

Chapter 4

Measures of distance between samples: Euclidean

We will be talking a lot about distances in this book. The concept of distance between two samples or between two variables is fundamental in multivariate analysis – almost everything we do has a relation with this measure. If we talk about a single variable we take this concept for granted. If one sample has a pH of 6.1 and another a pH of 7.5, the distance between them is 1.4: but we would usually call this the absolute difference. But on the pH line, the values 6.1 and 7.5 are at a distance apart of 1.4 units, and this is how we want to start thinking about data: points on a line, points in a plane, ... even points in a ten-dimensional space! So, given samples with not one measurement on them but several, how do we define distance between them. There are a multitude of answers to this question, and we devote three chapters to this topic. In the present chapter we consider what are called *Euclidean* distances, which coincide with our most basic physical idea of distance, but generalized to multidimensional points.

Contents

Pythagoras' theorem
Euclidean distance
Standardized Euclidean distance
Weighted Euclidean distance
Distances for count data
Chi-square distance
Distances for categorical data

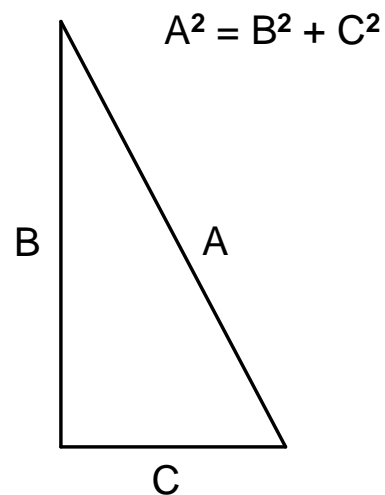


Pythagoras' theorem

The photo shows Michael in July 2008 in the town of Pythagorion, Samos island, Greece, paying homage to the one who is reputed to have made almost all the content of this book possible: ΠΥΘΑΓΟΡΑΣ Ο ΣΑΜΙΟΣ, Pythagoras the Samian. The illustrative geometric proof of Pythagoras' theorem stands carved on the marble base of the statue – it is this theorem that is at the heart of most of the multivariate analysis presented in this book, and particularly the graphical approach to data analysis that we are strongly promoting. When you see the word “square” mentioned in a statistical text (for example, chi square or least squares), you can be almost sure that the corresponding theory has some relation to this theorem. We first show the theorem in its simplest and most familiar two-dimensional form, before showing how easy it is to generalize it to multidimensional space. In a right-

angled triangle, the square on the hypotenuse (the side denoted by A in Exhibit 4.1) is equal to the sum of the squares on the other two sides (B and C); that is, $A^2 = B^2 + C^2$.

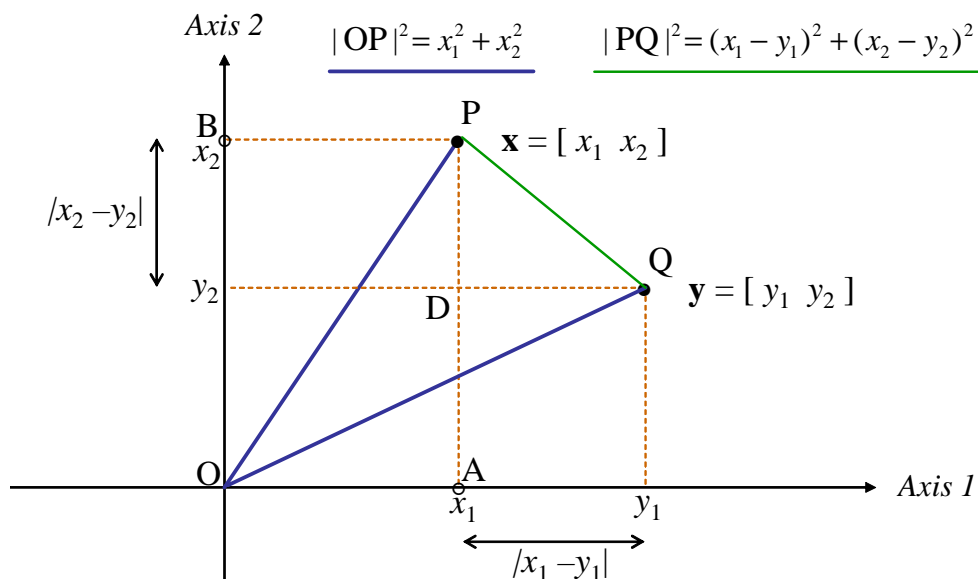
Exhibit 4.1 Pythagoras' theorem in the familiar right-angled triangle, and the monument to this triangle in the port of Pythagorion, Samos island, Greece, with Pythagoras himself forming one of the sides.



Euclidean distance

The immediate consequence of this is that the squared length of a vector $\mathbf{x} = [x_1 \ x_2]$ is the sum of the squares of its coordinates (see triangle OPA in Exhibit 4.2, or triangle OPB – $|\text{OP}|^2$ denotes the squared length of \mathbf{x} , that is the distance between point O and P); and the

Exhibit 4.2 Pythagoras' theorem applied to distances in two-dimensional space.



squared distance between two vectors $\mathbf{x} = [x_1 \ x_2]$ and $\mathbf{y} = [y_1 \ y_2]$ is the sum of squared differences in their coordinates (see triangle PQD in Exhibit 4.2; $|\text{PQ}|^2$ denotes the squared distance between points P and Q). To denote the distance between vectors \mathbf{x} and \mathbf{y} we can use the notation $d_{\mathbf{x},\mathbf{y}}$ so that this last result can be written as:

$$d_{\mathbf{x},\mathbf{y}}^2 = (x_1 - y_1)^2 + (x_2 - y_2)^2 \quad (4.1)$$

that is, the distance itself is the square root

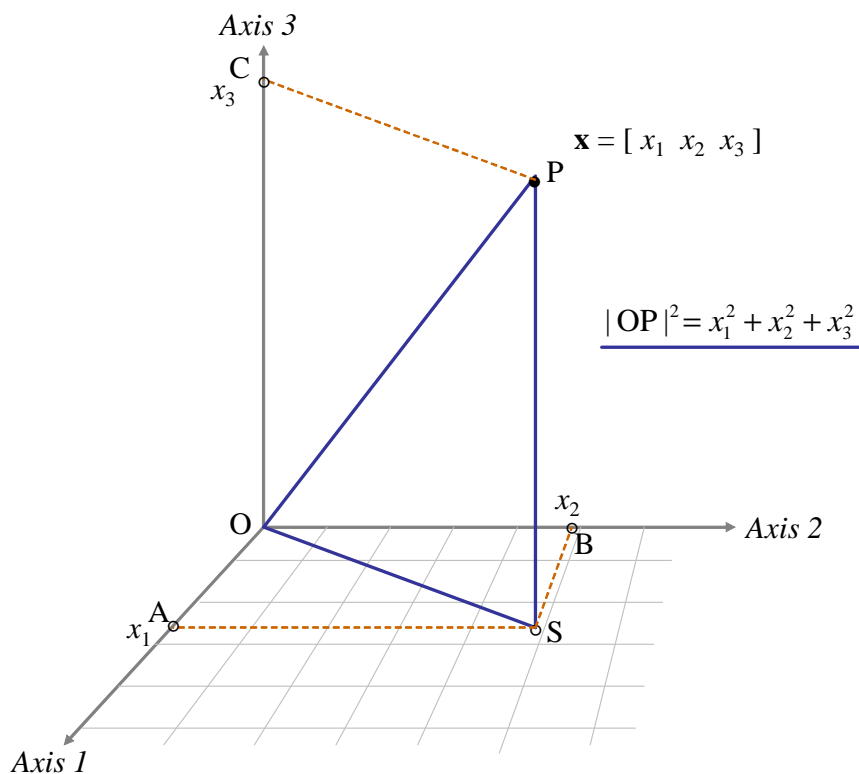
$$d_{\mathbf{x},\mathbf{y}} = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2} \quad (4.2)$$

What we called the squared length of \mathbf{x} , the distance between points P and O in Exhibit 4.2, is the distance between the vector $\mathbf{x} = [x_1 \ x_2]$ and the zero vector $\mathbf{0} = [0 \ 0]$ with coordinates all zero:

$$d_{\mathbf{x},\mathbf{0}} = \sqrt{x_1^2 + x_2^2} \quad (4.3)$$

which we could just denote by $d_{\mathbf{x}}$. The zero vector is called the *origin* of the space.

Exhibit 4.3 Pythagoras' theorem extended into three dimensional space



We move immediately to a three-dimensional point $\mathbf{x} = [x_1 \ x_2 \ x_3]$, shown in Exhibit 4.3. This figure has to be imagined in a room where the origin O is at the corner – to reinforce this idea ‘floor tiles’ have been drawn on the plane of axes 1 and 2, which is the ‘floor’ of the room. The three coordinates are at points A , B and C along the axes, and the angles AOB , AOC and COB are all 90° as well as the angle OSP at S , where the point P (depicting vector \mathbf{x}) is projected onto the ‘floor’. Using Pythagoras’ theorem twice we have:

$$|OP|^2 = |OS|^2 + |PS|^2 \quad (\text{because of right-angle at } S)$$

$$|OS|^2 = |OA|^2 + |AS|^2 \quad (\text{because of right-angle at } A)$$

and so

$$|OP|^2 = |OA|^2 + |AS|^2 + |PS|^2$$

that is, the squared length of \mathbf{x} is the sum of its three squared coordinates and so

$$d_{\mathbf{x}} = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

It is also clear that placing a point Q in Exhibit 4.3 to depict another vector \mathbf{y} and going through the motions to calculate the distance between \mathbf{x} and \mathbf{y} will lead to

$$d_{\mathbf{x},\mathbf{y}} = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2} \quad (4.4)$$

Furthermore, we can carry on like this into 4 or more dimensions, in general J dimensions, where J is the number of variables. Although we cannot draw the geometry any more, we can express the distance between two J -dimensional vectors \mathbf{x} and \mathbf{y} as:

$$d_{\mathbf{x},\mathbf{y}} = \sqrt{\sum_{j=1}^J (x_j - y_j)^2} \quad (4.5)$$

This well-known distance measure, which generalizes our notion of physical distance in two- or three-dimensional space to multidimensional space, is called the *Euclidean distance* (but often referred to as the ‘Pythagorean distance’ as well).

Standardized Euclidean distance

Let us consider measuring the distances between our 30 samples in Exhibit 1.1, using just the three continuous variables pollution, depth and temperature. What would happen if we applied formula (4.4) to measure distance between the last two samples, s_{29} and s_{30} , for example? Here is the calculation:

$$\begin{aligned} d_{s_{29},s_{30}} &= \sqrt{(6.0 - 1.9)^2 + (51 - 99)^2 + (3.0 - 2.9)^2} = \sqrt{16.81 + 2304 + 0.01} = \sqrt{2320.82} \\ &= 48.17 \end{aligned}$$

The contribution of the second variable depth to this calculation is huge – one could say that the distance is practically just the absolute difference in the depth values (equal to $|51-99| = 48$) with only tiny additional contributions from pollution and temperature. This is the problem of standardization discussed in Chapter 3 – the three variables are on completely different scales of measurement and the larger depth values have larger inter-sample differences, so they will dominate in the calculation of Euclidean distances.

Some form of standardization is necessary to balance out the contributions, and the conventional way to do this is to transform the variables so they all have the same variance of 1. At the same time we centre the variables at their means – this centring is not necessary for calculating distance, but it makes the variables all have mean zero and thus easier to compare. The transformation commonly called *standardization* is thus as follows:

$$\text{standardized value} = (\text{original value} - \text{mean}) / \text{standard deviation} \quad (4.5)$$

The means and standard deviations of the three variables are:

	<i>Pollution</i>	<i>Depth</i>	<i>Temperature</i>
mean	4.517	74.433	3.057
s.d.	2.141	15.615	0.281

leading to the table of standardized values given in Exhibit 4.4. These values are now on

Exhibit 4.4 Standardized values of the three continuous variables of Exhibit 1.1

SITE NO.	ENVIRONMENTAL VARIABLES		
	<i>Pollution</i>	<i>Depth</i>	<i>Temperature</i>
s1	0.132	-0.156	1.576
s2	-0.802	0.036	-1.979
s3	0.413	-0.988	-1.268
s4	1.720	-0.668	-0.557
s5	-0.288	-0.860	0.154
s6	-0.895	1.253	1.576
s7	0.039	-1.373	-0.557
s8	0.272	-0.860	0.865
s9	-0.288	-0.412	1.221
s10	2.561	-0.348	-0.201
s11	0.926	-1.116	0.865
s12	-0.335	0.613	0.154
s13	2.281	-1.373	-0.201
s14	0.086	0.549	-1.979
s15	1.020	1.637	-0.913
s16	-0.802	0.613	-0.201
s17	0.880	1.381	0.154
s18	-0.054	-0.028	-0.913
s19	-0.662	0.292	1.932
s20	0.506	-0.092	-0.201
s21	-0.101	-0.988	1.221
s22	-1.222	-1.309	-0.913
s23	-0.989	1.317	-0.557
s24	-0.101	-0.668	-0.201
s25	-1.175	1.445	-0.201
s26	-0.942	0.228	1.221
s27	-1.129	0.677	-0.201
s28	-0.522	1.125	0.865
s29	0.693	-1.501	-0.201
s30	-1.222	1.573	-0.557

comparable standardized scales, in units of standard deviation units with respect to the mean. For example, the value 0.693 would signify 0.693 standard deviations above the mean, and -1.222 would signify 1.222 standard deviations below the mean. The distance calculation thus aggregates squared differences in standard deviation units of each variable. As an example, the distance between the last two sites of Table is:

$$d_{s29,s30} = \sqrt{[0.693 - (-1.222)]^2 + [-1.501 - 1.573]^2 + [-0.201 - (-.557)]^2}$$

$$= \sqrt{3.667 + 9.449 + 0.127} = \sqrt{13.243} = 3.639$$

Pollution and temperature have higher contributions than before but depth still plays the largest role in this particular example, even after standardization. But this contribution is justified now, since it does show the biggest standardized difference between the samples. We call this the *standardized Euclidean distance*, meaning that it is the Euclidean distance calculated on standardized data. It will be assumed that standardization refers to the form defined by (4.5), unless specified otherwise.

We can repeat this calculation for all pairs of samples. Since the distance between sample A and sample B will be the same as between sample B and sample A, we can report these distances in a triangular matrix – Exhibit 4.5 shows part of this distance matrix, which contains a total of $\frac{1}{2} \times 30 \times 29 = 435$ distances.

Exhibit 4.5 Standardized Euclidean distances between the 30 samples, based on the three continuous environmental variables, showing part of the triangular distance matrix.

	s1	s2	s3	s4	s5	s6	...	s24	s25	s26	s27	s28	s29
s2	3.681												
s3	2.977	1.741											
s4	2.708	2.980	1.523										
s5	1.642	2.371	1.591	2.139									
s6	1.744	3.759	3.850	3.884	2.619								
s7	2.458	2.171	0.890	1.823	0.935	3.510							
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
s25	2.727	2.299	3.095	3.602	2.496	1.810	...	2.371					
s26	1.195	3.209	3.084	3.324	1.658	1.086	...	1.880	1.886				
s27	2.333	1.918	2.507	3.170	1.788	1.884	...	1.692	0.770	1.503			
s28	1.604	3.059	3.145	3.204	2.122	0.813	...	2.128	1.291	1.052	1.307		
s29	2.299	2.785	1.216	1.369	1.224	3.642	...	1.150	3.488	2.772	2.839	3.083	
s30	3.062	2.136	3.121	3.699	2.702	2.182	...	2.531	0.381	2.247	0.969	1.648	3.639

Readers might ask how all this has helped them – why convert a data table with 90 numbers to one that has 435, almost five times more? Were the histograms and scatterplots in Exhibits 1.2 and 1.4 not enough to understand these three variables? This is a good question, but we shall have to leave the answer to the Part 3 of the book, from Chapter 7 onwards, when we describe actual analyses of these distance matrices. At this early stage in the book, we can only ask readers to be patient – try to understand fully the concept of distance which will be the main thread to all the analytical methods to come.

Weighted Euclidean distance

The standardized Euclidean distance between two J -dimensional vectors can be written as:

$$d_{\mathbf{x},\mathbf{y}} = \sqrt{\sum_{j=1}^J \left(\frac{x_j}{s_j} - \frac{y_j}{s_j} \right)^2} \quad (4.6)$$

where s_j is the sample standard deviation of the j -th variable. Notice that we need not subtract the j -th mean from x_j and y_j because they will just cancel out in the differencing. Now (4.6) can be rewritten in the following equivalent way:

$$\begin{aligned} d_{\mathbf{x},\mathbf{y}} &= \sqrt{\sum_{j=1}^J \frac{1}{s_j^2} (x_j - y_j)^2} \\ &= \sqrt{\sum_{j=1}^J w_j (x_j - y_j)^2} \end{aligned} \quad (4.7)$$

where $w_j = 1/s_j^2$ is the inverse of the j -th variance. We think of w_j as a *weight* attached to the j -th variable: in other words, we compute the usual squared differences between the variables on their original scales, as we did in the (unstandardized) Euclidean distance, but then multiply these squared differences by their corresponding weights. Notice in this case how the weight of a variable with high variance is low, while the weight of a variable with low variance is high, which is another way of thinking about the compensatory effect produced by standardization. The weights of the three variables in our example are (to 4 significant figures) 0.2181, 0.004101 and 12.64 respectively, showing how much the depth variable is downweighted and the temperature variable upweighted: depth has over 3000 times the variance of temperature, so each squared difference in (4.7) is downweighted relatively by that much. We call (4.7) *weighted Euclidean distance*.

Distances for count data

So far we have looked at the distances between samples based on continuous data, now we consider distances on count data, for example the abundance data for the five taxa labelled **a**, **b**, **c**, **d** and **e** in Exhibit 1.1. First, notice that these five variables apparently do not have the problem of different measurement scales that we had for the continuous

environmental variables – all variables are counts. There are, however, different average frequencies of counts, and as we mentioned in Chapter 3, variances of count variables can be positively related to their means. The means and variances of these five variables are as follows:

	a	b	c	d	e
mean	13.47	8.73	8.40	10.90	2.97
variance	157.67	83.44	73.62	44.44	15.69

Variable **a** with the highest mean also has the highest variance, while **e** with the lowest mean has the lowest variance. Only **d** is out of line with the others, having smaller variance than **b** and **c** but a higher mean. Because this variance–mean relationship is a natural phenomenon for count variables, not one that is just particular any given example, some form of compensation of the variances needs to be performed, as before. It is not usual for count data to be standardized in the style of mean 0, variance 1, as was the case for continuous variables in (4.5). The most common ways of balancing the contributions are:

- a power transformation: usually square root $n^{1/2}$, but also double square root (i.e., fourth root $n^{1/4}$) when the variance increases faster than the mean (this situation is called ‘overdispersion’ in the literature);
- a ‘shifted log’ transformation: because of the many zeros in ecological count data, a positive number, usually 1, has to be added to the data before log-transforming; that is, $\log(1+n)$;
- chi-square distance: this is a weighted Euclidean distance of the form (4.7) which we discuss now.

The chi-square distance is special because it is at the heart of correspondence analysis, extensively used in ecological research. The first premise of this distance function is that it is calculated on relative counts, and not on the original ones, and the second is that it standardizes by the mean and not by the variance.

In our example, the count data are first converted into relative counts by dividing out the rows by their row totals so that each row contains relative proportions that add up to 1. These sets of proportions are called *profiles*, site profiles in this example – see Exhibit 4.6. The extra row at the end of Exhibit 4.6 gives the set of proportions called the *average profile*. These are the proportions calculated on the set of column totals, which are equal to 404, 262, 252, 327 and 89 respectively, with grand total 1334. Hence, $404/1334 = 0.303$, $262/1334 = 0.196$, etc. Chi-square distances are then calculated between the profiles, in a weighted Euclidean fashion, using the inverse of the average proportions as weights. Suppose c_j denotes the j -th element of the average profile, that is the abundance proportion of the j -th species in the whole data set. Then the *chi-square¹ distance*, denoted by χ , between two sites with profiles $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_J]$ and $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_J]$ is defined as:

$$\chi_{\mathbf{x},\mathbf{y}} = \sqrt{\sum_{j=1}^J \frac{1}{c_j} (x_j - y_j)^2} \quad (4.8)$$

¹ From the definition of this distance function it would have been better to call it the chi distance function, because it is not squared, as in the chi-square statistic! But the ‘chi-square’ epithet persists in the literature, so when we talk of its square we say the ‘squared chi-square distance’.

Exhibit 4.6 Profiles of the sites, obtained by dividing the rows of counts in Exhibit 1.1 by their respective row totals. The last row is the average profile, computed in the same way, as proportions of the column totals of the original table of counts.

SITE NO.	SPECIES PROPORTIONS				
	a	b	c	d	e
s1	0.000	0.074	0.333	0.519	0.074
s2	0.481	0.074	0.241	0.204	0.000
s3	0.000	0.370	0.333	0.296	0.000
s4	0.000	0.000	0.833	0.167	0.000
s5	0.342	0.132	0.079	0.263	0.184
s6	0.360	0.244	0.151	0.186	0.058
s7	0.321	0.214	0.000	0.393	0.071
s8	0.667	0.000	0.000	0.000	0.333
s9	0.315	0.130	0.185	0.259	0.111
s10	0.000	0.125	0.650	0.225	0.000
s11	0.000	0.276	0.276	0.207	0.241
s12	0.264	0.208	0.245	0.283	0.000
s13	0.000	0.000	0.760	0.000	0.240
s14	0.591	0.000	0.000	0.409	0.000
s15	0.154	0.000	0.385	0.462	0.000
s16	0.592	0.282	0.000	0.042	0.085
s17	1.000	0.000	0.000	0.000	0.000
s18	0.236	0.169	0.371	0.225	0.000
s19	0.053	0.132	0.316	0.421	0.079
s20	0.000	0.303	0.424	0.273	0.000
s21	0.444	0.000	0.000	0.222	0.333
s22	0.493	0.141	0.000	0.127	0.239
s23	0.146	0.171	0.024	0.415	0.244
s24	0.316	0.211	0.351	0.123	0.000
s25	0.395	0.321	0.000	0.284	0.000
s26	0.492	0.323	0.000	0.154	0.031
s27	0.333	0.236	0.000	0.347	0.083
s28	0.302	0.057	0.226	0.377	0.038
s29	0.423	0.000	0.269	0.308	0.000
s30	0.282	0.435	0.059	0.212	0.012
ave.	0.303	0.196	0.189	0.245	0.067

Exhibit 4.7 shows part of the 30×30 triangular matrix of chi-square distances. Once again, this is a large matrix with more numbers (435) than the original table of counts (150), and we shall see the benefit of calculating these distances from Part 3 onwards. For the moment, think of Exhibit 4.5 as a way of measuring similarities and differences between the 30 samples based on the (continuous) environmental data, while Exhibit 4.7 is the similar idea but based on the count data. But notice that the scale of distances in Exhibit 4.5 is not comparable to that of Exhibit 4.7, but the ordering of the values does have some meaning: for example, in Exhibit 4.5 the smallest standardized Euclidean distance (amongst those that we report there) is 0.381, between sites s30 and s25. In Exhibit 4.7 these two sites have one of the smallest chi-square distances as well. This means that these two sites are relatively similar in their environmental variables and also in their biological compositions. This might be something interesting, but we need to study all the pairwise distances, and not just an isolated one, in order to see if there is any connection between the biological abundances and the environmental variables (this will come later).

Exhibit 4.7 Chi-square distances between the 30 samples, based on the biological count data, showing part of the triangular distance matrix.

	s1	s2	s3	s4	s5	s6	...	s24	s25	s26	s27	s28	s29	s30
s2	1.139													
s3	0.855	1.137												
s4	1.392	1.630	1.446											
s5	1.093	0.862	1.238	2.008										
s6	1.099	0.539	0.887	1.802	0.597									
s7	1.046	0.845	1.081	2.130	0.573	0.555								
:	:	:	:	:	:	:	:							
:	:	:	:	:	:	:	:	:						
s25	1.312	0.817	1.057	2.185	0.858	0.495	...	0.917						
s26	1.508	0.805	1.224	2.241	0.834	0.475	...	0.915	0.338					
s27	1.100	0.837	1.078	2.136	0.520	0.489	...	0.983	0.412	0.562				
s28	0.681	0.504	0.954	1.572	0.724	0.613	...	0.699	0.844	0.978	0.688			
s29	0.951	0.296	1.145	1.535	0.905	0.708	...	0.662	0.956	1.021	0.897	0.340		
s30	1.330	0.986	0.846	2.101	0.970	0.535	...	0.864	0.388	0.497	0.617	1.001	1.142	

Distances for categorical data

In our introductory example we have only one categorical variable (sediment), so the question of computing distance is fairly trivial: if two samples have the same sediment then their distance is 0, and if its different it is 1. But what if there were several categorical variables, say K of them? There are several possibilities, one of the simplest being to simply extend the ‘matching’ idea and count how many matches and mismatches there are between samples, with optional averaging over variables. For example, suppose that there are five categorical variables, **C1** to **C5**, each with three categories, which we denote by $a/b/c$ and that there are two samples with the following characteristics:

	C1	C2	C3	C4	C5
sample 1	a	c	c	b	a
sample 2	b	c	b	a	a

:

Then the number of matches is 2 and the number of mismatches is 3, hence the distance between the two samples is 3 divided by 5, the number of variables, that is 0.6. This is called the *simple matching coefficient*. Sometimes this coefficient is expressed in terms of similarity, not dissimilarity, in which case it would be equal to 0.4, the relative number of matches – make sure you know which way it is being defined. Here we stick to distances, in other words dissimilarities or mismatches. Note that this coefficient is directly proportional to the squared Euclidean distance calculated between these data in dummy variable form, where each category defines a zero-one variable:

	C1a	C1b	C1c	C2a	C2b	C2c	C3a	C3b	C3c	C4a	C4b	C4c	C5a	C5b	C5c
sample 1	1	0	0	0	0	1	0	0	1	0	1	0	1	0	0
sample 2	0	1	0	0	0	1	0	1	0	1	0	0	1	0	0

The squared Euclidean distance sums the squared differences between these two vectors: if there is an agreement (there are two matches in this example) there is zero sum of squared differences, but if there is a discrepancy there are two differences, +1 and -1, which give a sum of squares of 2. So the sum of squared differences here is 6, and if this is expressed relative to the maximum discrepancy that can be achieved, namely 10 when there are no matches in the 5 variables, then this gives exactly the same value 0.6 as before.

There are several variations on the theme of the matching coefficient, and one of them is the chi-square distance for multivariate categorical data, which introduces a weighting of each category inverse to its mean value, as for profile data based on counts. Suppose that there are J categories in total (in the above example $J = 15$) and that the total occurrences of each category are denoted by n_1, \dots, n_J , with total $n = \sum_j n_j$ (since the totals for each variable equal the sample size, n will be the sample size times the number of variables). Then define c_j as follows: $c_j = n_j/n$ and use $1/c_j$ as weights in a weighted Euclidean distance between the samples coded in dummy variable form. The idea here is, as before, that the rarity of a category should count more than in the distance than a frequent category. Just like the chi-square distance function is at the heart of correspondence analysis of abundance data, so this form of the chi-square for multivariate categorical data is at the heart of *multiple correspondence analysis*. We do not treat multiple correspondence analysis specifically in this book, as it is more common in the social sciences where almost all the data are categorical, for example in questionnaire research.

SUMMARY: Measures of distance between samples: Euclidean

1. Pythagoras' theorem extends to vectors in multidimensional space: the squared length of a vector is the sum of squares of its coordinates.
2. As a consequence, squared distances between two vectors in multidimensional space are the sum of squared differences in their coordinates. This multidimensional distance is called the *Euclidean distance*, and is the natural generalization of our three-dimensional notion of physical distance to more dimensions.
3. When variables are on different measurement scales, standardization is necessary to balance the contributions of the variables in the computation of distance. The Euclidean distance computed on standardized variables is called the *standardized Euclidean distance*.
4. Standardization in the calculation of distances is equivalently thought of as *weighting* the variables – this leads to the notion of Euclidean distances with any choice of weights, called *weighted Euclidean distance*.
5. A particular weighted Euclidean distance applicable to count data is the *chi-square distance*, which is calculated between the relative counts for each sample, called *profiles*, and weights each variable by the inverse of the variable's overall mean count.