p-1} \mathbf{A}_{qs} \mathbf{U}_{p;sr} + \mathbf{A}_{qp} & : r = 0, \\ \sum_{s=0}^{p-1} \mathbf{A}_{qs} \mathbf{U}_{p;sr} + \sum_{s=p+1}^{p+2r} \mathbf{A}_{qs} \mathbf{U}_{p;sr} & : r > 0. \end{cases} \quad (4.21)$$

We are now set to prove:

Theorem 3

Provided the geometrical configuration of the contributing data cluster is such as to ensure that the generalized Gram matrix $\mathbf{A}_{[0;p-1,0;p-1]}$ remains positive-definite, the leading-order

terms in the eigenvector expansion associated with eigenvalues that asymptotically scale as ϵ^p are spatial polynomials of degree, p , that do not depend upon the form of the assumed background covariance; the number of distinct eigenvectors at this order is then equal to the number of columns in \mathbf{v}_p . \square

Proof of theorem 3

In the trivial case, $p=0$, the single term, $\mathbf{U}_{p;p0} = 1$, fully describes the leading, and in this case, only eigenvector. For $p > 0$, then since $\mathbf{P}_{p;q0} = \mathbf{0}$, $q < p$, and $\mathbf{U}_{p;q0} = \mathbf{0}$, $q > p$, we therefore have a system of equations of the generic structure:

$$\begin{aligned} \mathbf{A}_{00}\mathbf{U}_{p;00} + \mathbf{A}_{01}\mathbf{U}_{p;10} + \cdots + \mathbf{A}_{0,p-1}\mathbf{U}_{p;p-1,0} &= -\mathbf{A}_{0p}, \\ \mathbf{A}_{10}\mathbf{U}_{p;00} + \mathbf{A}_{11}\mathbf{U}_{p;10} + \cdots + \mathbf{A}_{1,p-1}\mathbf{U}_{p;p-1,0} &= -\mathbf{A}_{1p}, \\ &\dots \\ \mathbf{A}_{p-1,0}\mathbf{U}_{p;00} + \mathbf{A}_{p-1,1}\mathbf{U}_{p;10} + \cdots + \mathbf{A}_{p-1,p-1}\mathbf{U}_{p;p-1,0} &= -\mathbf{A}_{p-1,p}, \end{aligned}$$

Given that $\mathbf{A}_{[0:p-1,0:p-1]}$ is positive-definite, we are able to invert this system to obtain all the $\mathbf{U}_{p;q0}$, for $q < p$ and, $\mathbf{U}_{p;p0}$ from (4.19), which collectively define the \mathbf{U}_p to the leading degree, ϵ^p , in the asymptotic expansion without any dependency upon the coefficients b_r of the particular covariance model assumed. \square

Having determined all the coefficients, $\mathbf{U}_{p;q0}$, for a given p , it follows that $\mathbf{P}_{p;q0}$ is also known for any desired q by applying (4.21) with $r=0$. The $[p/2]$ sets of quantities, $\mathbf{P}_{p;p-2s,s}$ with $0 < s \leq [p/2]$, are now found from the simultaneous solution of the equations obtained by setting to zero the coefficients in the expansion on left-hand sides of (4.20) associated with $\epsilon^{p-s}\mathbf{v}_{p-2s}$, $s = 1, \dots, [p/2]$, and noting that the right-hand sides vanish for this range of $r = p - s$. Again, it is the positive-definiteness of the matrices in (4.6) and (4.7) that ensures that this is always possible. In order to evaluate the leading term in the eigen-matrix, $\mathbf{\Lambda}_{p;p}$, we then have:

$$\mathbf{\Lambda}_{p;p} = \sum_{s=0}^{[p/2]} \mathbf{b}_{p,p-2s} \mathbf{P}_{p;p-2s,s}. \quad (4.22)$$

Higher order terms in the expansions for the segregated subspaces of eigenvectors spanned by the \mathbf{U}_p , and of the corresponding eigen-matrices, $\mathbf{\Lambda}_p$, follow by delving deeper into the double sequence of terms in $\epsilon^r \mathbf{v}_q$ that are equated in each member, $\mathcal{E}(p; q, r)$, of (4.20) and each member, $\mathcal{P}(p; q, r)$, of (4.21). An outline of the generic procedure is given in Appendix C.

For raw data errors that are uncorrelated, so that \mathbf{W} is diagonal, then in the case of the important single dominant eigenvalue, the leading asymptotic term gives us

$$\begin{aligned} \mathbf{\Lambda}_{0;0} &= \mathbf{b}_{00} \mathbf{A}_{00} \\ &= \mathbf{b}_{00} \mathbf{v}_0^\top \mathbf{W} \mathbf{v}_0 \\ &= b_0 \text{trace}(\mathbf{W}). \end{aligned} \quad (4.23)$$

Since b_0 is just the background error variance, this result is giving us the intuitively reasonable result that the super-ob's dominant component has a weight which is, to leading approximation in the asymptotic parameter, simply the sum of the weights of the contributing data. The relationship connecting the dominant eigenvector, \mathbf{U}_0 , to the composition weight vector in the

first column of \mathbf{L} is that each component of this leading column of \mathbf{L} is proportional to the corresponding diagonal component of \mathbf{W} times the component of \mathbf{U}_0 and, since to leading order in ϵ , we have

$$\mathbf{U}_0 = \mathbf{v}_0 + \mathcal{O}(\epsilon), \quad (4.24)$$

the composition weights for these independent point data are, to leading approximation, proportional to these data's weights. This is again in accordance with the intuitive practice of generating scalar super-obs as weighted averages of the contributing data.

When we look at the remaining less dominant components, \mathbf{U}_p , for $p > 0$, the leading terms in the asymptotic expansions give us exactly the multipole ‘Type-2’ super-obs defined in Purser et al. (2000). Thus, the efforts of this section have led us to an asymptotic justification of the general multipole super-obs of Purser et al., described here as approximations to the optimal super-obs obtained by asymptotic expansions in the spatial coordinates. While these leading-order approximations to the composition weights of the multipole super-obs do not depend upon the covariance model for the background (theorem 3), the eigenvalues which give the effective precision weights of each component of the super-ob *do* depend upon the background covariance model, even at leading order in the expansion.

This study also informs us that, since the eigenvalues at each spatial degree, p , scale as ϵ^p , we can expect that the super-ob weights of any tight cluster [implying small ϵ , when the components of the \mathbf{v}_q are maintained to have values collectively of $\mathcal{O}(1)$] will be more strongly dominated by $\mu_0^2 \equiv \mathbf{\Lambda}_0$ than is the case when the cluster is of a looser form, extending out to distances more comparable with the scale set by the background error covariance. Thus, if super-obs are used in a multi-grid context, in which the assimilation calculations are carried out (at least partially) at a wide range of resolutions, it is at the coarsest resolutions that the non-trivial multipole super-ob representations are expected to have the greatest value, being less strongly dominated by the single largest eigenvalue of \mathbf{Q} in these cases.

The asymptotic analysis presented in Appendix C has a potential vulnerability for any $p > 0$ if we attempt to extend the expansion in cases where the leading contribution, $\mathbf{\Lambda}_{p;p}$, to the eigenmatrix is singular or ill-conditioned. This will occur if the geometrical configuration of the cluster of point data are such that they do not form a stencil that allows all spatial derivatives up to and including degree, p , to be expressed by these points alone, in a well-conditioned way. Unfortunately, however, real-world data are often distributed in geometrically regular ways that can certainly provoke troubles of this kind. For example, aircraft data lie along a (roughly) one-dimensional track through space and an attempt to combine even a very numerous (but spatially compact) set of such measurements into a two- or three-dimensional ‘dipole’ super-ob by the procedure we have described will fail; the stencil that raw data provide only allow the expression of the first-derivative *along* the track, not across it. In this case, the 2×2 matrix, $\mathbf{\Lambda}_{1;1}$, has a rank of one and, being non-invertible, does not allow a continuation of the asymptotic expansion of Appendix C to be carried out to a higher degree. It is therefore necessary to examine the rank of the leading matrix term, $\mathbf{\Lambda}_{p;p}$, at each p to decide whether the super-ob, as formulated, is of this degenerate type. In the example of aircraft data, such data must be super-obbed, if at all, recognizing that they are effectively confined to one dimension (even when their track is a curved one).

5. DISCUSSION

To an increasing degree, the data going into a modern operational data assimilation system are dominated by massive amounts of automatically gathered measurements, whether remotely sensed from satellites or radar, or deriving from the more conventional platforms of radiosondes and aircraft. The necessity of data reduction cannot be doubted, but the manner in which this reduction is carried out can adversely affect the assimilation if the adopted procedure does not intelligently take account of both the quality (precision) and the density of each type of data.

The first thing to note is that measurements that are essentially continuous in some index parameter without their errors being strongly correlated, should have their quality expressed quantitatively as a “precision density” in that index variable. The super-obbing theory validates this results since, as we have found from its first-order asymptotics, the effective weight of an aggregation of uncorrelated data bunched well within the characteristic scale of the background covariance (assumed smooth) is essentially additive. To give some examples, the quality of (Kelvin, K) temperature data measuring in a radiosonde profile where height is the index variable should be given as the precision density, K^{-2} per unit height; aircraft temperature data provided along a track with sufficient frequency to be essentially a continuous line as far as they affect a horizontal analysis should be assessed a precision density in units of $\text{K}^{-2}\text{km}^{-1}$; hyperspectral measurements of brightness temperatures, where the suitable index might be wavenumber, cm^{-1} , should be characterized by a precision per unit wavenumber, units: K^{-2}cm , and so on. Whether such data are super-obbed, or simply thinned, the proper specification of their original precision densities, together with the index densities of the resulting super-obs or thinned representatives, provide objective guidance for choosing the correct effective precisions of these artificial representatives of the underlying data-continuum.

Given the nature of the Krylov-based solution algorithms almost universally employed to solve variational assimilation problems, there is one surprising and paradoxically deleterious impact that very good data can have. When their precisions are specified to be vastly greater than the corresponding components of precision of the background field, and they are sufficiently numerous, they can stall the convergence of these algorithms owing to the very large condition numbers they engender. The method of super-obbing addresses both the problem of overwhelming amounts of data, and, at least in part, the issue of adverse effects on the condition number, provided the super-obs are constructed in a way that takes account of the information that might be safely under-weighted, without adverse effect, from the original (i.e., pre-super-obbed) data relative to the information already contained in the pre-existing background field.

This note is an attempt to formalize objective principles by which super-obs, from a range of sources, might be constructed so that valuable information, in an objectively measured sense, is retained to the greatest degree allowed by the inherent constraints imposed by projecting the information from each cluster or region of data onto a small and discrete number of its representative degrees of freedom. It is shown that, in a reasonably convincing asymptotic limit of smooth fields, the multipole construction described in Purser et al. (2000) is actually optimal in terms of minimizing information loss, and is not sensitive to the precise criterion amongst the set of “trace” measures of information adopted. The suggested empirically-motivated capping of eigenvalues associated with the relative precision of the super-obs is not so amenable to rigorous justification, but makes sense in the light of the known behavior of conjugate gradients-type

algorithms that struggle to converge when sufficiently many of their non-dimensionalized system matrix's largest eigenvalues greatly exceed unity.

An additional future advantage of intelligent super-obbing might be gained if observation-space preconditioners are adopted (as in the Navy's NAVDAS described by Daley and Barker, 2000), since in that case, the computational cost of accommodating the effects between distant tight clusters of data, treated as super-obs with only few degrees of freedom, will be considerably less than the corresponding cost of treating all the combinations of interactions between the very much more numerous pairs of the same data in their pre-super-obbed existence. This is a topic worthy of future investigation.

APPENDIX A

Defining logarithms of matrices

With

$$\mathbf{M} = \mathbf{DHC} \tag{A.1}$$

as in (2.10), we can perform a singular value decomposition:

$$\mathbf{M} = \mathbf{V}\boldsymbol{\mu}\mathbf{Z}^\top \tag{A.2}$$

where the columns of \mathbf{V} and \mathbf{Z} are orthogonal, and $\boldsymbol{\mu}$ is a diagonal matrix of positive elements. The logarithm of the symmetric and positive-definite combination, $\mathbf{I} + \mathbf{M}^\top\mathbf{M}$, is always defined and is:

$$\log(\mathbf{I} + \mathbf{M}^\top\mathbf{M}) = \mathbf{Z} \log(\mathbf{I} + \boldsymbol{\mu}^2)\mathbf{Z}^\top \tag{A.3}$$

and in the same way:

$$\log(\mathbf{I} + \mathbf{M}\mathbf{M}^\top) = \mathbf{V} \log(\mathbf{I} + \boldsymbol{\mu}^2)\mathbf{V}^\top. \tag{A.4}$$

It is straightforward to verify that,

$$\mathbf{H}^\top\mathbf{W}\mathbf{H}\mathbf{B} = \mathbf{B}^{-1}\mathbf{C}\mathbf{Z}\boldsymbol{\mu}^2\mathbf{Z}^\top\mathbf{C}^\top, \tag{A.5}$$

and so we also consistently define,

$$\log(\mathbf{I} + \mathbf{H}^\top\mathbf{W}\mathbf{H}\mathbf{B}) = \mathbf{B}^{-1}\mathbf{C}\mathbf{Z} \log(\mathbf{I} + \boldsymbol{\mu}^2)\mathbf{Z}^\top\mathbf{C}^\top. \tag{A.6}$$

The trace of a product of matrices is unchanged by any cyclic permutation of its factors, and hence:

$$\text{trace} \log(\mathbf{I} + \mathbf{M}^\top\mathbf{M}) = \text{trace} \log(\mathbf{I} + \mathbf{M}\mathbf{M}^\top) = \text{trace} \log(\mathbf{I} + \mathbf{H}^\top\mathbf{W}\mathbf{H}\mathbf{B}) = \text{trace} \log(\mathbf{I} + \boldsymbol{\mu}^2), \tag{A.7}$$

that is, the sum of the logarithms of the diagonal elements, $1 + \mu^2$, or the log-determinant of any of the matrices of which these diagonal elements are the non-unit eigenvalues.

APPENDIX B

Polynomial coefficient matrices, \mathbf{a}

In any number of dimensions

$$\mathbf{a}_{00} = [1], \quad (\text{B.1})$$

while, in two dimensions, second-degree terms are combined using the three coefficient matrices:

$$\mathbf{a}_{02} = \begin{bmatrix} 1, & 0, & 1 \end{bmatrix}, \quad (\text{B.2})$$

$$\mathbf{a}_{11} = \begin{bmatrix} -2, & 0 \\ 0, & -2 \end{bmatrix}, \quad (\text{B.3})$$

$$\mathbf{a}_{20} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}. \quad (\text{B.4})$$

The coefficient matrices of higher degrees can be obtained by the general recursion that uses the notion of the convolution product of matrices:

$$\mathbf{a}_{mn} = \sum_{i,j,k,l} (\mathbf{a}_{ij} \star \mathbf{a}_{kl}) \delta_{i+k,m} \delta_{j+l,n} \delta_{k+l,2}, \quad (\text{B.5})$$

where δ_{ij} denotes the Kronecker delta, and the convolution operator, \star , among the \mathbf{a} matrices is defined in the following way. First, for these matrices, we take the indices \mathbf{i}, \mathbf{j} , etc of their rows and columns to be *integer vectors* having as their own component of the corresponding powers of X, Y , and so on, referred to by these row and column indices. Let $\mathbf{a}_{\mathbf{i}\mathbf{j}}$ be the element of \mathbf{a} at such a vector-index row and column pair, \mathbf{i} and \mathbf{j} . Then the matrix convolution is defined, again using the Kronecker delta (but now generalized to vector indices):

$$\mathbf{a}''_{\mathbf{i}'\mathbf{j}''} = \sum_{\mathbf{i}, \mathbf{j}, \mathbf{i}', \mathbf{j}'} (\mathbf{a}'_{\mathbf{i}'\mathbf{j}'} \mathbf{a}_{\mathbf{i}\mathbf{j}}) \delta_{\mathbf{i}+\mathbf{i}', \mathbf{i}''} \delta_{\mathbf{j}+\mathbf{j}', \mathbf{j}''}. \quad (\text{B.6})$$

For example, applying (A.5) get the matrices \mathbf{a} of fourth-degree,

$$\mathbf{a}_{04} = \mathbf{a}_{02} \star \mathbf{a}_{02}, \quad (\text{B.7a})$$

$$\mathbf{a}_{13} = \mathbf{a}_{11} \star \mathbf{a}_{02} + \mathbf{a}_{02} \star \mathbf{a}_{11}, \quad (\text{B.7b})$$

$$\mathbf{a}_{22} = \mathbf{a}_{20} \star \mathbf{a}_{02} + \mathbf{a}_{11} \star \mathbf{a}_{11} + \mathbf{a}_{02} \star \mathbf{a}_{20}, \quad (\text{B.7c})$$

$$\mathbf{a}_{31} = \mathbf{a}_{20} \star \mathbf{a}_{11} + \mathbf{a}_{11} \star \mathbf{a}_{20}, \quad (\text{B.7d})$$

$$\mathbf{a}_{40} = \mathbf{a}_{20} \star \mathbf{a}_{20}, \quad (\text{B.7e})$$

and, in two dimensions, the definition of the matrix-convolution operator leads to the result:

$$\mathbf{a}_{04} = \mathbf{a}_{40}^{\text{T}} = \begin{bmatrix} 1, & 0, & 2, & 0, & 1 \end{bmatrix}, \quad (\text{B.8})$$

$$\mathbf{a}_{13} = \mathbf{a}_{31}^{\text{T}} = \begin{bmatrix} -4, & 0, & -4, & 0 \\ 0, & -4, & 0, & -4 \end{bmatrix}, \quad (\text{B.9})$$

$$\mathbf{a}_{22} = \begin{bmatrix} 6, & 0, & 2 \\ 0, & 8, & 0 \\ 2, & 0, & 6 \end{bmatrix}. \quad (\text{B.10})$$

APPENDIX C

Continuing the asymptotic expansions for the non-leading terms of \mathbf{U}_p and $\mathbf{\Lambda}_p$

It is already determined that $\mathbf{P}_{p;qr} = \mathbf{0}$ for all q and r such that $q + 2r < p$. After estimating the leading terms in \mathbf{U}_p and $\mathbf{\Lambda}_p$ the (possibly non-vanishing) values of $P_{p;qr}$ for all $q + 2r \leq p$ are known, as are those for any desired q for $r = 0$. The matrix quantities $\mathbf{U}_{p;q0}$ are also known for all q . We shall proceed by iteratively repeating a cycle of steps, each cycle indexed by $t \geq 1$, refining the estimates of \mathbf{U}_p and $\mathbf{\Lambda}$ by adding the next-degree of approximation to these quantities in the course of each cycle executed. We justify the progression of steps by an inductive argument that assumes that, at the beginning of the cycle, t , we know $\mathbf{P}_{p;qr}$ for any $q + 2r \leq p + 2t - 2$ and, for $r < t$ for any $q + 2r \leq p + 2t$. We also assume that all the $\mathbf{U}_{p;qr}$ are known for $r < t$. These facts are certainly true at the beginning of cycle, $t = 1$, so it is sufficient to show that, if they are true at the start of an arbitrary cycle, t , they necessarily remain true at the start of cycle, $t + 1$. However, before we extend the asymptotic expansion to higher degrees in the implicit parameter ϵ we need to ascertain that the initial approximation, $\mathbf{\Lambda}_{p;p}$, to each eigenmatrix is nondegenerate. If it possesses a null space, or even if it is ill-conditioned to inversion, we are forced to conclude that p th-degree super-ob components constructed from this cluster of data are collectively degenerate and there is not a unique way to extend the asymptotic series beyond the leading terms. In what follows, we therefore assume that the geometrical configuration does *not* suffer this defect.

The steps of cycle, t , are as follows.

Step 1:

Use $\{\mathcal{E}(p; p + 1 - 2s, p + t - s), \quad 1 \leq s \leq [(p + 1)/2]\}$

to solve simultaneously for $\{\mathbf{P}_{p;p+1-2s,t+s-1}, \quad 1 \leq s \leq [(p + 1)/2]\}$.

Step 2:

Use $\mathcal{E}(p; s, p + t)$ to solve for $\mathbf{U}_{p;s,t}$, for $p + 1 \leq s \leq p + 2t$.

Step 3:

Use $\{\mathcal{P}(p; s, t), \quad 0 \leq s < p\}$ to solve simultaneously for $\mathbf{U}_{p;s,t}$, for $0 \leq s < p$.

Step4:

Use $\mathcal{P}(p; p, t)$ to solve for $\mathbf{P}_{p;p,t}$, and use $\{\mathcal{P}(p; p + 1 + s, t - [s/2]), \quad 0 \leq s \leq 2t + 2\}$

to solve for $\mathbf{P}_{p;p+1+s,t-[s/2]}, \quad 0 \leq s \leq 2t + 2$.

Step 5:

Use $\{\mathcal{P}(p; p - s - 1, p + t - s), \quad 1 \leq s \leq [p/2]\}$ to solve simultaneously for

$\{\mathbf{P}_{p;p-s-1,t+s}, \quad 1 \leq s \leq [p/2]\}$.

Step 6:

Use $\mathcal{E}(p; p, p + t)$ to solve for $\mathbf{\Lambda}_{p;p+t}$.

Increment t by 1, go back to Step 1 and repeat as necessary, or:

Step 7:

End. \square

Obviously, some of these iterative steps are not exercised when $p = 0$. An examination of the simultaneous linear equations we need to solve at Step 1 and at Step 5 reveals that, again, we are assuming the non-singularity of the block-matrices of (4.6) and (4.7), while an

examination of the simultaneous linear equation we need to solve at Step 3 reveals that we are assuming the same condition of invertibility of the generalized Gram matrix, $\mathbf{A}_{[0:p-1,0:p-1]}$, as was already assumed in the asymptotic derivation of the leading coefficients, $\mathbf{U}_{p;q0}$, $q < p$, in section 4. It is at Step 2 that we require $\mathbf{\Lambda}_{p;p}$ to be invertible and well-conditioned. Given all these conditions we find that our cycle of derivations justify the continuation of the inductive process. Therefore, if the conditions are met that allow us to initially solve for the leading asymptotic terms of \mathbf{U}_p and non-degenerate $\mathbf{\Lambda}_p$, then we can certainly continue the cycle of derivations for the asymptotic terms up to any higher degree we choose.

REFERENCES

- Courant, R. and D. Hilbert 1989 *Methods of Mathematical Physics, Volume I*. Wiley, New York.
- Courtier, P. 1997 Dual formulation of four-dimensional variational assimilation. *Quart. J. Roy. Meteor. Soc.*, **123**, 2449–2461.
- Cover, T. M., and J. A. Thomas 1991 *Elements of Information Theory*. Wiley, New York, 576 pp.
- Daley, R., and E. Barker 2000 *NAVDAS Source Book 2000*, Naval Research Laboratory, Atmospheric Dynamics and Prediction Branch, Marine Meteorology Division, Monterey, CA 93943-5502. 153 pp.
- Eyre, J. R. 1990 The information content of data from satellite sounding systems: a simulation study. *Quart. J. Roy. Meteor. Soc.*, **116**, 401–434.
- Kiefer, J. 1974 General equivalence theory for optimum designs (approximate theory). *Ann. Statist.*, **2**, 849–879.
- Lorenc, A. C. 1981 A global three dimensional multivariate statistical interpolation scheme. *Mon. Wea. Rev.*, **109**, 701–721.
- Pázman, A. 1986 *Foundations of Optimum Experimental Design*. Reidel, 228 pp.
- Peckham, G. 1974 The information content of remote measurements of the atmospheric temperature by satellite IR radiometry and optimum radiometer configurations. *Quart. J. Roy. Meteor. Soc.*, **100**, 406–419.
- Purser, R. J. 1998 Structural preconditioning in optimal analysis. NOAA/NCEP Office Note 422. 20 pp.
- Purser, R. J., and H.-L. Huang 1993 Estimating the effective data density in a satellite retrieval or an objective analysis. *J. Appl. Meteor.*, **32**, 1092–1107.
- Purser, R. J., D. F. Parrish, and M. Masutani 2000 Meteorological observational data compression; an alternative to conventional “super-obbing”. NOAA/NCEP Office Note 430. 12pp.
- Rodgers, C. D. 1998 Information content and optimisation of high spectral resolution remote measurements. *Adv. Space Res.*, **21**, 361–367.
- Shannon, C. E. 1948 The mathematical theory of communication. *Bell Syst. Technol. J.*, **27**, 379–423, 623–656.
- Shannon, C. E., and W. Weaver 1949 *The Mathematical Theory of Communication*. University of Illinois Press, Chicago.
- van der Vorst, H. A. 2003 *Iterative Krylov Methods for Large Linear Systems*. Cambridge University Press.
- Wahba, G. 1985 Design criteria and eigensequence plots for satellite-computed tomography. *J. Atmos. Oceanic Technol.*, **2**, 125–132.