# On Criteria for Measures of Compositional Difference[1]

## John Aitchison[2]

*Simple perceptions about the nature of compositions lead through logical necessity to certain forms of analysis of compositional data. In this paper the consequences of essential requirements of scale, perturbation and permutation invariance, together with that of subcompositional dominance, are applied to the problem of characterizing change and measures of difference between two compositions. It will be shown that one strongly advocated scalar measure of difference fails these tests of logical necessity, and that one particular form of scalar measure of difference (the sum of the squares of all possible logratio differences in the components of the two compositions), although not unique, emerges as the simplest and most tractable satisfying the criteria.*

## INTRODUCTION

Simple perceptions of the nature of compositions lead through logical necessity to certain specific forms of compositional data analysis. Unfortunately recent publications in *Mathematical Geology* (Watson and Philip (1989) and Watson (1990, 1991) in reply to Aitchison (1990, 1991)), have tended to obscure the essential simplicity that can be achieved in the study of patterns of variability of compositional data sets. The aim of this paper is to correct this trend by spelling out, in clear and precise mathematical language, a minimum set of criteria for one particular object of interest, namely some sensible scalar measure of the difference between two compositions, and by revealing how these criteria limit the apparent range of possible scalar measures. In the course of this journey the angular measure of Watson and Philip (1989), will be found to fail these criteria in two fundamental ways. Among measures which satisfy the criteria one scalar measure of compositional difference, namely the sum of the squares

of all possible logratio differences in the components, will be seen to emerge as the simplest and most tractable.

It should be emphasized that the perceptions which lead to these forms are general, not particular to geological applications. Thus any geological perceptions would be additional to those considered here. The fact that some recent publications have already digressed from what are logical necessities for general compositional data analysis precludes their consideration from any particularity imposed by additional, geological considerations.

In what follows an attempt has been made to set out the logical steps in as simple mathematical terms as possible. Since the aim of the paper is, however, to bring mathematical precision and deduction to the topic, there is no real alternative to the use of fundamental mathematical concepts such as group operations and invariance under groups of transformations.

## THE NATURE OF COMPOSITIONS

A $D$-part composition is a $D$-dimensional positive vector $x = (x_1, \ldots, x_D)$, where each component $x_i (i = 1, \ldots, D)$ is measured in the same units. Consider the following three vectors:

| Units of measurement | A | B | C | D |
|---|---|---|---|---|
| gm/specimen | 20.80 | 16.12 | 11.44 | 3.64 |
| oz/specimen | 0.733 | 0.569 | 0.404 | 0.128 |
| proportion by weight | 0.40 | 0.31 | 0.22 | 0.07 |

The parts A, B, C, D are four minerals and the three vectors represent the mineral composition measured in gm per specimen, oz per specimen, and as proportion by weight of the same rock specimen. It is widely recognized that the particular unit chosen should make no difference in any statistical analysis of such specimens. The components of each of these vectors are equally scaled versions of the components of one of the other vectors. It is convenient to formalize this simple perception about compositions for further use.

*Equivalent Compositions.*   Two compositions $x$ and $X$ are to be regarded as equivalent, written $x \sim X$, if there is some $a > 0$, a scale factor, such that $X = ax$. For example, if $x$ and $X$ are the first two vectors above then $a = 0.03527$, the ratio of the ounce to the gram. The equivalence relationship is, of course, a symmetrical one since $x = a^{-1}X$.

*Composition Class.*   For any composition $x$ the set of compositions equivalent to $x$ can be collected into a composition class $c = \{ax: a > 0\}$, so that the whole of $R_+^D$, the positive orthant of D-dimensional real space, is partitioned by the equivalence operation into disjoint composition classes.

*Standardized Compositions.*   Within each composition class $c$ there is a

unique standardized composition $u$ whose components have a unit sum, $u_1 + \ldots + u_D = 1$, determined by

$$u = x/(x_1 + \ldots + x_D)$$

where $x$ is any of the equivalent compositions in $c$. The third vector above is the standardized form of the other two vectors. It is sometimes convenient to use the standardized composition in the definition of its composition class, for example, as

$$c = \{au: a > 0\}.$$

The geometric picture is presented in Fig. 1 for 3-part compositions. The composition $x$ determines as its composition class $c$ the ray through the origin and $x$ and the standardized composition $u$ is the intersection of the ray with the plane $x_1 + x_2 + x_3 = 1$.

## SCALE INVARIANT FUNCTIONS

A first perception in the study of compositions is that it should not matter which specific composition $x$ is selected from $c$ as representing the physical entity. Thus, in compositional data analysis any construct or function $f$ of a composition must satisfy the requirement of scale invariance:

$$f(ax) = f(x) \text{ for every } x \in c \text{ and for every } a > 0 \qquad (1)$$
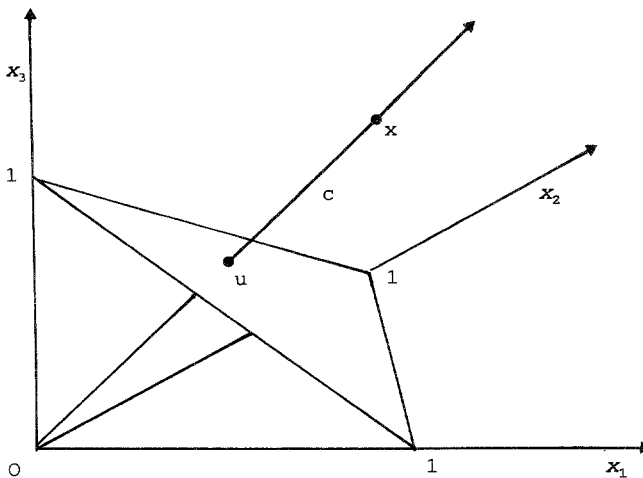


**Fig. 1.** Representation of a composition $x$, its composition class $c$, and its standardized composition $u$.

An immediate and important characterization of a scale invariant function is obtained by substituting the special value $a = 1/x_D$ in Eq. (1), showing that

$$f(x_1, \ldots, x_D) = f(x_1/x_D, \ldots, x_d/x_D, 1) \tag{2}$$

where $d = D - 1$. Thus any scale invariant function of a composition is expressible as a function of a set of $d = D - 1$ ratios, for example

$$r_i = x_i/x_D \qquad (i = 1, \ldots, d) \tag{3}$$

of any of the equivalent compositions in $c$. It follows that any meaningful statement about a composition must be expressible in terms of ratios of components. Another way of seeing this is that the inverse of the transformation Eq. (3) in standardized compositional terms is

$$u_i = \frac{r_i}{r_1 + \ldots + r_d + 1} \qquad (i = 1, \ldots, d) \tag{4a}$$

$$u_D = \frac{1}{r_1 + \ldots + r_d + 1} \tag{4b}$$

Since any statement made about a composition can be expressed in terms of its standardized composition $(u_1, \ldots, u_D)$ it can be converted, simply by replacing each $u_i$ by the appropriate expression in ratios from Eq. (4). For example, the component statements that

$$u_1 < u_2, \quad u_1 < u_D, \quad u_1 < 0.4$$

translate into the corresponding ratio statements

$$r_1 < r_2, \quad r_1 < 1, \quad 3r_1 - 2(r_2 + \ldots + r_d) < 2$$

This consequence of scale invariance, which is often simply expressed by the statement that compositions provide only information about the relative magnitudes of the components, seems so obvious that it is surprising to find misunderstanding of it within the pages of *Mathematical Geology*. For example, Watson (1990) makes the following unsupported assertion:

> to suggest that "Any discussion of the variability of a composition . . . can be expressed in terms of ratios . . . of components" is in any way axiomatic or self-evident (Aitchison, 1990), exhibits a failure to investigate even elementary examples.

and yet elsewhere (Watson and Philip, 1989) the same author, recognizing the importance of scale invariance, acknowledges the ratio characterization (italics denote insertions to relate to the terminology of this paper):

> The magnitude of any component does not have any significance in itself, but only in its proportion relative to other components. For example, if a thin section were inadvertently counted twice, the magnitudes of the components would be doubled, but the compositional axis (*the compositional class* c) remains the same (within expected count-

ing error). . . . Any set of $(D)$ numbers that define the same compositional axis can be used to describe that composition.

There are many scale invariant functions of a composition, such as

$$x_i/x_D; \quad x_i/(x_1^2 + \ldots + x_D^2)^{1/2}; \quad \log\{x_i/(x_1 \ldots x_D)^{1/D}\}$$
$$x_i/(x_1 + \ldots + x_D)$$

Note also that the function defining any component $u_i$ of a standardized composition in terms of any composition $x$ of its compositional class is scale invariant. Some examples of functions which are not scale invariant are the following:

$$x_i; \quad x_1 + x_2; \quad \log x_i$$

## MATHEMATICAL NOTE

The structure of the above problem is common in many areas of investigation in statistical theory where it is important to ensure that functions and procedures are *invariant* under some relevant group of transformations of the data. The relevant group here consists of the scale transformations of compositions and the composition classes are the orbits of the group. A main aim in such investigations is to identify a *maximal invariant*, defined as an invariant function $f$ which is such that $f(x) = f(X)$ implies that $x \sim X$ or equivalently that there is some $a > 0$ such that $X = ax$. The central role of a maximal invariant is that every invariant function can be expressed as a function of a maximal invariant (Lehmann, 1959, p. 216). In the above problem, it is trivial to show that the vector function

$$f(x) = (x_1/x_D, \ldots, x_d/x_D)$$

is a maximal invariant, corroborating the conclusion above that scale invariance requires that statements about compositions must be expressible in terms of ratios of the components. The set of ratios considered is not confined to the set at Eq. (3) above. Any set, such as $x_1/x_2, x_2/x_3, \ldots, x_d/x_D$, or $x_i/(x_1 \ldots x_D)^{1/D}$ ($i = 1, \ldots, D$), also maximal invariants, which determine a composition class, will serve equally well.

At this stage, it is of interest to note that there are many possible functions $f(x, X)$ of two compositions $x$ and $X$, which are scale invariant, with $f(ax, AX) = f(x, X)$ for every $a > 0$, $A > 0$, and which might serve as a scalar measure of difference between $x$ and $X$. Some examples, with $\Sigma_i$ and $\Sigma_{i<j}$ denoting summation over the integers $i = 1, \ldots, D$ and $i = 1, \ldots d; j = i + 1, \ldots,$ $D$, respectively, are the following.

$$\arccos\left\{ \sum_i x_i X_i \middle/ \left( \sum_i x_i^2 \sum_i X_i^2 \right)^{1/2} \right\} \tag{5}$$

$$\sum_{i<j} (x_i X_j - x_j X_i)^2 \middle/ \left( \sum_i x_i^2 \sum_i X_i^2 \right) \tag{6}$$

$$\sum_{i<j} \{ (x_i/x_j)/(X_i/X_j) - (x_j/x_i)/(X_j/X_i) \}^2 \tag{7}$$

$$\sum_{i<j} |\log \{ (x_i/x_j)/(X_i/X_j) \}| \tag{8}$$

$$\sum_{i<j} [\log \{ (x_i/x_j)/(X_i/X_j) \}]^2 \tag{9}$$

Among the above functions, Eq. (5) is the angular measure of Watson and Philip (1989), the angle between the two composition classes regarded as directions, and Eq. (9) is a scalar measure equivalent to that of Aitchison (1986a, p. 193, 1986b).

It is clear that scale invariance alone leaves a large class of possible scalar measures of difference. There are, however, further criteria which are necessary and it will be seen that two other natural groups of transformations play a central role in what follows in describing important invariance requirements of compositional data analysis, particularly in relation to the search for sensible measures of difference between compositions.

## PERTURBATION AS THE MEASURE OF CHANGE BETWEEN TWO COMPOSITIONS

In Euclidean space, two vectors $z$ and $Z$ can always be fully related by asking what transformation is required to change $z$ into $Z$. The answer is in the operation of a translation $t$ where $Z = t + z$, or equivalently by the inverse translation $z = -t + Z$. Moreover this relationship between $z$ and $Z$ is the same as that between $z^*$ and $Z^*$ if and only if $z^*$ and $Z^*$ are equal translations of $z$ and $Z$. Any definition of a difference or a distance measure must thus be such that the measure is the same for $t + z$, $t + Z$ as for $z$, $Z$ for every translation $t$. Technically, this is a requirement of invariance under the group of translations. In the consideration of the differences between compositions or composition classes the obvious first questions are whether there is an operation on a composition $x$, analogous to translation of a vector in $R^D$, which transforms it into $X$, and whether this can be used to characterize the relationship or difference between two compositions. The answers are to be found in the perturbation operator defined by Aitchison (1986a, Section 2.8) and already used in a geological application by Woronow and Love (1990). The argument is only

slightly more complicated than that for Euclidean space because of the equivalence relationships between compositions.

The perturbation operator can be motivated by the following observation. For any two equivalent compositions $x$ and $X$, in the same compositional class, there is a scale relationship

$$(X_1, \ldots, X_D) = (ax_1, \ldots, ax_D) \tag{10}$$

for some $a > 0$, where each component of $x$ is scaled by the same factor to obtain the corresponding component of $X$. For any two compositions $x$ and $X$ in different compositional classes $c$ and $C$ a similar, but differential, scaling relationship

$$(X_1, \ldots, X_D) = (q_1 x_1, \ldots, q_D x_D) \tag{11}$$

can always be found, simply by taking $q_i = X_i/x_i$ $(i = 1, \ldots, D)$. If the operation between the positive perturbing vector $q = (q_1, \ldots, q_D)$ and the composition $x$ is denoted by $\circ$ so that

$$q \circ x = (q_1 x_1, \ldots, q_D x_D) \tag{12}$$

then the above relationship is denoted by $X = q \circ x$.

The form of this relationship suggests a method by which the difference between two composition classes can be fully characterized. In Eq. (11) the perturbing vector $q$ is of the same mathematical form as a composition and it is convenient to define its perturbation class $p = \{bq: b > 0\}$ in the same way as a composition class can be defined in relation to any one of its compositions. The set of all perturbations $q \in R_+^D$ form a group under the $\circ$ operation: in particular $q$ has an inverse $q^{-1} = (1/q_1, \ldots, 1/q_D)$, the identity perturbation is $e = (1, \ldots, 1)$ and the operation is commutative with $q_1 \circ q_2 = q_2 \circ q_1$ for any two perturbations $q_1$ and $q_2$. Some simple obvious rules of the operator $\circ$, such as $(Q \circ q^{-1}) \circ q = Q \circ (q^{-1} \circ q) = Q \circ e = Q$ are required in what follows. A consequence of the similarity between perturbations and compositions is that a perturbation can in fact be defined in terms of compositions; for example $U \circ u^{-1}$ is the perturbation $(U_1/u_1, \ldots, U_D/u_D)$.

The relation between any two compositions $x \in c = \{au: a > 0\}$ and $X \in C = \{AU: A > 0\}$ can always be expressed as a perturbation operation $X = q \circ x$, where $q$ is a perturbation in the perturbation class $p = \{bU \circ u^{-1}: b > 0\}$. This is very easily proved since for some $a > 0, A > 0, x = au, X = AU$ so that $X = AU = AU \circ u^{-1} \circ u = (A/a) U \circ u^{-1} \circ x$ and $(A/a) U \circ u^{-1}$ belongs to the perturbation class $p$. Thus the perturbation class defined by $U \circ u^{-1}$ or equivalently $X \circ x^{-1}$, where $x, X$ are any compositions in $c, C$, characterizes the change from $c$ to $C$. The change from $X$ to $x$ is simply the inverse perturbation class defined by $u \circ U^{-1}$.

An outstanding question is whether this role of the group of perturbations

in characterizing change in compositional classes is unique. Is there some other group $T$ of transformations $t: R^D_+ \rightarrow R^D_+$, which could equally serve the purpose of describing such changes. Such a group would require to satisfy the following conditions. The first simply describes in mathematical terms the role of a transformation.

*T1.*   For every $x, X \in R^D_+$ there is a unique $t \in T$ such that $X = t(x)$ and $x = t^{-1}(X)$.

The second requires that if $X = t(x)$ and $x$ is scaled to $ax$ in the same composition class then a corresponding scaling from $X$ to $aX$ is required in the transformed composition class, expressed formally as follows.

*T2.*   For every $x \in R^D_+$ and every $a > 0$, $t(ax) = at(x)$.

In mathematical terminology, T2 limits the transformations to homogeneous functions of $x$ of degree one. Perturbations are such homogeneous functions, and it requires only one more reasonable condition to reduce $T$ to the group of perturbations. Indeed there is a choice of two conditions, which are set out as T3a and T3b below.

Suppose that $X = t(x)$ records the change from $x$ to $X$. Suppose further that a scientist chooses to record the parts of a composition in different units; for example, in terms of the mineral composition at the beginning of this paper, the first, . . . , fourth parts might be recorded in g, oz, kg, tons per specimen. This would amount to a differential scaling, say $(q_1 x_1, \ldots, q_D x_D)$ or a perturbation $q \circ x$ of the original composition. However eccentric the scientist's behavior may be, it is surely reasonable to require that if a transformation $t$ is applied to the differentially scaled composition $q \circ x$ then for consistency the resulting vector $t(q \circ x)$ must be the appropriately scaled version of $X$, namely $q \circ X = q \circ t(x)$. This leads to the following condition.

*T3a.*   For every $x \in R^D_+$ and every perturbation $q$, $t(q \circ x) = q \circ t(x)$.

This is a generalization of T2 since perturbations of the special form $q = ae$, where $e$ is the identity perturbation, are scale transformations.

Condition T3a immediately determines the nature of $t$ by the following argument. Since perturbations and compositions are equivalent mathematically, T3a must hold when $q$ and $x$ are interchanged, giving $t(q \circ x) = t(q) \circ x$ with use of the commutative property of perturbations. On setting $q = e$, the identity perturbation, this gives

$$t(x) = t(e) \circ x$$

Since, for given $t$, $t(e) \in R^D_+$ and does not depend on $x$ the transformation $t$ is clearly a perturbation.

An alternative condition is to require that each component of a transformation $t$ should be a separate function of the corresponding component of $x$.

*T3b.*   The $i$th component $t_i$ of the vector function $t$ is a function only of $x_i$, in other words, $t(x) = \{t_1(x_1), \ldots, t_D(x_D)\}$.

A simple mathematical argument then leads again to perturbations. Differentiating the relationship $t_i(ax_i) = at_i(x_i)$ with respect to $a$ and then setting $a = 1$ gives a differential equation $t_i'(x_i)/t_i(x_i) = 1/x_i$ for $t_i$, for which the general solution is $t_i(x_i) = q_ix_i$, where $q_i$ is a constant. Thus, for $t$ to satisfy T3b, it must be a perturbation.

Thus acceptance of either T3a or T3b as a reasonable condition means that any measure of difference between composition classes $c$ and $C$ must be expressible in terms of one or other or both of the perturbation classes determined by $U \circ u^{-1}$ and $u \circ U^{-1}$.

## PERTURBATION INVARIANT FUNCTIONS

If two compositions $x$ and $X$ are operated on by the same perturbation $q$ to become $x^* = q \circ x$ and $X^* = q \circ X$ then the change from $x$ to $X$ is the same as from $x^*$ to $X^*$, since $X^* = (X \circ x^{-1}) \circ x^*$. Any measure of difference between composition classes must therefore be invariant under the operation of applying the same perturbation to each class. Thus any function purporting to measure difference must satisfy the criterion that, for every $x \in c$ and $X \in C$,

$$f(q \circ x, q \circ X) = f(x, X) \quad \text{for every perturbation } q \tag{13}$$

and is naturally termed a perturbation invariant function. As described above the important next step is to identify a maximal invariant: $f$ is a maximal invariant if, for every $x \in c$ and $X \in C$, $f(x^*, X^*) = f(x, X)$ implies that there exists a perturbation $q$ such that $x^* = q \circ x$ and $X^* = q \circ X$. It is easy to see that $f(x, X) = x \circ X^{-1}$ is a maximal invariant with $q = X^* \circ X^{-1}$ acting as the required perturbation. Thus any perturbation invariant function, and hence any measure of difference between composition classes, must be a function of the ratios $(x_1/X_1, \ldots, x_D/X_D)$. Since any such function is, *a fortiori*, a function of the reciprocals $X \circ x^{-1}$ there is no need to take account of the apparent asymmetry in taking a particular direction of change.

## PERMUTATION INVARIANT FUNCTIONS

There is a further simple invariance property that is required of any scalar measure of difference between two composition classes. If two scientists arrange the parts of two compositions in different orders and apply the formula for a scalar measure of difference to their arrangements, then they should arrive at the same value for the measure. Mathematically the operation here involves the group of permutations, and requires that $f(Px, PX) = f(x, X)$ for every $x \in c$, $X \in C$ and for every permutation $P$. The rather obvious though important consequence of this permutation invariant requirement is that $f(x, X)$ must be symmetric in the components of both $x$ and $X$.

A similar criterion arises with respect to the two compositions themselves, namely that the measure must not depend on the order in which $x$ and $X$ are considered, so that $f(x, X) = f(X, x)$. The simple consequence of this is the obvious one that $x$ and $X$ must be interchangeable in the specification of $f$.

## FEASIBLE SCALAR MEASURES OF DIFFERENCE

The various criteria considered above can now be brought together and their joint consequences applied to identify which scalar measures of difference may be regarded as satisfactory and which fail one or more of the criteria. A scalar measure of difference $f(x, X)$ between two compositions $x$ and $X$ is feasible if, for every $x$ and $X$, it satisfies the following criteria:

C1. Positivity. $f(x, X) > 0$ if $x$ and $X$ are not equivalent.

C2. Zero difference between equivalent compositions. $f(x, X) = 0$ if $x \sim X$.

C3. Interchangability of compositions. $f(x, X) = f(X, x)$.

C4. Scale invariance. $f(ax, AX) = f(x, X)$ for every $a > 0$, $A > 0$.

C5. Perturbation invariance. $f(q \circ x, q \circ X) = f(x, X)$ for every perturbation $q$.

C6. Permutation invariance. $f(Px, PX) = f(x, X)$ for every permutation $P$.

Although further obvious criteria will be introduced in the next section it is worth investigating how these five criteria delimit the range of functions to feasible scalar functions $f$.

C4 requires that $f$ must be expressible in terms of intracompositional ratios of the form $x_i/x_j$ and $X_i/X_j$. C5 requires that $f$ should be expressible in terms of the inter-compositional ratios of the form $x_i/X_i$ or, equivalently, $X_i/x_i$. Thus C4 and C5 together require that $f$ must be expressible in terms of ratios of the form $(x_i/x_j)/(X_i/X_j)$. Moreover, C6 requires that $f$ should use these ratios of ratios in a symmetric way.

Examination of the simplest possible situation, that of defining a scalar measure of difference between two-part composition classes, throws considerable light on what are feasible functions. Here there is only one ratio of ratios, namely $r/R$, where $r = x_1/x_2$, $R = X_1/X_2$, or its inverse, so that $f(x, X)$ can be expressed as $h(r/R)$, a function of a single variable. So far only C4 and C5 have been used. The symmetry requirement of C6 imposes the condition $h(r/R) = h(R/r)$ on $h$. The meaning of this can be expressed more familiarly in terms of $y = \log(r/R)$, and an equivalent function $H(y) = h(r/R)$. Since $r/R = \exp(y)$, $R/r = \exp(-y)$ the restriction on $h$ translates into the restriction $H(y) = H(-y)$ on $H$, that is $H$ must be an even function of $y$.

Once this stage of refinement is reached it seems sensible to place a premium on simplicity and tractability and to ask what is the simplest polynomial in $y$ which will meet the requirements. The evenness of $H$ means that the simplest polynomial would be of the form $H(y) = a + by^2$. Then C2 requires that $h(1) = 0$, or equivalently that $H(0) = 0$, so that $a = 0$ leading to $H(y) = y^2$, on dropping the now unnecessary coefficient $b$. Thus in terms of the original compositions the simplest measure of distance between two-part compositions satisfying all the criteria is

$$f(x_1, x_2; X_1, X_2) = \{\log(x_1/x_2) - \log(X_1/X_2)\}^2 \tag{14}$$

It is easily seen that Eqs. (5) and (6), although satisfying C1–4 and C5, do not satisfy the perturbation invariance requirement C5. That these are not feasible scalar measures of difference is easily seen by considering two-part compositions in the notation above. For Eq. (5) becomes arccos $[(rR + 1)/\{(r^2 + 1)(R^2 + 1)\}^{1/2}]$ and Eq. (6) becomes $(r - R)^2/\{(r^2 + 1)(R^2 + 1)\}$, and neither of those can be expressed in terms of the ratio $r/R$. On the other hand, it is easy to check that the functions defined by Eqs. (7)–(9) are admissible. For example, for two-part compositions, they can be expressed, respectively, as $(r/R - R/r)^2$, $|\log(r/R)|$, $\{\log(r/R)\}^2$, where the last is the simple form discovered in Eq. (14).

Indeed it can be argued that Eq. (9), a generalization of Eq. (14), emerges as the simplest feasible form by an easy extension of the argument applied above for two-part compositions. For this purpose, it is easiest to consider a maximal scalar and perturbation invariant in its simplest form

$$r_i = \{x_i/g(x)\}/\{X_i/g(X)\} \tag{15}$$

where the divisors $g(x) = (x_1 \ldots x_D)^{1/D}$ and $g(X) = (X_1 \ldots X_D)^{1/D}$ are the geometric means of the compositional components. Then any scale, perturbation and permutation invariant scalar measure of difference $f(x, X)$ must be expressible as a symmetric function of $r = (r_1 \ldots, r_D)$, say $h(r)$. Then C3 requires that $h(r) = h(r^{-1})$ and this can be translated, as before, into a more familiar type of requirement on $H(z)$, where $z = \log(r)$, namely $H(z) = H(-z)$. As before, from a wish for simplicity, suppose that the form is limited to a symmetric polynomial in $z$ of second degree:

$$H(z) = a + b \sum_i z_i + c \sum_i z_i^2 + d \sum_{i \neq j} z_i z_j$$

The evenness relationship gives $b = 0$, C2 requires $H(0) = 0$ and so $a = 0$. The $H(z)$ can be written in the form

$$H(z) = (d/2) \left( \sum_i z_i \right)^2 + (c - d/2) \sum_i z_i^2$$

Since $\Sigma_i \, z_i = 0$ it follows that $H(z)$ may be taken as

$$\sum_i z_i^2 = \sum_i \, [\log \{x_i / g(x)\} \, - \, \log \{X_i / g(X)\}]^2 \qquad (16)$$

It is easy to show that, apart from a constant factor, this is identical to form Eq. (9).
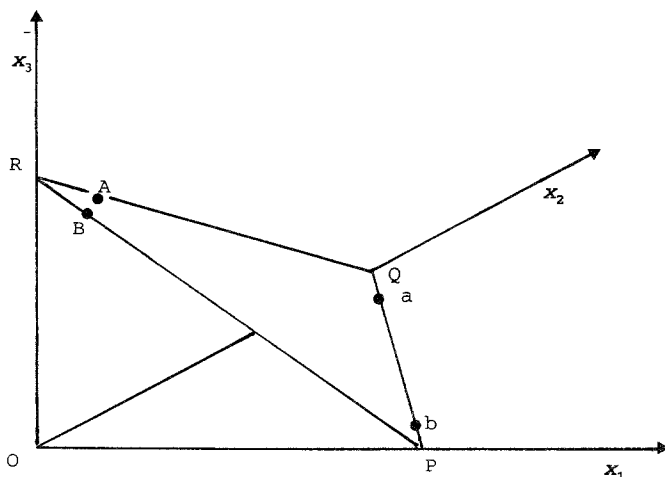
## SUBCOMPOSITIONAL DOMINANCE

There is a further requirement of any feasible scalar measure of difference $f(x, X)$ between two compositions $x$ and $X$ which must be added to C1–C6 above. Consider again two scientists A and B, the first recording a full $D$-part composition $x = (x_1, \ldots , x_D)$ but the second recording only a $D^*$-part subcomposition, namely a subset, say $x^*$, of $D^*$ of the components of $x$. The adoption of a method of determining a scalar measure of difference uses functions $f_D$ for the full $D$-part compositions and $f_{D^*}$ for the $D^*$-part subcomposition. Now the first scientist clearly cannot see less difference in the full compositions than the second scientist sees in the subcompositions; the second scientist has fewer ratios that may vary. Thus, for any proper scalar measure of difference there is a inequality requirement, whereby full compositional differences must dominate subcompositional differences:

C7.  Subcompositional dominance. $f_D(x, X) \geq f_{D^*}(x^*, X^*)$.

A simple example quickly demonstrates that the angular difference Eq. (5) does not satisfy this necessary dominance requirement. For the two three-part compositions (0.01, 0.09, 0.90) and (0.09, 0.01, 0.90), represented by the points A and B in the triangle PQR in Fig. 2, the angular measure AOB between the full compositions is arccos (0.992) = 0.125 radians. The subcompositions consisting of the first two parts are (0.1, 0.9) and (0.9, 0.1), which can be represented as the points $a$ and $b$, the projections from $R$ of the points $A$ and $B$ on to the line segment $PQ$, with angular difference arccos (0.220) = 1.35 radians, larger than for the full compositions. This unsatisfactory feature of the greater angular measure of the subcomposition is clear from Fig. 2.

Since Eq. (6) is a monotonic function of Eq. (5) through the simple relation Eq. (6) = $\sin^2$ Eq. (5), it follows that Eq. (6) does not possess the property of subcompositional dominance. The scalar measures defined by Eqs. (7)–(9) satisfy the subcompositional dominance requirement since they are of the general form

$$f_D(x, X) = \sum_{i < j} f_2(x_i, x_j; X_i, X_j)$$

**Fig. 2.** Representation of two compositions $A(0.01, 0.09, 0.90)$, $B(0.09, 0.01, 0.90)$ and the corresponding subcompositions $a(0.1, 0.9)$, $b(0.9, 0.1)$ of the first two parts. The fact that angle $aOb >$ angle $AOB$ demonstrates the unsatisfactory nature of the angular measure of compositional difference.

and $f_{D*}(x^*, X^*)$ consists of only a subset of the two-part difference terms in the sum and so cannot be greater than $f_D(x, X)$.

## DISTANCE BETWEEN TWO COMPOSITIONS

So far no attempt has been made to require that a scalar measure $f(x, X)$ of difference should have the usual property attributed to mathematical distance, namely the so-called triangular inequality. In the notation of this paper this would require that for any three compositions $x$, $X$, $\xi$ the following inequality must hold;

$$f(x, \xi) + f(\xi, X) \geq f(x, X)$$

To consider whether this is a reasonable requirement, suppose that the changes between the compositions are represented by perturbations $\xi = p \circ x$, $X = q \circ \xi$, $X = r \circ x$. Thus the direct perturbation $r$ from $x$ to $X$ is equivalent to the combination $p \circ q$ of the perturbations through the intermediate composition $\xi$. In reducing the full dimension $d$ of the change implicit in the perturbation relationship to that of a scalar difference function it seems reasonable to require that the difference between $x$ and $X$ assessed directly as $f(x, X)$ should not be greater than the sum of the differences assessed through an intermediate composition, namely $f(x, \xi) + f(\xi, X)$. It is relatively easy to establish that the use of the square roots of either Eqs. (8) or (9), both of which are feasible on criteria

C1–C7, satisfy this distance triangular property. On the other hand, neither Eq. (7) nor its square root satisfy the distance triangular property, as can easily be established by computing the measures for the case

$$D = 3, \quad x = (0.4, 0.2, 0.4), \quad \xi = (0.25, 0.25, 0.50), \quad X = (0.1, 0.1, 0.8)$$

## DISCUSSION

In this paper, consideration has been deliberately confined to one aspect of compositional data analysis, namely the way in which perceptions of the nature of compositions delimit the choice of scalar measures of difference between compositions. Although the various criteria considered easily dismiss some proposed measures of difference, in particular the angular measure so strongly defended as the "unique" measure of difference between two compositions by Watson and Philip (1989) and Watson (1990, 1991), they do not reduce the class to any unique measure. Among those which survive scrutiny by all the criteria it seems that, judged on the basis of simplicity, the logratio difference measure Eq. (9) is by far the simplest.

There are, however, other aspects of compositional data analysis which interact with the study of difference measures considered in this paper. Difference measures concentrate their view of the pattern of variability of compositions on the study of variation between compositions. Many problems require the study of compositions from the dual standpoint of characterizing the variability within compositions. In addition to criteria of invariance there is, within this wider framework (Aitchison, 1992), an important criterion of subcompositional coherence, requiring that there should be agreement between a scientist who studies the full composition and another who confines attention to subcompositions in their statements about the subcompositions. Such criteria limit the range of statistical procedures and, in particular, restrict the definition of meaningful measures of dependence between components of a composition. For example, such considerations suggest that the simplest way of satisfying these criteria is to regard

$$\sum_{i < j} \text{var} \{\log (x_i/x_j)\} \tag{17}$$

as a measure of the total variability of a distribution describing compositional variability.

For a compositional data set consisting of $N$ $D$-dimensional compositions

$$x_r = (x_{r1}, \ldots, x_{rD}) \quad (r = 1, \ldots, N)$$

a measure of total variability could then be obtained by the replacement of the variances in Eq. (17) by their sample estimates

$$(N - 1)^{-1} \sum_{r=1}^{N} \{\log (x_{ri}/x_{rj}) - m_{ij}\}^2 \tag{18}$$

where $m_{ij} = N^{-1} \sum_{r=1}^{N} \log (x_{ri}/x_{rj})$, to obtain a measure $T_1$ of total variability

$$T_1 = (N - 1)^{-1} \sum_{i<j} \sum_{r} \{\log (x_{ri}/x_{rj} - m_{ij}\}^2 \tag{19}$$

From the inter-compositional viewpoint of this paper another reasonable measure $T_2$ of total variability is the sum of the scalar measure of differences Eq. (9) for each possible pair of compositions in the data set:

$$T_2 = \sum_{r<s} \sum_{i<j} \{\log (x_{ri}/x_{rj}) - \log (x_{si}/x_{sj})\}^2 \tag{20}$$

It is a trivial exercise to establish that $T_1$ and $T_2$ are identical, apart from a factor dependent only on $D$ and $N$. Thus, measures of total variability based on the scalar difference Eq. (9) between compositions and based on estimated logratio variances measuring variability within compositions are in conformity, a surely desirable feature.

## ACKNOWLEDGMENT

## REFERENCES

Aitchison, J., 1986a, The Statistical Analysis of Compositional Data: Chapman and Hall, London.
Aitchison, J., 1986b, CODA: A Microcomputer Package for the Statistical Analysis of Compositional Data: Chapman and Hall, London.
Aitchison, J., 1990, Comment on "Measures of Variability for Geological Data" by D. F. Watson and G. M. Philip: Math. Geol., v. 22, p. 223–226.
Aitchison, J., 1991, Delusions of Uniqueness and Ineluctability: Math. Geol., v. 23, p. 275–277.
Aitchison, J., 1992, Principles of Compositional Data Analysis, in Anderson, T. W., Olkin, I. and Fang, K. T. (Eds.), Proceedings of the International Symposium of Multivariate Analysis and its Applications, Institute of Mathematical Statistics. To appear.
Lehmann, E. L., 1959, Testing Statistical Hypotheses: New York, Wiley.
Watson, D. F., 1990, Reply to Comment on "Measures of Variability for Geological Data" by D. F. Watson and G. M. Philip: Math. Geol., v. 22, p. 227–231.
Watson, D. F., 1991, Reply to "Delusions of Uniqueness and Ineluctability" by J. Aitchison: Math. Geol., v. 23, p. 279.
Watson, D. F., and Philip, G. M., 1989, Measures of Variability for Geological Data: Math. Geol., v. 21, p. 233–254.
Woronow, A., and Love, K. M., 1990, Quantifying and Testing Differences Among Means of Compositional Data Suites: Math. Geol., v. 22, p. 837–852.