

## 3.8. DATA CLUSTERING AND VISUALIZATION

shows the default minimum spanning tree with 12 links. In Fig. 3.8.9(f) the scree plot indicates that three clusters will account for more than 95% of the data variability. The steep initial slope is a clear indication of good cluster estimation. The silhouettes are shown in Fig. 3.8.9(g–i). These were discussed in Section 3.8.5.1. In Fig. 3.8.9(j) the default parallel-coordinates plot for the same data is shown, and in Fig. 3.8.9(k) there is another view taken from the grand tour. These two plots validate the clustering and also indicate that there is no significant error introduced into the MMDS plot by truncating it into three dimensions.

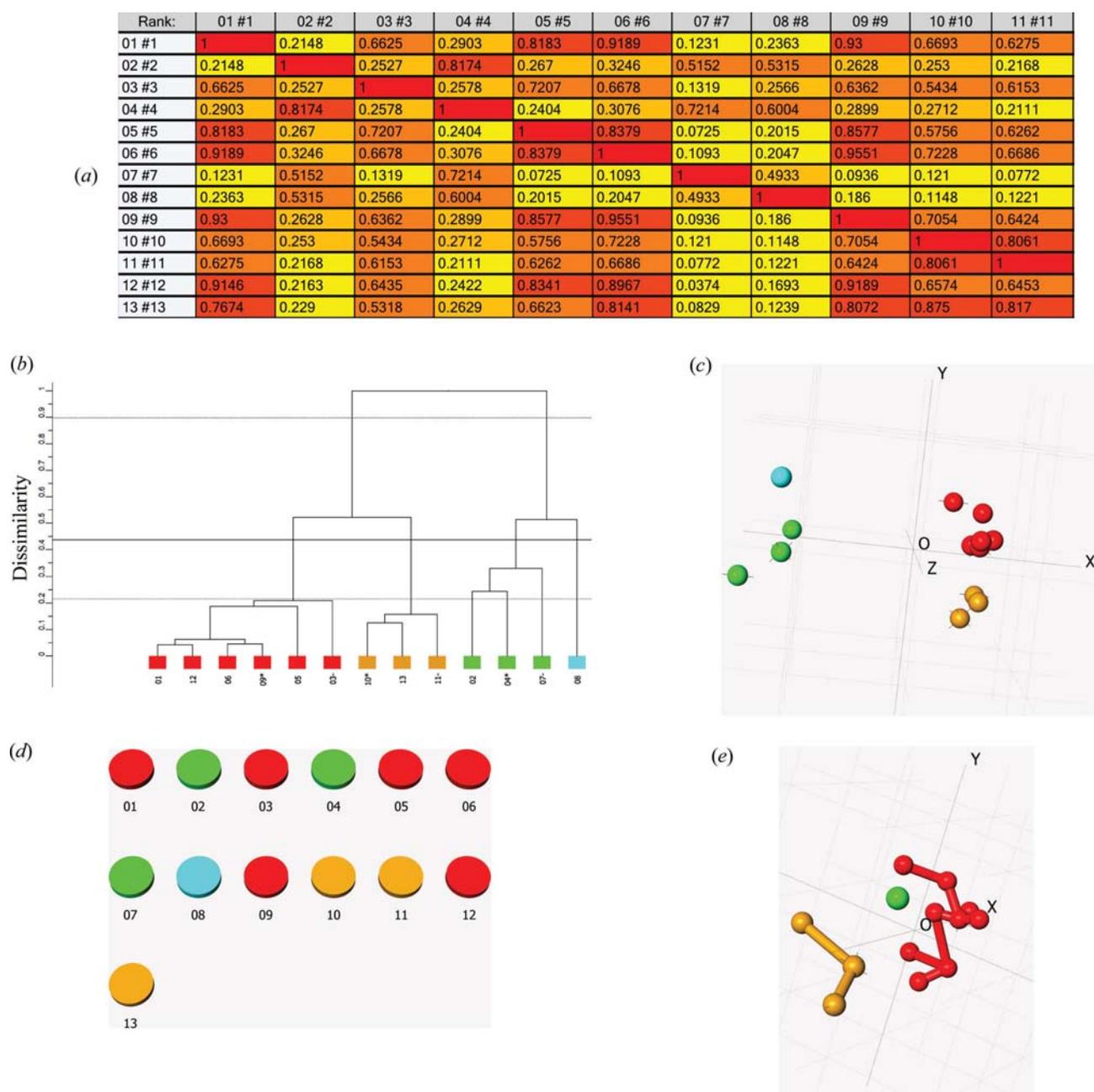
## 3.8.6.1.1. Aspirin data with amorphous samples included

As a demonstration of the handling of data from amorphous samples, five patterns for amorphous samples were included in the aspirin data and the clustering calculation was repeated. The results are shown in Fig. 3.8.10. Fig. 3.8.10(a) shows the

dendrogram. It can be seen that the amorphous samples are positioned as isolated clusters on the right-hand end. They also appear as an isolated cluster in the MMDS plot and the parallel-coordinates plots, as shown in Figs. 3.8.10(b) and (c). It could be argued that these samples should be treated as a single, five-membered cluster rather than five individuals, but we have found that this confuses the clustering algorithms, and it is clearer to the user if the data from amorphous samples are presented as separate classes.

## 3.8.6.2. Phase transitions in ammonium nitrate

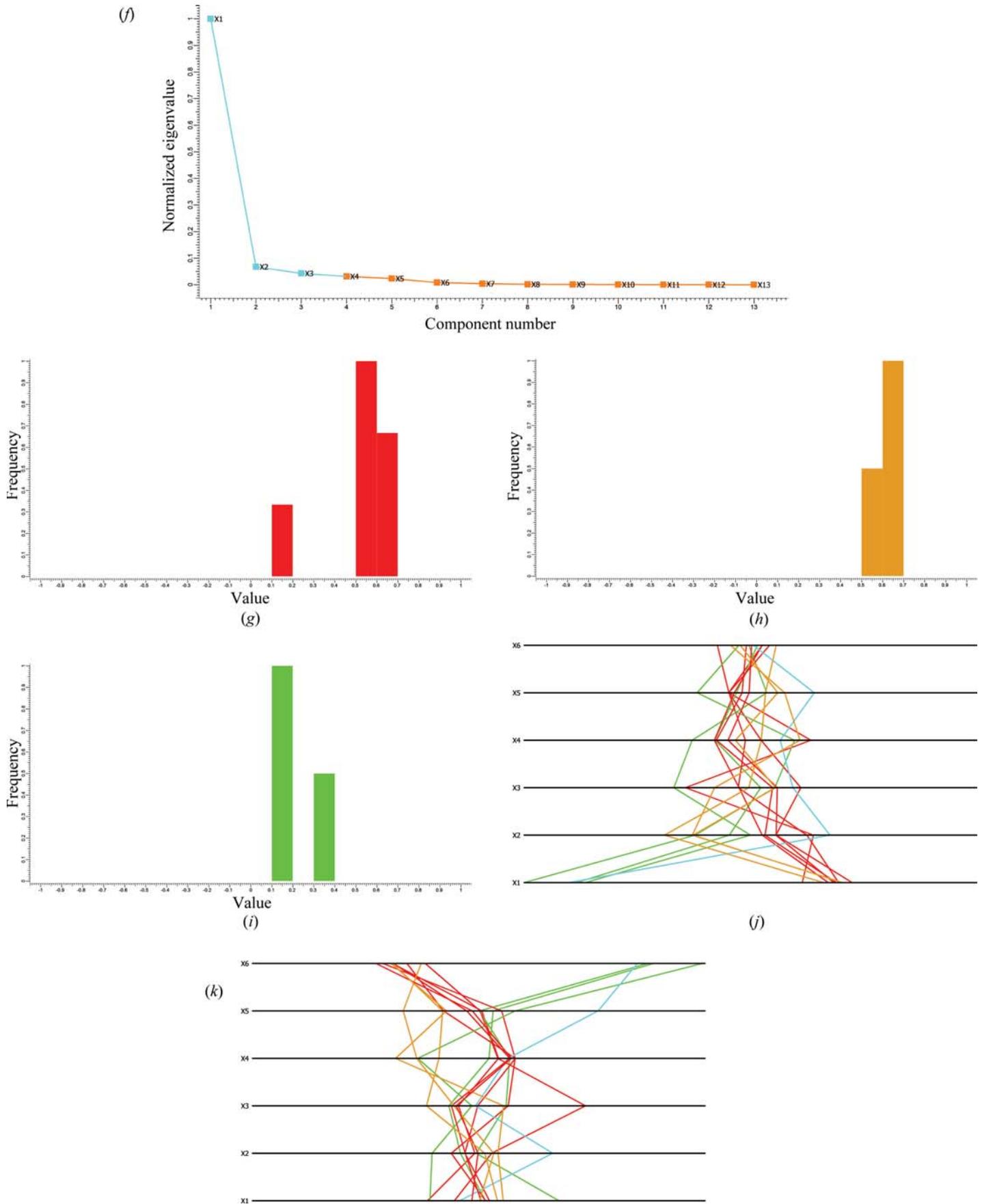
Ammonium nitrate exhibits temperature-induced phase transformations. Between 256 and 305 K it crystallizes in the orthorhombic space group  $Pmmm$  with  $a = 5.745$ ,  $b = 5.438$ ,  $c = 4.942$  Å and  $Z = 2$ ; from 305 to 357 K it crystallizes in  $Pbnm$  with



**Figure 3.8.9**

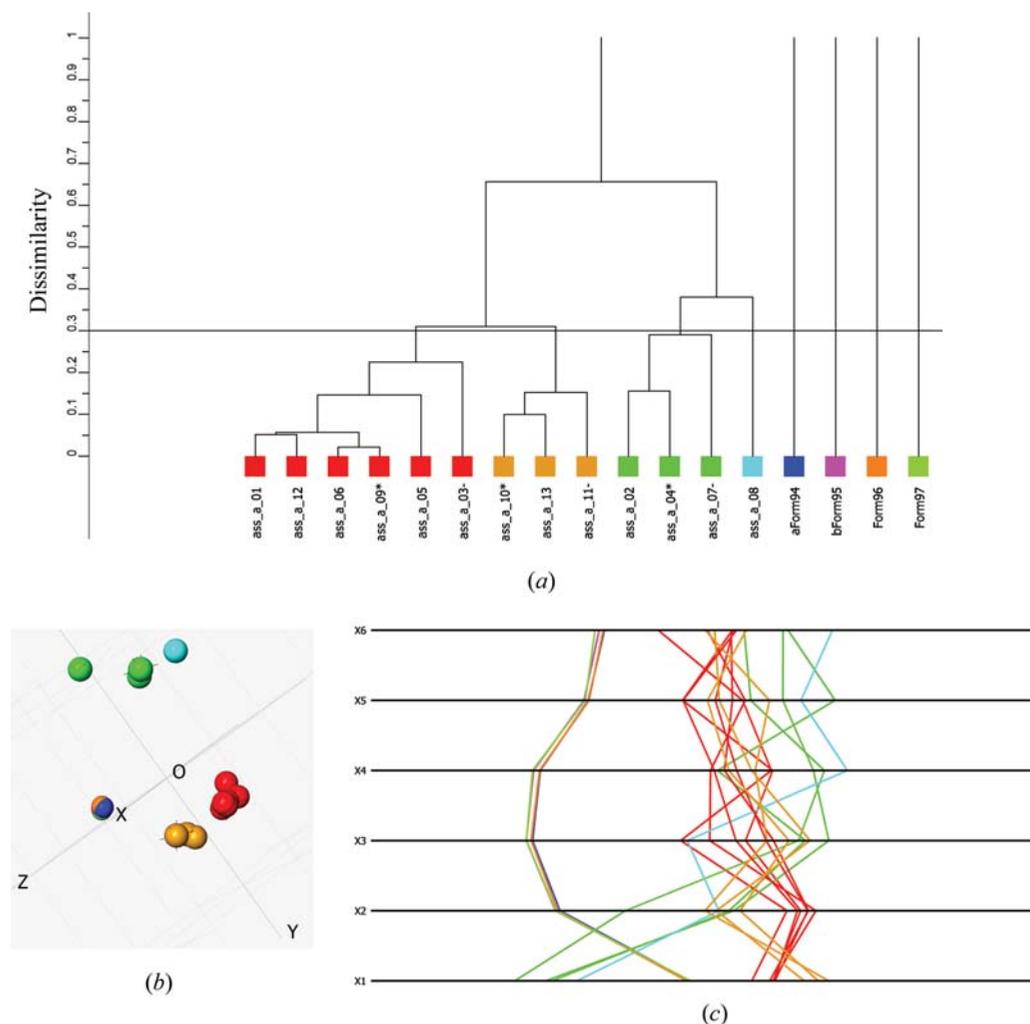
The complete cluster analysis for the aspirin samples. (a) The correlation matrix, which is the source of all the clustering results. The entries are colour coded: the darker the shade, the higher the correlation. (b) The dendrogram. The colours assigned to the samples are used in all the visualization tools. (c) The corresponding MMDS plot. The clustering defined by the dendrogram is well defined. (d) The pie-chart view. (e) The minimum spanning tree.

### 3. METHODOLOGY



**Figure 3.8.9 (continued)**

The complete cluster analysis for the aspirin samples (continued). (f) The scree plot. It indicates that three clusters explain 95% of the variance of the distance matrix derived from (a). (g–i) The silhouettes for the red, the orange and the green clusters, respectively. These are discussed in detail in the caption to Fig. 3.8.8. (j) The default parallel-coordinates plot. The clusters are well maintained into the 4th, 5th and 6th dimensions. (k) Another view of the parallel coordinates using the grand tour. The clustering remains well maintained in higher dimensions.

**Figure 3.8.10**

The aspirin data including data from five amorphous samples. (a) The resulting dendrogram and (b) the corresponding MMDS plot. (c) The parallel-coordinates plot.

$a = 7.14$ ,  $b = 7.65$ ,  $c = 5.83$  Å with  $Z = 4$ ; between 357 and 398 K it crystallizes in the tetragonal space group  $P\bar{4}2_1m$  with  $a = 5.719$ ,  $c = 4.932$  Å,  $Z = 2$ , and above 398 K it transforms to the cubic space group  $Pm\bar{3}m$  with  $a = 4.40$  Å and  $Z = 1$ . PXRD data containing 75 powder patterns taken at intervals of 3 K starting at 203 K using a D5000 Siemens diffractometer and Cu  $K\alpha$  radiation with a  $2\theta$  range of 10–100° were used (Herrmann & Engel, 1997). Fig. 3.8.11(a) shows the data in the  $2\theta$  range 17–45°.

The visualization of these data following cluster analysis is shown in Fig. 3.8.11(b) using an MMDS plot on which has been superimposed a line showing the route followed by the temperature increments. The purple line follows the transition from a mixture of forms IV and V at low temperature (red) through form IV (yellow), form II (blue) and finally form I at high temperature (green). This is an elegant and concise representation of the data in a single diagram.

### 3.8.7. Quantitative analysis with high-throughput PXRD data without Rietveld refinement

Since mixtures are so common in high-throughput experiments, and indeed in many situations with multiple data sets, it is useful to have a method of automatic quantitative analysis. The quality of data that results from high-throughput crystallography makes it unlikely that an accuracy better than 5–10% can be achieved but, nonetheless, the identification of mixtures can be carried out by whole-profile matching. First a database of  $N$  pure phases is

created, or, if that is not possible, then the most representative patterns with appropriate safeguards can be used. Assume that there is a sample pattern,  $S$ , which is considered to be a mixture of up to  $N$  components.  $S$  comprises  $m$  data points,  $S_1, S_2, \dots, S_m$ . The  $N$  patterns can be considered to make up fractions  $p_1, p_2, p_3, \dots, p_N$  of the sample pattern. The best possible combination of the database patterns to fit the sample pattern is required. A system of linear equations can be constructed in which  $x_{11}$  is measurement point 1 of pattern 1 *etc.*:

$$\begin{aligned} x_{11}p_1 + x_{12}p_2 + x_{13}p_3 + \dots + x_{1N}p_N &= S_1, \\ x_{21}p_1 + x_{22}p_2 + x_{23}p_3 + \dots + x_{2N}p_N &= S_2, \\ &\vdots \\ x_{m1}p_1 + x_{m2}p_2 + x_{m3}p_3 + \dots + x_{mN}p_N &= S_m. \end{aligned} \quad (3.8.26)$$

Writing these in matrix form, we get

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1N} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & x_{m3} & \dots & x_{mN} \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_N \end{bmatrix} = \begin{bmatrix} S_1 \\ S_2 \\ \vdots \\ S_N \end{bmatrix} \quad (3.8.27)$$

or

$$\mathbf{xp} = \mathbf{S}. \quad (3.8.28)$$

A solution for  $S$  that minimizes