

3. Log-ratio transformation

In most aspects of metagenomic analysis we are interested in comparing the content of two or more samples, either to determine how identical or how different they are or to explore if certain variables associated with the samples can explain the variance in the set of samples. In order to make this comparison on a quantitative level, we need to define a distance between two samples.

3.1 Linear algebra

Before we can define the distance between two compositions, we need to define the simplex as a normed vector space. Let us first recall how real vector spaces in \mathbb{R}^D are defined. In any real space, that is for non-compositional data, distances are given by the Euclidean metric,

Definition 3.1.1 — Euclidean distance.

$$d_e(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_e = \sqrt{\sum_{i=1}^D (q_i - p_i)^2}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^D, \quad (3.1)$$

where subscript e denotes the Euclidean distance. The metric is induced by the norm,

Definition 3.1.2 — Euclidean Norm.

$$\|\mathbf{x}\|_e = \sqrt{\mathbf{x} \cdot \mathbf{x}} \quad \mathbf{x} \in \mathbb{R}^D, \quad (3.2)$$

which, for real vectors, is the length of a vector \mathbf{x} , where the dot product is defined as,

Definition 3.1.3 — Dot product (Euclidean inner product).

$$\mathbf{x} \cdot \mathbf{y} = \langle \mathbf{x}, \mathbf{y} \rangle_e = \sum_{i=1}^D x_i y_i, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^D. \quad (3.3)$$

These formulas are familiar to anyone who have studied ordinary linear algebra, but they are not appropriate when working with compositional data as illustrated by the following example.

■ **Example 3.1** Consider the four compositions

$$\begin{bmatrix} 5 \\ 65 \\ 30 \end{bmatrix} \quad \begin{bmatrix} 10 \\ 60 \\ 30 \end{bmatrix} \quad \begin{bmatrix} 50 \\ 20 \\ 30 \end{bmatrix} \quad \begin{bmatrix} 55 \\ 15 \\ 30 \end{bmatrix}$$

Using definition 3.1.1, we can determine the Euclidean distance between the first two, and the last two compositions. In both cases $d_e \approx 7.07$, meaning that from a Euclidean geometric point of view, the first two and the last two compositions are equally far apart. However, if we look at the proportions, as is appropriate to do for compositions, the first part doubles between the first two compositions, whereas for the last two compositions, it only increases by 10%. The second part decreases by 7.7% between the two first compositions while it decreases by 30% between the second two. From a compositional point of view, the latter two are much closer than the former two. ■

A vector space is defined as a set V which is closed under operations of addition and scalar multiplication. ‘Closed under’ means that the addition of two elements of V should yield a new element which is also part of V and for scalar multiplication, that $\lambda \mathbf{x} \in V$ for any $\mathbf{x} \in V$ and $\lambda \in \mathbb{R}$. For real vectors \mathbf{r}, \mathbf{q} in a real vector space on \mathbb{R}^D , the following algebraic rules apply,

Definition 3.1.4 — Sum and scalar multiplication.

$$\mathbf{r} + \mathbf{q} = (r_1 + q_1, r_2 + q_2, \dots, r_D + q_D) \quad (3.4)$$

$$\alpha \mathbf{x} = (\alpha x_1, \alpha x_2, \dots, \alpha x_D) \quad (3.5)$$

These are rules we are familiar with from linear algebra and this is how we know real vectors behave. However, as we saw in the previous chapter, this kind of addition does not work on compositions as there are no neutral element nor inverse elements. What about scalar multiplication? This is exactly how we defined equivalence classes in the previous chapter, which means that scalar multiplying a compositions by a real number does not yield a new composition, but the same. It is clear that we can not use these operations to define the simplex as a vector space

3.2 The Aitchison geometry

We need find a metric on the compositional sampling space, the simplex, that gives it Euclidean properties. The resulting geometry is called Aitchison geometry, named after the Scottish mathematician who first defined it in the early 1980’s.

3.2.1 Vector space properties

We should define two new operations which are analogous to addition and scalar multiplication, which operates on elements in the simplex to yield a new element in the simplex. They should obey the following 8 axioms:

$$1. \quad \mathbf{v} + \mathbf{w} = \mathbf{w} + \mathbf{v} \quad (\text{commutative law of addition}) \quad (3.6)$$

$$2. \quad (\mathbf{v} + \mathbf{w}) + \mathbf{x} = \mathbf{v} + (\mathbf{w} + \mathbf{x}) \quad (\text{associative law of addition}) \quad (3.7)$$

$$3. \quad \mathbf{v} + \mathbf{0} = \mathbf{v} \quad (\text{additive identity law}) \quad (3.8)$$

$$4. \quad \mathbf{v} + (-\mathbf{v}) = \mathbf{0} \quad (\text{additive inverse law}) \quad (3.9)$$

$$5. \quad r(\mathbf{v} + \mathbf{w}) = r\mathbf{v} + r\mathbf{w} \quad (\text{distributive law}) \quad (3.10)$$

$$6. \quad (r + s)\mathbf{v} = r\mathbf{v} + s\mathbf{v} \quad (\text{distributive law}) \quad (3.11)$$

$$7. \quad r(s\mathbf{v}) = (rs)\mathbf{v} \quad (\text{associative law of multiplication}) \quad (3.12)$$

$$8. \quad 1\mathbf{v} = \mathbf{v} \quad (\text{scalar identity law}) \quad (3.13)$$

for $\mathbf{v}, \mathbf{w}, \mathbf{x} \in \mathcal{S}^D$ and $r, s \in \mathbb{R}$. It should be noted that the $\mathbf{0}$ vector in axiom 3, is not necessarily a vector with zeros in each entry (which does not exist in the simplex). It is the generalized neutral element, which, when added to any other element, returns the same element. Likewise, the minus sign in axiom 4 does not necessarily mean the negative valued vector, but can be any inverse element, which when added to the element itself yields the neutral element. We can quickly convince ourselves that ordinary addition and scalar multiplication break several of these axioms. Instead we have two other operations, perturbation and powering, that works. They are defined as,

Definition 3.2.1 — Perturbation and powering.

$$\mathbf{x} \oplus \mathbf{y} = \mathcal{C}(x_1 y_1, x_2 y_2, \dots, x_D y_D) \quad (3.14)$$

$$\alpha \odot \mathbf{x} = \mathcal{C}(x_1^\alpha, x_2^\alpha, \dots, x_D^\alpha) \quad (3.15)$$

The circle around the plus and multiplication signs symbolizes that the operations are analogous to addition and scalar multiplication, but not the same. It is left as an exercise to check that they obey the 8 axioms of vector spaces. The meaning of perturbation and powering is hard to visualize, but they are used to scale and center a composition, something that we will encounter in a later lecture.

Just like for Euclidean real vector spaces, we can define an inner product,

Definition 3.2.2 — Aitchison inner product.

$$\langle \mathbf{x}, \mathbf{y} \rangle_a = \frac{1}{2D} \sum_{i=1}^D \sum_{j=1}^D \ln \frac{x_i}{x_j} \ln \frac{y_i}{y_j} \quad \mathbf{x}, \mathbf{y} \in \mathcal{S}^D \quad (3.16)$$

for compositions on \mathcal{S}^D . This function obeys the three rules, scalar invariance, permutation invariance, and subcompositional coherence. The inner product gives us the norm,

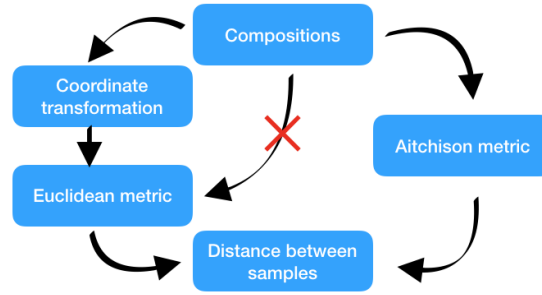


Figure 3.1: The two routes to obtaining the distance between compositions.

Definition 3.2.3 — Aitchison norm.

$$\|\mathbf{x}\|_a = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_a} = \sqrt{\frac{1}{2D} \sum_{i=1}^D \sum_{j=1}^D \left(\ln \frac{x_i}{x_j} \right)^2} \quad \mathbf{x} \in \mathcal{S}^D \quad (3.17)$$

and from the norm we get the Aitchison distance between compositions,

Definition 3.2.4 — Aitchison distance.

$$d_a(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} \ominus \mathbf{y}\|_a = \sqrt{\frac{1}{2D} \sum_{i=1}^D \sum_{j=1}^D \left(\ln \frac{x_i}{x_j} - \ln \frac{y_i}{y_j} \right)^2} \quad \mathbf{x}, \mathbf{y} \in \mathcal{S}^D. \quad (3.18)$$

The generalized minus in the above definition denotes perturbation by the inverse. Together with the Aitchison geometry, the simplex becomes a linear vector space known as a (finite dimensional) real Hilbert space, where standard Euclidean properties are valid, such as the triangle inequality, Pythagoras theorem, and the Cauchy-Schwartz inequality. We denote this vector space $(\mathcal{S}^D, \oplus, \odot)$ as opposed to ordinary real vector spaces $(\mathbb{R}^D, +, \cdot)$.

3.3 Transformations

Unfortunately, the Aitchison distance is somewhat cumbersome to work with due to the double summations. We will therefore introduce certain coordinate transformations, which will transform the compositions from the simplex and onto the real Euclidean space, where the ordinary Euclidean metric applies. This means that rather than calculating the distance between compositions using the above defined Aitchison distance, we transform the compositions from \mathcal{S}^D onto a subspace of \mathbb{R}^D and then use all the normal methods that apply to ordinary real vectors. The reason why this is a preferred route is that many statistical and explorative data analysis methods implicitly assumes Euclidean distances and they come pre-implemented in many modern programming languages, such as R, Python, Matlab, etc. So in order to not having to reimplement everything using Aitchison geometry, we can just transform our compositional data and work with the packages as they are.

As of today, we know of three different coordinate transformations which allows us to apply Euclidean metric after transforming the data. Each transformation has pros and cons

and in the end, it is up to the analyst to choose the most appropriate for the problem at hand. There are no strict rules to which transformation should be used in a given situation, but sometimes one choice is more obvious than another.

3.3.1 Additive logratio transformation (ALR)

The first and most intuitive transformation is the *Additive logratio transformation*. This transformation maps a composition from \mathcal{S}^D onto \mathbb{R}^{D-1} , that is, to a subspace of the D -dimensional real space. Formally it is defined as,

Definition 3.3.1 — Additive logratio (ALR) coordinates. Given a D -part composition $\mathbf{x} = (x_1, x_2, \dots, x_D)$ and using x_D as reference part, the alr transformation is given by

$$\text{alr}(\mathbf{x}) = \left(\ln \frac{x_1}{x_D}, \ln \frac{x_2}{x_D}, \dots, \ln \frac{x_{D-1}}{x_D} \right) = \boldsymbol{\zeta}. \quad (3.19)$$

and its inverse

$$\text{alr}^{-1}(\boldsymbol{\zeta}) = \mathcal{C}(\exp(\zeta_1), \exp(\zeta_2), \dots, \exp(\zeta_{D-1}), 1) = \mathbf{x}. \quad (3.20)$$

In the ALR transformation, one part (by convention, usually the last part of the composition) is chosen as denominator. The interpretation of this is, that all parts are given relative to one part, i.e., for a nutrition table of a food product, fat per protein, carbohydrates per protein, salt per protein, etc. The choice of denominator is up to the analyst and must be chosen to be meaningful with respect to the question at hand.

The reference part becomes zero in ALR-space, because the ratio x_D/x_D is 1 and the logarithm of 1 is zero. ALR is also not a unique transformation, since it depends on the choice of denominator part. Taking the logarithm ensures that the algebraic operations perturbation and powering on the simplex \mathcal{S}^D translates into the algebraic operations sum and multiplication in \mathbb{R}^{D-1}

$$\text{alr}(\alpha \odot \mathbf{x} \oplus \beta \odot \mathbf{y}) = \alpha \cdot \text{alr}(\mathbf{x}) + \beta \cdot \text{alr}(\mathbf{y}).$$

This relation makes the ALR transformation an isomorphism and it ensures that real space operations sum and multiplication can be used, but not Euclidean distance. The reason for this is that ALR does not provide an isometry between \mathcal{S}^D and \mathbb{R}^{D-1} and therefore, distances between vectors are not invariant under the ALR transformation. The ALR transformation is frequently used in genomics, even if used implicitly (e.g., so called log(FPKM) values), and it is easy to interpret, but one should be careful with ALR if any kind of intersample distance or metrics are involved in the analysis. The rule is that analysis methods can be applied to ALR values if they are *affine equivalent*. Methods are affine equivalent if the results are the same after translating, rotating, or scaling the data, so that the result is translated, rotated or scaled as well. This means that only algebraic vector space operations are involved and no metric concepts. It is not always obvious when this is the case.

3.3.2 Centered logratio transformation (CLR)

Another transformation is the *Centered logratio transformation*. The CLR transformation uses a fixed constant as denominator, rather than a specific part and is formally defined as,

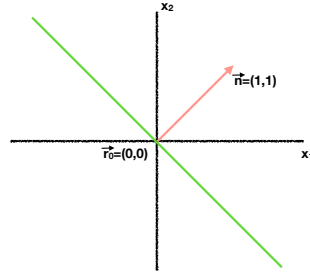


Figure 3.2: The green line is the subplane of \mathbb{R}^2 in which all 2-part CLR values exist. The plane is defined as the space going through $(0,0)$ and having the vector $(1,1)$ as normal. This generalizes to arbitrary dimensions.

Definition 3.3.2 — Centered logratio (CLR) values. Given a D -part composition $\mathbf{x} = (x_1, x_2, \dots, x_D)$, CLR values are given by,

$$\text{clr}(\mathbf{x}) = \left(\ln \frac{x_1}{g_m(\mathbf{x})}, \ln \frac{x_2}{g_m(\mathbf{x})}, \dots, \ln \frac{x_D}{g_m(\mathbf{x})} \right) = \boldsymbol{\zeta}. \quad (3.21)$$

where g_m is the geometric mean

$$g_m(\mathbf{x}) = \left(\prod_{i=1}^D x_i \right)^{1/D} = \exp \left(\frac{1}{D} \sum_{i=1}^D \ln x_i \right) \quad (3.22)$$

and its inverse

$$\text{clr}^{-1}(\boldsymbol{\zeta}) = \mathcal{C}(\exp(\zeta_1), \exp(\zeta_2), \dots, \exp(\zeta_D)) = \mathbf{x}. \quad (3.23)$$

The definition is very similar to ALR. The only difference is that the denominator in the ratio is the geometric mean of \mathbf{x} , rather than a specific part. We could, in principle use any constant as denominator, due to the fact that compositions are equivalent classes, but when using the geometric mean, the sum of CLR values will always be 0. The interpretation of CLR values is that each part is given relative to the mean of all parts: positive values are larger than the mean and negative part are smaller than the mean. A CLR transformed vector will, in general, have as many parts as the original composition, so the CLR transform maps coordinates from \mathcal{S}^D onto \mathbb{R}^D . However, because the sum of a CLR vector is always 0, the vector is constrained to a subspace of \mathbb{R}^D as illustrated in Fig. 3.2. CLR values, always fall on a $(D-1)$ -dimensional plane in \mathcal{R}^D defined by the point $\mathbf{r}_0 = (0, 0, \dots, 0)$ and the normal vector $\mathbf{n} = (1, 1, \dots, 1)$. Because of this constraint, CLR values do not span \mathbb{R}^D and they are therefore not coordinates in any orthonormal basis of \mathbb{R}^D . Hence we refer to CLR values as *CLR coefficients*, rather than CLR coordinates.

CLR values obey the rules

$$\text{clr}(\alpha \odot \mathbf{x} \oplus \beta \odot \mathbf{y}) = \alpha \cdot \text{clr}(\mathbf{x}) + \beta \cdot \text{clr}(\mathbf{y}) \quad (3.24)$$

$$\langle \mathbf{x}, \mathbf{y} \rangle_a = \langle \text{clr}(\mathbf{x}), \text{clr}(\mathbf{y}) \rangle_e \quad (3.25)$$

$$d_a(\mathbf{x}, \mathbf{y}) = d_e(\text{clr}(\mathbf{x}), \text{clr}(\mathbf{y})), \quad (3.26)$$

which means that CLR transformation is both an isomorphism and an isometry and it therefore conserves the metric properties of a composition.

The downside of the CLR transformation is that, due to the fact that they are confined to a subplane of \mathbb{R}^D , the sum of its values is zero, which makes the determinant of the covariance matrix of a set of compositions, zero. This means that certain statistical methods that rely on the covariance matrix being non-singular cannot be applied to CLR values. When the determinant of the covariance matrix is zero, it means that the parts are perfectly correlated, which is true by construction for CLR values, as illustrated in the following example.

■ **Example 3.2** Consider the 2 2-part compositions

$$\mathbf{x} = \begin{pmatrix} 40 \\ 60 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 20 \\ 80 \end{pmatrix}$$

From eq. 3.22 we can calculate the geometric means,

$$g_m(\mathbf{x}) \approx 49, \quad g_m(\mathbf{y}) \approx 40$$

and then the CLR transform,

$$\text{clr}(\mathbf{x}) = \begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}, \quad \text{clr}(\mathbf{y}) = \begin{pmatrix} -0.69 \\ 0.69 \end{pmatrix}$$

The CLR values are symmetric around 0 (as they should be), there is only 1 independent variable (they are confined to a 1D plane in 2D space), and there is a perfect correlation between the parts (one goes up and the other goes down, by the exact same amount). ■

Another problem with the CLR transformation is that the values are not subcompositional coherent, because the geometric mean will change if parts are removed, and thus the CLR values will generally be different for different subcompositions. This gives rise to a *very important* limitation of the CLR transformation: we cannot pre-calculate CLR values and chose to analyze a subset of the values. CLR values will have to be calculated for a given subcomposition and they are only valid within the scope of that exact subcomposition.

3.3.3 Orthonormal coordinates

Coordinates of a vector are expressed with respect to a orthonormal basis, typically the canonical basis of \mathbb{R}^D , $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_D\} = \{[1, 0, \dots, 0], [0, 1, \dots, 0], \dots, [0, 0, \dots, 1]\}$. Any vector $\mathbf{v} \in \mathbb{R}^D$ can be written in the form,

$$\mathbf{v} = v_1 \mathbf{e}_1 + v_2 \mathbf{e}_2 + \dots + v_D \mathbf{e}_D = \sum_{i=1}^D v_i \cdot \mathbf{e}_i. \quad (3.27)$$

The canonical basis of \mathbb{R}^D is not, however, a basis with respect to \mathcal{S}^D since the vectors are not part of the simplex themselves. Just look at all the zeros. But we can make the canonical basis into a spanning set for \mathcal{S}^D by taking the closure of the exponentials,

$$\begin{aligned}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D\} &= \{\mathcal{C}(\exp(\mathbf{e}_1)), \mathcal{C}(\exp(\mathbf{e}_2)), \dots, \mathcal{C}(\exp(\mathbf{e}_D))\} \\ &= \{\mathcal{C}[e, 1, \dots, 1], \mathcal{C}[1, e, \dots, 1], \dots, \mathcal{C}[1, 1, \dots, e]\}.\end{aligned}\quad (3.28)$$

The vectors \mathbf{w}_i span the simplex, but they are not linearly independent, and therefore not a basis, since there are D vectors in the set and the dimensionality of \mathcal{S}^D is $D - 1$.

We can express compositions in the same form as eq. 3.27 by replacing operations by their simplex analogous,

$$\begin{aligned}\mathbf{x} &= \bigoplus_{i=1}^D \ln \frac{x_i}{g_m(\mathbf{x})} \odot \mathbf{w}_i \\ &= \ln \frac{x_1}{g_m(\mathbf{x})} \odot [e, 1, \dots, 1] \oplus \ln \frac{x_2}{g_m(\mathbf{x})} \odot [1, e, \dots, 1] \oplus \dots \oplus \ln \frac{x_D}{g_m(\mathbf{x})} \odot [1, 1, \dots, e] \\ &= \left[\frac{x_1}{g_m(\mathbf{x})}, \frac{x_2}{g_m(\mathbf{x})}, \dots, \frac{x_D}{g_m(\mathbf{x})} \right] \\ &= [x_1, x_2, \dots, x_D],\end{aligned}\quad (3.29)$$

where we have used the definition of compositions as equivalence classes in the final step. However, $\{\mathbf{w}\}$ is not a basis and therefore, $\ln \frac{x_i}{g_m(\mathbf{x})}$, which are recognized as CLR values, are not coordinates. We can drop any one of the \mathbf{w}_i 's to get a proper basis, but, by taking the Aitchison dot product between any two vectors in this basis, it is easily seen that it is not an orthogonal basis. We will explore how to build a orthonormal basis in Sect. 3.3.5, but for now let us assume that such a basis exist and is denoted $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{D-1}\}$. Then we can define the *contrast matrix*,

Definition 3.3.3 — Contrast matrix. Given an orthonormal basis of the simplex \mathcal{S}^D , $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{D-1}\}$, the *contrast matrix* Ψ is a $(D, D - 1)$ -matrix where each row $\Psi_i = \text{clr}(\mathbf{e}_i)$, $i = 1, 2, \dots, D - 1$. Each row is a logcontrast.

An orthonormal basis has the property that the dot-product between basis vectors is zero, while the dot-product between a basis vector and itself is 1,

$$\langle \mathbf{e}_i, \mathbf{e}_j \rangle_a = \langle \text{clr}(\mathbf{e}_i), \text{clr}(\mathbf{e}_j) \rangle_e = \delta_{ij}, \quad (3.30)$$

where δ_{ij} is the Kronecker-delta which is a function that equals 0 for $i \neq j$ and 1 for $i = j$. From this it follows that,

$$\Psi \Psi^T = \mathbf{I}_{D-1},$$

that is, the identity matrix in $D - 1$ dimensions.

3.3.4 Isometric logratio transformation (ILR)

We can now introduce a third transformation which is known as the *Isometric logratio transformation*. As the name suggests, it too conserves distances (it is an isometry) and contrary to CLR, it does not have a singular covariance matrix. It has all the benefits of

both ALR and CLR and none of the downsides. This comes at a cost however: ILR values are expressed as coordinates in an orthonormal basis of the simplex, meaning that the ILR coordinates can be difficult to interpret.

As discussed above, we can obtain coordinates by expressing a compositions in an orthonormal basis according to Eq. 3.29. We can furthermore replace the Aichison operators by the Euclidean analogous, by using the CLR transform,

$$\text{CLR}(\mathbf{x}) = \sum_{i=1}^D x_i^* \cdot \text{CLR}(\mathbf{e}_i) = \mathbf{x}^* \cdot \mathbf{\Psi} \iff \mathbf{x}^* = \text{CLR}(\mathbf{x}) \cdot \mathbf{\Psi}^T \quad (3.31)$$

This transformation is called ILR and is basically obtained by multiplying the CLR transformation with a contrast matrix $\mathbf{\Psi}$ with dimensions $(D-1, D)$,

Definition 3.3.4 — Isometric logratio (ILR) coordinates. Given a D -part composition $\mathbf{x} = (x_1, x_2, \dots, x_D)$ and a contrast matrix $\mathbf{\Psi}_{D-1, D}$ based on the basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{D-1}\}$, ILR coordinates are given by

$$\text{ilr}(\mathbf{x}) = \text{clr}(\mathbf{x}) \cdot \mathbf{\Psi}^T = \mathbf{x}^* \quad (3.32)$$

and its inverse

$$\text{ilr}^{-1}(\mathbf{x}^*) = \mathcal{C}(\exp(\mathbf{x}^* \mathbf{\Psi})) = \mathbf{x} \quad (3.33)$$

ILR coordinates obey the same properties as CLR coefficients (3.24) but they are true coordinates in \mathbb{R}^{D-1} . The downside is, that each part in the ILR-transformed vector is a linear combination of parts from the composition. ILR transforms are also not unique, since infinitely many orthonormal bases can be constructed. In any case, before the ILR transform can be applied, a basis has to be constructed. One option is to use *balances* as basis vectors. These are formed by applying a sequential binary partition of the parts of a composition. Many computational CoDa packages have the option to use this basis for ILR transformation.

3.3.5 Balances

A popular choice of basis to use in ILR transformation are balances. Balances are formed from the sequential binary partition of the composition and represent groups of features. This partition is a hierarchical division of the parts, splitting the parts into two groups, each of which are again split into two groups and so on and until there are $D - 1$ balances. An example of a binary partition table of a six-part composition is

x_1	x_2	x_3	x_4	x_5	x_6
+1	+1	+1	-1	-1	-1
+1	+1	-1	0	0	0
+1	-1	0	0	0	0
0	0	0	+1	+1	-1
0	0	0	+1	-1	0

The coordinates of this partition are called balances and the vectors are called balancing elements. In this example we have chosen to make the first balance an even split between

parts (x_1, x_2, x_3) versus (x_4, x_5, x_6) , but we could just as well have chosen to split (x_1, x_2) versus versus (x_3, x_4, x_5, x_6) .

When we have decided on a binary partition table, we can form the balances by counting the number of positive and negative entries per row. Let us denote the number of plus-signs by r and the number of minus-signs by s and let us refer to the row in the partition table by k . In the example above, for the first row, $k = 1$, $r = 3$ and $s = 3$. For $k = 2$, $r = 2$ and $s = 1$, and so on. A balance is defined as the (normalized) logratio of the geometric mean of the two groups (plus- and minus group),

$$b_k = \sqrt{\frac{rs}{r+s}} \ln \frac{(x_{i_1} x_{i_2} \dots x_{i_r})^{1/r}}{(x_{j_1} x_{j_2} \dots x_{j_s})^{1/s}}, \quad (3.34)$$

where the square root in front is the normalizing factor. We can use the logarithm rules to formulate this as,

$$\begin{aligned} b_k &= \ln \frac{(x_{i_1} x_{i_2} \dots x_{i_r})^{a_+}}{(x_{j_1} x_{j_2} \dots x_{j_s})^{a_-}} \\ &= \ln(x_{i_1} x_{i_2} \dots x_{i_r})^{a_+} - \ln(x_{j_1} x_{j_2} \dots x_{j_s})^{a_-} = \sum_{j=1}^D a_{kj} \ln x_j \end{aligned} \quad (3.35)$$

where

$$a_+ = +\frac{1}{r} \sqrt{\frac{rs}{r+s}}, \quad a_- = -\frac{1}{s} \sqrt{\frac{rs}{r+s}}. \quad (3.36)$$

We recognize Eq. 3.35 as a logcontrast as introduced in Chapter 2, Eq. 2.5. So if we form a matrix out of the balances, then each row k is a logcontrast and it is also a CLR transform (they sum to zero), they are a linearly independent spanning set and they are normalized. Hence they form an orthogonal basis and the matrix is a contrast matrix Ψ . For the example partition table above, the contrast matrix of balances is,

x_1	x_2	x_3	x_4	x_5	x_6
$+\frac{1}{3} \sqrt{\frac{3 \cdot 3}{3+3}}$	$+\frac{1}{\sqrt{6}}$	$+\frac{1}{\sqrt{6}}$	$-\frac{1}{\sqrt{6}}$	$-\frac{1}{\sqrt{6}}$	$-\frac{1}{\sqrt{6}}$
$+\frac{1}{2} \sqrt{\frac{2 \cdot 1}{2+1}}$	$+\frac{1}{\sqrt{6}}$	$-\sqrt{\frac{2}{3}}$	0	0	0
$+\frac{1}{1} \sqrt{\frac{1}{2}}$	$-\frac{1}{\sqrt{2}}$	0	0	0	0
0	0	0	$+\frac{1}{\sqrt{6}}$	$+\frac{1}{\sqrt{6}}$	$-\sqrt{\frac{2}{3}}$
0	0	0	$+\frac{1}{1} \sqrt{\frac{1}{2}}$	$-\frac{1}{\sqrt{2}}$	0

The sequential binary partition is not unique in the sense that we can decide ourselves how many parts should go into each group. In the above example, we chose to divide the composition into two groups of three parts each in the first step, but we could just as well have chose one group with two parts and one with four. It is up to the analyst to chose a partition which optimizes the interpretability of the result.

3.3.6 Example transformations

We will end this chapter by logratio transforming a composition, using all three different logratio transforms.

■ **Example 3.3** Consider the composition $\mathbf{x} = (25, 30, 45)$, which is closed to 100.

For the ALR transform, we need to choose a part. In this example, we choose the third part x_3 .

$$\text{ALR}(\mathbf{x}) = \ln \begin{pmatrix} 25/45 \\ 30/45 \end{pmatrix} \approx \begin{pmatrix} -0.59 \\ -0.41 \end{pmatrix} \quad (3.37)$$

For the CLR transform, we need to calculate the geometric mean,

$$g_m(\mathbf{x}) = (25 \cdot 30 \cdot 45)^{1/3} \approx 32.32, \quad (3.38)$$

which then gives the CLR values,

$$\text{CLR}(\mathbf{x}) = \ln \begin{pmatrix} 25/32.32 \\ 30/32.32 \\ 45/32.32 \end{pmatrix} \approx \begin{pmatrix} -0.26 \\ -0.07 \\ 0.33 \end{pmatrix} \quad (3.39)$$

For the ILR transform, we need to build an orthonormal basis. We choose a basis made from balances, using the following sequential partition table,

x_1	x_2	x_3
+1	+1	-1
+1	-1	0

From this table we calculate the balances to form the contrast matrix,

$$\Psi = \begin{pmatrix} \frac{1}{2}\sqrt{\frac{1.2}{1+2}} & \frac{1}{2}\sqrt{\frac{1.2}{1+2}} & -\frac{1}{1}\sqrt{\frac{1.2}{1+2}} \\ \frac{1}{1}\sqrt{\frac{1.1}{1+1}} & -\frac{1}{1}\sqrt{\frac{1.1}{1+1}} & 0 \end{pmatrix} \approx \begin{pmatrix} 0.41 & 0.41 & -0.82 \\ 0.71 & -0.71 & 0 \end{pmatrix}$$

The ILR coordinates are obtained by matrix multiplying the CLR values by the transposed contrast matrix,

$$\text{ILR}(\mathbf{x}) = \begin{pmatrix} -0.26 \\ -0.07 \\ 0.33 \end{pmatrix} \cdot \begin{pmatrix} 0.41 & 0.41 & -0.82 \\ 0.71 & -0.71 & 0 \end{pmatrix}^T \approx \begin{pmatrix} -0.41 \\ -0.13 \end{pmatrix} \quad (3.40)$$

If we plot the ILR coordines in a Cartesian coordinate system, the axis would represent the logarithm of $x_1 + x_2 - x_3$ and $x_1 - x_2$, respectively. ■

3.4 Exercises

Exercise 3.1 Consider the two vectors $\mathbf{x} = [0.7, 0.5, 0.8]$ and $\mathbf{y} = [0.25, 0.75, 0.5]$. Perturb one vector by the other, with and without previous closure. Are there any differences? ■

Exercise 3.2 Compute the Aitchison inner product of $\mathbf{x} = \mathcal{C}[0.7, 0.4, 0.8]$ and $\mathbf{y} = \mathcal{C}[2, 8, 1]$. Are they orthogonal? ■

Exercise 3.3 Consider the two six-part compositions, given in percentages,

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} 3.74 & 9.35 & 16.82 & 18.69 & 23.36 & 28.04 \\ 9.35 & 28.04 & 16.82 & 3.74 & 18.69 & 23.36 \end{pmatrix}.$$

Calculate the Aitchison norms and the Aitchison inner product. What is the angle between the two compositions?

HINT: $\cos(\theta) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$ ■

Exercise 3.4 Compute the Aitchison norm of $\mathbf{x} = \mathcal{C}[0.7, 0.4, 0.8]$ and call it a . Compute $\alpha \odot \mathbf{x}$ with $\alpha = 1/a$. Compute the Aitchison norm of the resulting composition. How do you interpret the result? ■

Exercise 3.5 Redo Exercise 2.5 from Chapter 2, but using the Aitchison distance. Is it subcompositionally dominant? ■

Exercise 3.6 Using the data from Exercise 2.3, Compute the CLR coefficients. Verify that the sum of the transformed components equals zero. ■

Exercise 3.7 Using the data from Exercise 2.3, apply the ALR transformation to the compositions. Plot the transformed data in \mathbb{R}^2 . Now do it using a different part as denominator for the ALR transformation. Compare the results. ■

Exercise 3.8 Show that the exponential of the canonical basis,

$$\begin{aligned} \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{D-1}\} &= \{\mathcal{C}(\exp(\mathbf{e}_1)), \mathcal{C}(\exp(\mathbf{e}_2)), \dots, \mathcal{C}(\exp(\mathbf{e}_{D-1}))\} \\ &= \{\mathcal{C}[e, 1, \dots, 1], \mathcal{C}[1, e, \dots, 1], \dots, \mathcal{C}[1, \dots, e, 1]\} \end{aligned}$$

is not orthogonal on \mathcal{S}^D . ■

Exercise 3.9 Build a sequential partition table of a three part composition and calculate the corresponding balances. Write out the contrast matrix. ■

Exercise 3.10 Using the contrast matrix from the previous exercise, redo Exercise 3.7, with the ILR transformation. Plot the result. Compare with a different contrast matrix. ■