

where R is a bound on the norm of the feature vectors $\phi(\mathbf{x})$ and γ is the margin obtained by the corresponding hard margin batch algorithm. This gives the following corollary.

Corollary 8.16 Fix $\delta > 0$. Suppose the batch ranking algorithm with $\nu = 1/\ell$ has margin γ on the training set

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\}$$

drawn independently at random according to a distribution \mathcal{D} and contained in a ball of radius R about the origin. Then with probability at least $1 - \delta$ over the draw of the set S , the generalisation error of the ranking function $r_{\alpha, \mathbf{b}}(\mathbf{x})$ obtained by running the on-line ranking algorithm on S in batch mode is bounded by

$$P_{\mathcal{D}}(r_{\alpha, \mathbf{b}}(\mathbf{x}) \neq y) \leq \frac{2}{\ell} \left(\frac{(|Y| - 1)(R^2 + 1)}{\gamma^2} \ln \ell + \ln \frac{\ell}{2\delta} \right),$$

provided

$$\frac{(|Y| - 1)(R^2 + 1)}{\gamma^2} \leq \frac{\ell}{2}.$$

8.2 Discovering cluster structure in a feature space

Cluster analysis aims to discover the internal organisation of a dataset by finding structure within the data in the form of 'clusters'. This generic word indicates separated groups of similar data items. Intuitively, the division into clusters should be characterised by within-cluster similarity and between-cluster (external) dissimilarity. Hence, the data is broken down into a number of groups composed of similar objects with different groups containing distinctive elements. This methodology is widely used both in multivariate statistical analysis and in machine learning.

Clustering data is useful for a number of different reasons. Firstly, it can aid our understanding of the data by breaking it into subsets that are significantly more uniform than the overall dataset. This could assist for example in understanding consumers by identifying different 'types' of behaviour that can be regarded as prototypes, perhaps forming the basis for targeted marketing exercises. It might also form the initial phase of a more complex data analysis. For example, rather than apply a classification algorithm to the full dataset, we could use a separate application for each cluster with the intention of rendering the local problem within a single cluster easier to solve accurately. In general we can view the clustering as making the data

simpler to describe, since a new data item can be specified by indicating its cluster and then its relation to the cluster centre.

Each application might suggest its own criterion for assessing the quality of the clustering obtained. Typically we would expect the quality to involve some measure of fit between a data item and the cluster to which it is assigned. This can be viewed as the pattern function of the cluster analysis. Hence, a stable clustering algorithm will give assurances about the expected value of this fit for a new randomly drawn example. As with other pattern analysis algorithms this will imply that the pattern of clusters identified in the training set is not a chance occurrence, but characterises some underlying property of the distribution generating the data.

Perhaps the most common choice for the measure assumes that each cluster has a centre and assesses the fit of a point by its squared distance from the centre of the cluster to which it is assigned. Clearly, this will be minimised if new points are assigned to the cluster whose centre is nearest. Such a division of the space creates what is known as a *Voronoi diagram* of regions each containing one of the cluster centres. The boundaries between the regions are composed of intersecting hyperplanes each defined as the set of points equidistant from some pair of cluster centres.

Throughout this section we will adopt the squared distance criterion for assessing the quality of clustering, initially based on distances in the input space, but subsequently generalised to distances in a kernel-defined feature space. In many ways the use of kernel methods for clustering is very natural, since the kernel defines pairwise similarities between data items, hence providing all the information needed to assess the quality of a clustering. Furthermore, using kernels ensures that the algorithms can be developed in full generality without specifying the particular similarity measure being used.

Ideally, all possible arrangements of the data into clusters should be tested and the best one selected. This procedure is computationally infeasible in all but very simple examples since the number of all possible partitions of a dataset grows exponentially with the number of data items. Hence, efficient algorithms need to be sought. We will present a series of algorithms that make use of the distance in a kernel-defined space as a measure of dissimilarity and use simple criteria of performance that can be used to drive practical, efficient algorithms that approximate the optimal solution.

We will start with a series of general definitions that are common to all approaches, before specifying the problem as a (non-convex) optimisation problem. We will then present a greedy algorithm to find sub-optimal solu-

tions (local minima) and a spectral algorithm that can be solved globally at the expense of relaxing the optimisation criterion.

8.2.1 Measuring cluster quality

Given an unlabelled set of data

$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\},$$

we wish to find an assignment of each point to one of a finite – but not necessarily prespecified – number N of classes. In other words, we seek a map

$$f : S \rightarrow \{1, 2, \dots, N\}.$$

This partition of the data should be chosen among all possible assignments in such a way as to solve the measure of clustering quality given in the following computation.

Computation 8.17 [Cluster quality] The clustering function should be chosen to optimise

$$f = \operatorname{argmin}_f \sum_{i,j: f_i=f(\mathbf{x}_i)=f(\mathbf{x}_j)=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2, \quad (8.6)$$

where we have as usual assumed a projection function ϕ into a feature space F , in which the kernel κ computes the inner product

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle.$$

We will use the short notation $f_i = f(\mathbf{x}_i)$ throughout this section. Figure 8.2 shows an example of a clustering of a set of data into two clusters with an indication of the contributions to (8.6). As indicated above this is not the most general clustering criterion that could be considered, but we begin by showing that it does have a number of useful properties and does subsume some apparently more general criteria. A first criticism of the criterion is that it does not seem to take into account the between-cluster separation, but only the within-cluster similarity. We might want to consider a criterion that balanced both of these two factors

$$\min_f \left\{ \sum_{i,j: f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 - \lambda \sum_{i,j: f_i \neq f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \right\}. \quad (8.7)$$

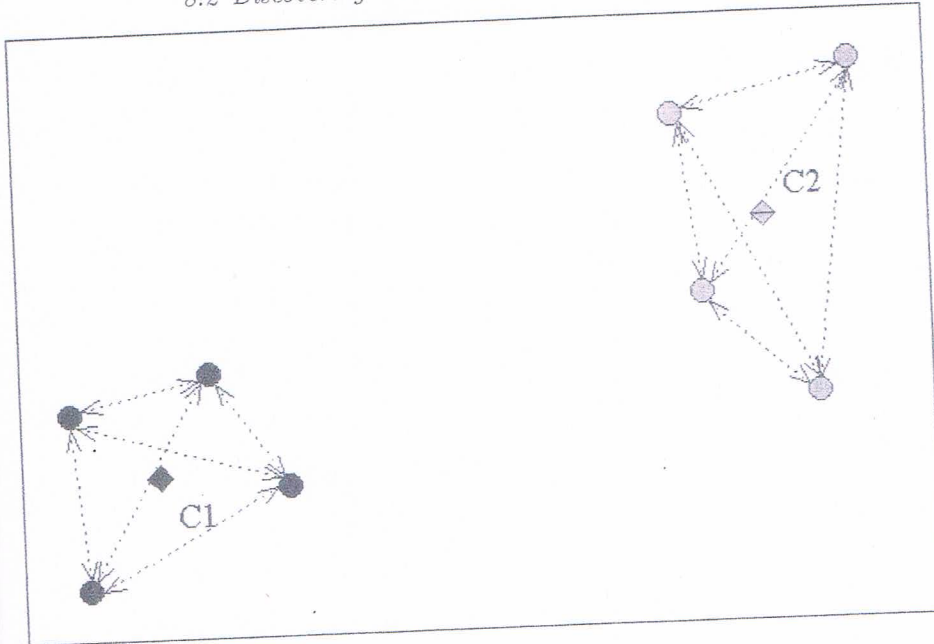


Fig. 8.2. An example of a clustering of a set of data.

However, observe that we can write

$$\begin{aligned}
 \sum_{i,j:f_i \neq f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 &= \sum_{i,j=1}^{\ell} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \\
 &\quad - \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \\
 &= A - \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2,
 \end{aligned}$$

where A is constant for a given dataset. Hence, equation (8.7) can be expressed as

$$\begin{aligned}
 &\min_f \left\{ \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 - \lambda \sum_{i,j:f_i \neq f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \right\} \\
 &= \min_f \left\{ (1 + \lambda) \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 - \lambda A \right\},
 \end{aligned}$$

showing that the same clustering function f solves the two optimisations (8.6) and (8.7). These derivations show that minimising the within-cluster

distances for a fixed number of clusters automatically maximises the between-cluster distances.

There is another nice property of the solution of the optimisation criterion (8.6). If we simply expand the expression, we obtain

$$\begin{aligned}
 \text{opt} &= \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \\
 &= \sum_{k=1}^N \sum_{i:f_i=k} \sum_{j:f_j=k} \langle \phi(\mathbf{x}_i) - \phi(\mathbf{x}_j), \phi(\mathbf{x}_i) - \phi(\mathbf{x}_j) \rangle \\
 &= \sum_{k=1}^N 2 \left(|f^{-1}(k)| \sum_{i:f_i=k} \kappa(\mathbf{x}_i, \mathbf{x}_i) - \sum_{i:f_i=k} \sum_{j:f_j=k} \kappa(\mathbf{x}_i, \mathbf{x}_j) \right) \\
 &= \sum_{k=1}^N 2 |f^{-1}(k)| \sum_{i:f_i=k} \|\phi(\mathbf{x}_i) - \mu_k\|^2,
 \end{aligned}$$

where the last line follows from (5.4) of Chapter 5 expressing the average-squared distance of a set of points from their centre of mass, and

$$\mu_k = \frac{1}{|f^{-1}(k)|} \sum_{i \in f^{-1}(k)} \phi(\mathbf{x}_i) \quad (8.8)$$

is the centre of mass of those examples assigned to cluster k , a point often referred to as the *centroid* of the cluster. This implies that the optimisation criterion (8.6) is therefore also equivalent to the criterion

$$f = \operatorname{argmin}_f \sum_{k=1}^N \left(\sum_{i:f_i=k} \|\phi(\mathbf{x}_i) - \mu_k\|^2 \right) = \operatorname{argmin}_f \sum_{i=1}^{\ell} \left\| \phi(\mathbf{x}_i) - \mu_{f(\mathbf{x}_i)} \right\|^2, \quad (8.9)$$

that seeks a clustering of points minimising the sum-squared distances to the centres of mass of the clusters. One might be tempted to assume that this implies the points are assigned to the cluster whose centroid is nearest. The following theorem shows that indeed this is the case.

Theorem 8.18 *The solution of the clustering optimisation criterion*

$$f = \operatorname{argmin}_f \sum_{i,j:f_i=f_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2$$

of Computation 8.17 can be found in the form

$$f(\mathbf{x}_i) = \operatorname{argmin}_{1 \leq k \leq N} \|\phi(\mathbf{x}_i) - \mu_k\|,$$

where μ_j is the centroid of the points assigned to cluster j .

Proof Let μ_k be as in equation (8.8). If we consider a clustering function g defined on S that assigns points to the nearest centroid

$$g(\mathbf{x}_i) = \operatorname{argmin}_{1 \leq k \leq N} \|\phi(\mathbf{x}_i) - \mu_k\|,$$

we have, by the definition of g

$$\sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \mu_{g(\mathbf{x}_i)}\|^2 \leq \sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \mu_{f(\mathbf{x}_i)}\|^2. \quad (8.10)$$

Furthermore, if we let

$$\hat{\mu}_k = \frac{1}{|g^{-1}(k)|} \sum_{i \in g^{-1}(k)} \phi(\mathbf{x}_i)$$

it follows that

$$\sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \hat{\mu}_{g(\mathbf{x}_i)}\|^2 \leq \sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \mu_{g(\mathbf{x}_i)}\|^2 \quad (8.11)$$

by Proposition 5.2. But the left-hand side is the value of the optimisation criterion (8.9) for the function g . Since f was assumed to be optimal we must have

$$\sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \hat{\mu}_{g(\mathbf{x}_i)}\|^2 \geq \sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \mu_{f(\mathbf{x}_i)}\|^2,$$

implying with (8.10) and (8.11) that the two are in fact equal. The result follows. \square

The characterisation given in Proposition 8.18 also indicates how new data should be assigned to the clusters. We simply use the natural generalisation of the assignment as

$$f(\mathbf{x}) = \operatorname{argmin}_{1 \leq k \leq N} \|\phi(\mathbf{x}) - \mu_k\|.$$

Once we have chosen the cost function of Computation 8.17 and observed that its test performance is bound solely in terms of the number of centres and the value of equation (8.6) on the training examples, it is clear that any clustering algorithm must attempt to minimise the cost function. Typically we might expect to do this for different numbers of centres, finally selecting the number for which the bound on $\mathbb{E}_{\mathcal{D}} \min_{1 \leq k \leq N} \|\phi(\mathbf{x}) - \mu_k\|^2$ is minimal.

Hence, the core task is given a fixed number of centres N find the partition into clusters which minimises equation (8.6). In view of Proposition 8.18, we therefore arrive at the following clustering optimisation strategy.

Computation 8.19 [Clustering optimisation strategy] The clustering optimisation strategy is given by

input	$S = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$, integer N
process	$\mu = \operatorname{argmin}_{\mu} \sum_{i=1}^{\ell} \min_{1 \leq k \leq N} \ \phi(\mathbf{x}_i) - \mu_k\ ^2$
output	$f(\cdot) = \operatorname{argmin}_{1 \leq k \leq N} \ \phi(\cdot) - \mu_k\ $

Figure 8.3 illustrates this strategy by showing the distances (dotted arrows) involved in computing the sum-squared criterion. The minimisation of this sum automatically maximises the indicated distance (dot-dashed arrow) between the cluster centres.

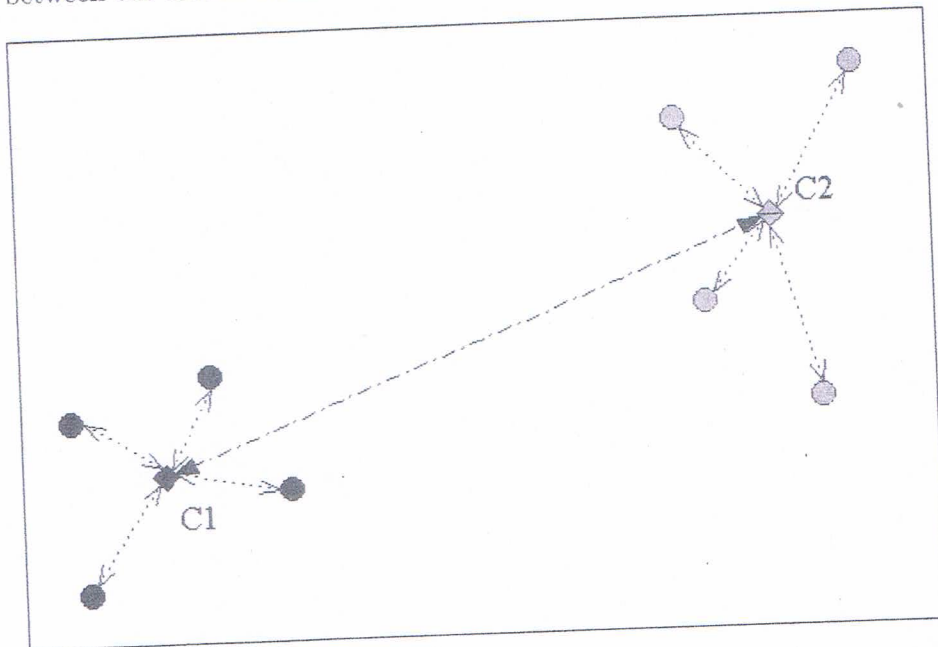


Fig. 8.3. The clustering criterion reduces to finding cluster centres to minimise sum-squared distances.

Remark 8.20 [Stability analysis] Furthermore we can see that this strategy suggests an appropriate pattern function for our stability analysis to

estimate

$$\mathbb{E}_{\mathcal{D}} \min_{1 \leq k \leq N} \|\phi(\mathbf{x}) - \mu_k\|^2;$$

the smaller the bound obtained the better the quality of the clustering achieved. ■

Unfortunately, unlike the optimisation problems that we have described previously, this problem is not convex. Indeed the task of checking if there exists a solution with value better than some threshold turns out to be NP-complete. This class of problems is generally believed not to be solvable in polynomial time, and we are therefore forced to consider heuristic or approximate algorithms that seek solutions that are close to optimal.

We will describe two such approaches in the next section. The first will use a greedy iterative method to seek a local optimum of the cost function, hence failing to be optimal precisely because of its non-convexity. The second method will consider a relaxation of the cost function to give an approximation that can be globally optimised. This is reminiscent of the approach taken to minimise the number of misclassification when applying a support vector machine to non-separable data. By introducing slack variables and using their 1-norm to upper bound the number of misclassifications we can approximately minimise this number through solving a convex optimisation problem.

The greedy method will lead to the well-known k -means algorithm, while the relaxation method gives spectral clustering algorithms. In both cases the approaches can be applied in kernel-defined feature spaces.

Remark 8.21 [Kernel matrices for well-clustered data] We have seen how we can compute distances in a kernel-defined feature space. This will provide the technology required to apply the methods in these spaces. It is often more natural to consider spaces in which unrelated objects have zero inner product. For example using a Gaussian kernel ensures that all distances between points are less than $\sqrt{2}$ with the distances becoming larger as the inputs become more orthogonal. Hence, a good clustering is achieved when the data in the same cluster are concentrated close to the prototype and the prototypes are nearly orthogonal. This means that the kernel matrix for clustered data – assuming without loss of generality that the data are sorted by cluster – will be a perturbation of a block-diagonal matrix with

one block for each cluster

$$\begin{pmatrix} B_1 & 0 & 0 & 0 \\ 0 & B_2 & 0 & 0 \\ 0 & 0 & B_3 & 0 \\ 0 & 0 & 0 & B_4 \end{pmatrix}$$

Note that this would not be the case for other distance functions where, for example, negative inner products were possible. ■

On between cluster distances There is one further property of this minimisation that relates to the means of the clusters. If we consider the covariance matrix of the data, we can perform the following derivation

$$\begin{aligned} \ell \mathbf{C} &= \sum_{i=1}^{\ell} (\phi(\mathbf{x}_i) - \phi_S) (\phi(\mathbf{x}_i) - \phi_S)' \\ &= \sum_{i=1}^{\ell} (\phi(\mathbf{x}_i) - \mu_{f_i} + \mu_{f_i} - \phi_S) (\phi(\mathbf{x}_i) - \mu_{f_i} + \mu_{f_i} - \phi_S)' \\ &= \sum_{i=1}^{\ell} (\phi(\mathbf{x}_i) - \mu_{f_i}) (\phi(\mathbf{x}_i) - \mu_{f_i})' \\ &\quad + \sum_{k=1}^N \left(\sum_{i:f_i=k} (\phi(\mathbf{x}_i) - \mu_k) \right) (\mu_k - \phi_S)' \\ &\quad + \sum_{k=1}^N (\mu_k - \phi_S) \sum_{i:f_i=k} (\phi(\mathbf{x}_i) - \mu_k)' \\ &\quad + \sum_{i=1}^{\ell} (\mu_{f_i} - \phi_S) (\mu_{f_i} - \phi_S)' \\ &= \sum_{i=1}^{\ell} (\phi(\mathbf{x}_i) - \mu_{f_i}) (\phi(\mathbf{x}_i) - \mu_{f_i})' \\ &\quad + \sum_{k=1}^N |f^{-1}(k)| (\mu_k - \phi_S) (\mu_k - \phi_S)'. \end{aligned}$$

Taking traces of both sides of the equation we obtain

$$\text{tr}(\ell \mathbf{C}) = \sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - \mu_{f_i}\|^2 + \sum_{k=1}^N |f^{-1}(k)| \|\mu_k - \phi_S\|^2.$$

The first term on the right-hand side is just the value of the Computation 8.19 that is minimised by the clustering function, while the value of the left-hand side is independent of the clustering. Hence, the clustering criterion automatically maximises the trace of the second term on the right-hand side. This corresponds to maximising

$$\sum_{k=1}^N |f^{-1}(k)| \|\mu_k - \phi_S\|^2;$$

in other words the sum of the squares of the distances from the overall mean of the cluster means weighted by their size. We again see that optimising the tightness of the clusters automatically forces their centres to be far apart.

8.2.2 Greedy solution: k -means

Proposition 8.18 confirms that we can solve Computation 8.17 by identifying centres of mass of the members of each cluster. The first algorithm we will describe attempts to do just this and is therefore referred to as the k -means algorithm. It keeps a set of cluster centroids C_1, C_2, \dots, C_N that are initialised randomly and then seeks to minimise the expression

$$\sum_{i=1}^{\ell} \|\phi(\mathbf{x}_i) - C_{f(\mathbf{x}_i)}\|^2, \quad (8.12)$$

by adapting both f as well as the centres. It will converge to a solution in which C_k is the centre of mass of the points assigned to cluster k and hence will satisfy the criterion of Proposition 8.18.

The algorithm alternates between updating f to adjust the assignment of points to clusters and updating the C_k giving the positions of the centres in a two-stage iterative procedure. The first stage simply moves points to the cluster whose cluster centre is closest. Clearly this will reduce the value of the expression in (8.12). The second stage repositions the centre of each cluster at the centre of mass of the points assigned to that cluster. We have already analysed this second stage in Proposition 5.2 showing that moving the cluster centre to the centre of mass of the points does indeed reduce the criterion of (8.12).

Hence, each stage can only reduce the expression (8.12). Since the number of possible clusterings is finite, it follows that after a finite number of iterations the algorithm will converge to a stable clustering assignment provided ties are broken in a deterministic way. If we are to implement in a dual form we must represent the clusters by an indicator matrix \mathbf{A} of dimension $\ell \times N$

containing a 1 to indicate the containment of an example in a cluster

$$\mathbf{A}_{ik} = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ is in cluster } k; \\ 0 & \text{otherwise.} \end{cases}$$

We will say that the *clustering is given* by matrix \mathbf{A} . Note that each row of \mathbf{A} contains exactly one 1, while the column sums give the number of points assigned to the different clusters. Matrices that have this form will be known as *cluster matrices*. We can therefore compute the coordinates of the centroids C_k as the N columns of the matrix

$$\mathbf{X}'\mathbf{A}\mathbf{D},$$

where \mathbf{X} contains the training example feature vectors as rows and \mathbf{D} is a diagonal $N \times N$ matrix with diagonal entries the inverse of the column sums of \mathbf{A} , indicating the number of points ascribed to that cluster. The distances of a new test vector $\phi(\mathbf{x})$ from the centroids is now given by

$$\begin{aligned} \|\phi(\mathbf{x}) - C_k\|^2 &= \|\phi(\mathbf{x})\|^2 - 2\langle\phi(\mathbf{x}), C_k\rangle + \|C_k\|^2 \\ &= \kappa(\mathbf{x}, \mathbf{x}) - 2(\mathbf{k}'\mathbf{A}\mathbf{D})_k + (\mathbf{D}\mathbf{A}'\mathbf{X}\mathbf{X}'\mathbf{A}\mathbf{D})_{kk}, \end{aligned}$$

where \mathbf{k} is the vector of inner products between $\phi(\mathbf{x})$ and the training examples. Hence, the cluster to which $\phi(\mathbf{x})$ should be assigned is given by

$$\operatorname{argmin}_{1 \leq k \leq N} \|\phi(\mathbf{x}) - C_k\|^2 = \operatorname{argmin}_{1 \leq k \leq N} (\mathbf{D}\mathbf{A}'\mathbf{K}\mathbf{A}\mathbf{D})_{kk} - 2(\mathbf{k}'\mathbf{A}\mathbf{D})_k,$$

where \mathbf{K} is the kernel matrix of the training set. This provides the rule for classifying new data. The update rule consists in reassigning the entries in the matrix \mathbf{A} according to the same rule in order to redefine the clusters.

Algorithm 8.22 [Kernel k -means] Matlab code for the kernel k -means algorithm is given in Code Fragment 8.3. ■

Despite its popularity, this algorithm is prone to local minima since the optimisation is not convex. Considerable effort has been devoted to finding good initial guesses or inserting additional constraints in order to limit the effect of this fact on the quality of the solution obtained. In the next section we see two relaxations of the original problem for which we can find the global solution.

8.2.3 Relaxed solution: spectral methods

In this subsection, rather than relying on gradient descent methods to tackle a non-convex problem, we make a convex relaxation of the problem in order

```

% original kernel matrix stored in variable K
% clustering given by a ell x N binary matrix A
% and cluster allocation function f
% d gives the distances to cluster centroids
A = zeros(ell,N);
f = ceil(rand(ell,1)* N);
for i=1,ell
    A(i,f(i)) = 1;
end
change = 1;
while change = 1
    change = 0;
    E = A * diag(1./sum(A));
    Z = ones(ell,1)* diag(E'*K*E) - 2*K*E;
    [d, ff] = min(Z, [], 2);
    for i=1,ell
        if f(i) ~= ff(i)
            A(i,ff(i)) = 1;
            A(i, f(i)) = 0;
            change = 1;
        end
    end
    f = ff;
end

```

Code Fragment 8.3. Matlab code to perform k -means clustering.

to obtain a closed form approximation. We can then study the approximation and statistical properties of its solutions.

Clustering into two classes We first consider the simpler case when there are just two clusters. In this relaxation we represent the cluster assignment by a vector $\mathbf{y} \in \{-1, +1\}^\ell$, that associates to each point a $\{-1, +1\}$ label. For the two classes case the clustering quality criterion described above is minimised by maximising

$$\sum_{y_i \neq y_j} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2.$$

Assuming that the data is normalised and the sizes of the clusters are equal this will correspond to minimising the so-called cut cost

$$2 \sum_{y_i \neq y_j} \kappa(\mathbf{x}_i, \mathbf{x}_j) = \sum_{i,j=1}^{\ell} \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i,j=1}^{\ell} y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j),$$

subject to $\mathbf{y} \in \{-1, +1\}^\ell$,

since it measures the kernel 'weight' between vertices in different clusters. Hence, we must solve

$$\begin{array}{ll} \max & \mathbf{y}'\mathbf{K}\mathbf{y} \\ \text{subject to} & \mathbf{y} \in \{-1, +1\}^\ell. \end{array}$$

We can relax this optimisation by removing the restriction that \mathbf{y} be a binary vector while controlling its norm. This is achieved by maximising the Raleigh quotient (3.2)

$$\max \frac{\mathbf{y}'\mathbf{K}\mathbf{y}}{\mathbf{y}'\mathbf{y}}.$$

As observed in Chapter 3 this is solved by the eigenvector of the matrix \mathbf{K} corresponding to the largest eigenvalue with the value of the quotient equal to the eigenvalue λ_1 . Hence, we obtain a lower bound on the cut cost of

$$0.5 \left(\sum_{i,j=1}^{\ell} \kappa(\mathbf{x}_i, \mathbf{x}_j) - \lambda_1 \right),$$

giving a corresponding lower bound on the value of the sum-squared criterion. Though such a lower bound is useful, the question remains as to whether the approach can suggest useful clusterings of the data. A very natural way to do so in this two-cluster case is simply to threshold the vector \mathbf{y} hence converting it to a binary clustering vector. This naive approach can deliver surprisingly good results though there is no a priori guarantee attached to the quality of the solution.

Remark 8.23 [Alternative criterion] It is also possible to consider minimising a ratio between the cut size and a measure of the size of the clusters. This leads through similar relaxations to different eigenvalue problems. For example if we let \mathbf{D} be the diagonal matrix with entries

$$D_{ii} = \sum_{j=1}^{\ell} K_{ij},$$

then useful partitions can be derived from the eigenvectors of

$$\mathbf{D}^{-1}\mathbf{K}, \mathbf{D}^{-1/2}\mathbf{K}\mathbf{D}^{-1/2} \text{ and } \mathbf{K} - \mathbf{D}$$

with varying justifications. In all cases thresholding the resulting vectors delivers the corresponding partitions. Generally the approach is motivated using the Gaussian kernel with its useful properties discussed above. ■

Multiclass clustering We now consider the more general problem of multiclass clustering. We start with an equivalent formulation of the sum-of-squares minimization problem as a trace maximization under special constraints. By successively relaxing these constraints we are led to an approximate algorithm with nice global properties.

Consider the derivation and notation introduced to obtain the program for k -means clustering. We can compute the coordinates of the centroids C_k as the N columns of the matrix

$$\mathbf{X}'\mathbf{A}\mathbf{D},$$

where \mathbf{X} is the data matrix, \mathbf{A} a matrix assigning points to clusters and \mathbf{D} a diagonal matrix with inverses of the cluster sizes on the diagonal. Consider now the matrix

$$\mathbf{X}'\mathbf{A}\mathbf{D}\mathbf{A}'.$$

It has ℓ columns with the i th column a copy of the cluster centroid corresponding to the i th example. Hence, we can compute the sum-squares of the distances from the examples to their corresponding cluster centroid as

$$\begin{aligned} \|\mathbf{X}'\mathbf{A}\mathbf{D}\mathbf{A}' - \mathbf{X}'\|_F^2 &= \|(\mathbf{I}_\ell - \mathbf{A}\mathbf{D}\mathbf{A}')\mathbf{X}\|_F^2 \\ &= \text{tr}(\mathbf{X}'(\mathbf{I}_\ell - \mathbf{A}\mathbf{D}\mathbf{A}')\mathbf{X}) \\ &= \text{tr}(\mathbf{X}\mathbf{X}') - \text{tr}(\sqrt{\mathbf{D}}\mathbf{A}'\mathbf{X}\mathbf{X}'\mathbf{A}\sqrt{\mathbf{D}}), \end{aligned}$$

since

$$(\mathbf{I}_\ell - \mathbf{A}\mathbf{D}\mathbf{A}')^2 = (\mathbf{I}_\ell - \mathbf{A}\mathbf{D}\mathbf{A}')$$

as $\mathbf{A}'\mathbf{A}\mathbf{D} = \sqrt{\mathbf{D}}\mathbf{A}'\mathbf{A}\sqrt{\mathbf{D}} = \mathbf{I}_N$, indicating that $(\mathbf{I}_\ell - \mathbf{A}\mathbf{D}\mathbf{A}')$ is a projection matrix. We have therefore shown the following proposition.

Proposition 8.24 *The sum-squares cost function $ss(\mathbf{A})$ of a clustering, given by matrix \mathbf{A} can be expressed as*

$$ss(\mathbf{A}) = \text{tr}(\mathbf{K}) - \text{tr}(\sqrt{\mathbf{D}}\mathbf{A}'\mathbf{K}\mathbf{A}\sqrt{\mathbf{D}}),$$

where \mathbf{K} is the kernel matrix and $\mathbf{D} = \mathbf{D}(\mathbf{A})$ is the diagonal matrix with the inverses of the column sums of \mathbf{A} .

Since the first term is not affected by the choice of clustering, the proposition leads to the Computation.

Computation 8.25 [Multiclass clustering] We can minimise the cost $ss(\mathbf{A})$ by solving

$$\begin{array}{ll} \max_{\mathbf{A}} & \text{tr}(\sqrt{\mathbf{D}}\mathbf{A}'\mathbf{K}\mathbf{A}\sqrt{\mathbf{D}}) \\ \text{subject to} & \mathbf{A} \text{ is a cluster matrix and } \mathbf{D} = \mathbf{D}(\mathbf{A}). \end{array}$$

Remark 8.26 [Distances of the cluster centres from the origin] We can see that $\text{tr}(\sqrt{\mathbf{D}}\mathbf{A}'\mathbf{K}\mathbf{A}\sqrt{\mathbf{D}})$ is the sum of the squares of the cluster centroids, relating back to our previous observations about the sum-square criterion corresponding to maximising the sum-squared distances of the centroids from the overall centre of mass of the data. This shows that this holds for the origin as well and since an optimal clustering is invariant to translations of the coordinate axes, this will be true of the sum-squared distances to any fixed point. ■

We have now arrived at a constrained optimisation problem whose solution solves to the min-squared clustering criterion. Our aim now is to relax the constraints to obtain a problem that can be optimised exactly, but whose solution will not correspond precisely to a clustering. Note that the matrices \mathbf{A} and $\mathbf{D} = \mathbf{D}(\mathbf{A})$ satisfy

$$\sqrt{\mathbf{D}}\mathbf{A}'\mathbf{A}\sqrt{\mathbf{D}} = \mathbf{I}_N = \mathbf{H}'\mathbf{H}, \quad (8.13)$$

where $\mathbf{H} = \mathbf{A}\sqrt{\mathbf{D}}$. Hence, only requiring that this $\ell \times N$ matrix satisfies equation (8.13), we obtain the relaxed maximization problem given in the computation.

Computation 8.27 [Relaxed multiclass clustering] The relaxed maximisation for multiclass clustering is obtained by solving

$$\begin{array}{ll} \max & \text{tr}(\mathbf{H}'\mathbf{K}\mathbf{H}) \\ \text{subject to} & \mathbf{H}'\mathbf{H} = \mathbf{I}_N, \end{array}$$

The solutions of Computation 8.27 will not correspond to clusterings, but may lead to good clusterings. The important property of the relaxation is that it can be solved in closed-form.

Proposition 8.28 The maximum of the trace $\text{tr}(\mathbf{H}'\mathbf{K}\mathbf{H})$ over all $\ell \times N$ matrices satisfying $\mathbf{H}'\mathbf{H} = \mathbf{I}_N$ is equal to the sum of the first N eigenvalues

of \mathbf{K}

$$\max_{\mathbf{H}'\mathbf{H}=\mathbf{I}_N} \text{tr}(\mathbf{H}'\mathbf{K}\mathbf{H}) = \sum_{k=1}^N \lambda_k$$

while the \mathbf{H}^* realising the optimum is given by

$$\mathbf{V}_N \mathbf{Q},$$

where \mathbf{Q} is an arbitrary $N \times N$ orthonormal matrix and \mathbf{V}_N is the $\ell \times N$ matrix composed of the first N eigenvectors of \mathbf{K} . Furthermore, we can lower bound the sum-squared error of the best clustering by

$$\begin{aligned} \min_{\mathbf{A} \text{ clustering matrix}} \text{ss}(\mathbf{A}) &= \text{tr}(\mathbf{K}) - \max_{\mathbf{A} \text{ clustering matrix}} \text{tr}(\sqrt{\mathbf{D}} \mathbf{A}' \mathbf{K} \mathbf{A} \sqrt{\mathbf{D}}) \\ &\geq \text{tr}(\mathbf{K}) - \max_{\mathbf{H}'\mathbf{H}=\mathbf{I}_N} \text{tr}(\mathbf{H}'\mathbf{K}\mathbf{H}) = \sum_{k=N+1}^{\ell} \lambda_k. \end{aligned}$$

Proof Since $\mathbf{H}'\mathbf{H} = \mathbf{I}_N$ the operator $\mathbf{P} = \mathbf{H}\mathbf{H}'$ is a rank N projection. This follows from the fact that $(\mathbf{H}\mathbf{H}')^2 = \mathbf{H}\mathbf{H}'\mathbf{H}\mathbf{H}' = \mathbf{H}\mathbf{I}_N\mathbf{H}' = \mathbf{H}\mathbf{H}'$ and $\text{rank } \mathbf{H} = \text{rank } \mathbf{H}'\mathbf{H} = \text{rank } \mathbf{I}_N = N$. Therefore

$$\mathbf{I}_N \mathbf{H}'\mathbf{K}\mathbf{H} \mathbf{I}_N = \mathbf{H}'\mathbf{H}\mathbf{H}'\mathbf{X}\mathbf{X}'\mathbf{H}\mathbf{H}'\mathbf{H} = \|\mathbf{H}'\mathbf{P}\mathbf{X}\|_F^2 = \|\mathbf{P}\mathbf{X}\|_F^2.$$

Hence, we seek the N -dimensional projection of the columns of \mathbf{X} that maximises the resulting sum of the squared norms. Treating these columns as the training vectors and viewing maximising the projection as minimising the residual we can apply Proposition 6.12. It follows that the maximum will be realised by the eigensubspace spanned by the N eigenvectors corresponding to the largest eigenvalues for the matrix $\mathbf{X}\mathbf{X}' = \mathbf{K}$. Clearly, the projection is only specified up to an arbitrary $N \times N$ orthonormal transformation \mathbf{Q} . \square

At first sight we have not gained a great deal by moving to the relaxed version of the problem. It is true that we have obtained a strict lower bound on the quality of clustering that can be achieved, but the matrix \mathbf{H} that realises that lower bound does not in itself immediately suggest a method of performing a clustering that comes close to the bound. Furthermore, in the case of multi-class clustering we do not have an obvious simple thresholding algorithm for converting the result of the eigenvector analysis into a clustering as we found in the two cluster examples. We mention three possible approaches.

Re-clustering One approach is to apply a different clustering algorithm in the reduced N -dimensional representation of the data, when we map each example to the corresponding row of the matrix \mathbf{V}_N , possibly after performing a renormalisation.

Eigenvector approach We will describe a different method that is related to the proof of Proposition 8.28. Consider the choice $\mathbf{H}^* = \mathbf{V}_N$ that realises the optimum bound of the proposition. Let $\mathbf{W} = \mathbf{V}_N \sqrt{\mathbf{\Lambda}_N}$ be obtained from \mathbf{V}_N by multiplying column i by $\sqrt{\lambda_i}$, $i = 1, \dots, N$. We now form the cluster matrix \mathbf{A} by setting the largest entry in each row of \mathbf{W} to 1 and the remaining entries to 0.

QR approach An alternative approach is inspired by a desire to construct an approximate cluster matrix which is related to \mathbf{V}_N by an orthonormal transformation

$$\mathbf{A} \approx \mathbf{V}_N \mathbf{Q},$$

implying that

$$\mathbf{V}'_N \approx \mathbf{Q} \mathbf{A}'.$$

If we perform a QR decomposition of \mathbf{V}'_N we obtain

$$\mathbf{V}'_N = \hat{\mathbf{Q}} \mathbf{R}$$

with $\hat{\mathbf{Q}}$ an $N \times N$ orthogonal matrix and \mathbf{R} an $N \times \ell$ upper triangular. By assigning vector i to the cluster index by the row with largest entry in the column i of matrix \mathbf{R} , we obtain a cluster matrix $\mathbf{A}' \approx \mathbf{R}$, hence giving a value of $\text{ss}(\mathbf{A})$ close to that given by the bound.

8.3 Data visualisation

Visualisation refers to techniques that can present a dataset

$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$$

in a manner that reveals some underlying structure of the data in a way that is easily understood or appreciated by a user. A clustering is one type of structure that can be visualised. If for example there is a natural clustering of the data into say four clusters, each one grouped tightly around a separate centroid, our understanding of the data is greatly enhanced by displaying the four centroids and indicating the tightness of the clustering around them. The centroids can be thought of as prototypical data items.

Hence, for example in a market analysis, customers might be clustered into a set of typical types with individual customers assigned to the type to which they are most similar.

In this section we will consider a different type of visualisation. Our aim is to provide a two- or three-dimensional 'mapping' of the data, hence displaying the data points as dots on a page or as points in a three-dimensional image. This type of visualisation is of great importance in many data mining tasks, but it assumes a special role in kernel methods, where we typically embed the data into a high-dimensional vector space. We have already considered measures for assessing the quality of an embedding such as the classifier margin, correlation with output values and so on. We will also in later chapters be looking into procedures for transforming prior models into embeddings. However once we arrive at the particular embedding, it is also important to have ways of visually displaying the data in the chosen feature space. Looking at the data helps us get a 'feel' for the structure of the data, hence suggesting why certain points are outliers, or what type of relations can be found. This can in turn help us pick the best algorithm out of the toolbox of methods we have been assembling since Chapter 5. In other words being able to 'see' the relative positions of the data in the feature space plays an important role in guiding the intuition of the data analyst.

Using the first few principal components, as computed by the PCA algorithm of Chapter 6, is a well-known method of visualisation forming the core of the classical multidimensional scaling algorithm. As already demonstrated in Proposition 6.12 PCA minimises the sum-squared norms of the residuals between the subspace representation and the actual data.

Naturally if one wants to look at the data from an angle that emphasises a certain property, such as a given dichotomy, other projections can be better suited, for example using the first two partial least squares features. In this section we will assume that the feature space has already been adapted to best capture the view of the data we are concerned with but typically using a high-dimensional kernel representation. Our main concern will therefore be to develop algorithms that can find low-dimensional representations of high-dimensional data.

Multidimensional scaling This problem has received considerable attention in multivariate statistics under the heading of multidimensional scaling (MDS). This is a series of techniques directly aimed at finding optimal low-dimensional embeddings of data mostly for visualisation purposes. The starting point for MDS is traditionally a matrix of distances or similarities

rather than a Gram matrix of inner products or even a Euclidean embedding. Indeed the first stages of the MDS process aim to convert the matrix of similarities into a matrix of inner products. For metric MDS it is assumed that the distances correspond to embeddings in a Euclidean space, while for non-metric MDS these similarities can be measured in any way. Once an approximate inner product matrix has been formed, classical MDS then uses the first two or three eigenvectors of the eigen-decomposition of the resulting Gram matrix to define two- or three-dimensional projections of the points for visualisation. Hence, if we make use of a kernel-defined feature space the first stages are no longer required and MDS reduces to computing the first two or three kernel PCA projections.

Algorithm 8.29 [MDS for kernel-embedded data] The MDS algorithm for data in a kernel-defined feature space is as follows:

input	Data $S = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$, dimension $k = 2, 3$.
process	$\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$, $i, j = 1, \dots, \ell$ $\mathbf{K} - \frac{1}{\ell} \mathbf{j} \mathbf{j}' \mathbf{K} - \frac{1}{\ell} \mathbf{K} \mathbf{j} \mathbf{j}' + \frac{1}{\ell^2} (\mathbf{j}' \mathbf{K} \mathbf{j}) \mathbf{j} \mathbf{j}'$, $[\mathbf{V}, \mathbf{\Lambda}] = \text{eig}(\mathbf{K})$ $\alpha^j = \frac{1}{\sqrt{\lambda_j}} \mathbf{v}_j$, $j = 1, \dots, k$. $\tilde{\mathbf{x}}_i = \left(\sum_{j=1}^k \alpha_j^j \kappa(\mathbf{x}_i, \mathbf{x}) \right)_{j=1}^k$
output	Display transformed data $\tilde{S} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_\ell\}$.

Visualisation quality We will consider a further method of visualisation strongly related to MDS, but which is motivated by different criteria for assessing the quality of the representation of the data.

We can define the problem of visualisation as follows. Given a set of points

$$S = \{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_\ell)\}$$

in a kernel-defined feature space F with

$$\phi : X \longrightarrow F,$$

find a projection τ from X into \mathbb{R}^k , for small k such that

$$\|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\| \approx \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|, \text{ for } i, j = 1, \dots, \ell.$$

We will use τ_s to denote the projection onto the s th component of τ and with a slight abuse of notation, as a vector of these projection values indexed by the training examples. As mentioned above it follows from Proposition 6.12

that the embedding determined by kernel PCA minimises the sum-squared residuals

$$\sum_{i=1}^{\ell} \|\tau(\mathbf{x}_i) - \phi(\mathbf{x}_i)\|^2,$$

where we make τ an embedding into a k -dimensional subspace of the feature space F . Our next method aims to control more directly the relationship between the original and projection distances by solving the following computation.

Computation 8.30 [Visualisation quality] The quality of a visualisation can be optimised as follows

$$\begin{aligned} \min_{\tau} E(\tau) &= \sum_{i,j=1}^{\ell} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\|^2 \\ &= \sum_{i,j=1}^{\ell} \kappa(\mathbf{x}_i, \mathbf{x}_j) \|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\|^2, \\ \text{subject to } \|\tau_s\| &= 1, \quad \tau_s \perp \mathbf{j}, \quad s = 1, \dots, k, \\ \text{and } \tau_s &\perp \tau_t, \quad s, t = 1, \dots, k. \end{aligned} \tag{8.14}$$

Observe that it follows from the constraints that

$$\begin{aligned} \sum_{i,j=1}^{\ell} \|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\|^2 &= \sum_{i,j=1}^{\ell} \sum_{s=1}^k (\tau_s(\mathbf{x}_i) - \tau_s(\mathbf{x}_j))^2 \\ &= \sum_{s=1}^k \sum_{i,j=1}^{\ell} (\tau_s(\mathbf{x}_i) - \tau_s(\mathbf{x}_j))^2 \\ &= 2 \sum_{s=1}^k \left(\ell \sum_{i=1}^{\ell} \tau_s(\mathbf{x}_i)^2 - \sum_{i,j=1}^{\ell} \tau_s(\mathbf{x}_i) \tau_s(\mathbf{x}_j) \right) \\ &= 2\ell k - 2 \sum_{s=1}^k \sum_{i=1}^{\ell} \tau_s(\mathbf{x}_i) \sum_{j=1}^{\ell} \tau_s(\mathbf{x}_j) = 2\ell k. \end{aligned}$$

It therefore follows that, if the data is normalised, solving Computation 8.30 corresponds to minimising

$$E(\tau) = \sum_{i,j=1}^{\ell} \left(1 - 0.5 \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \right) \|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\|^2$$

$$= 2\ell k - \sum_{i,j=1}^{\ell} 0.5 \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \|\tau(\mathbf{x}_i) - \tau(\mathbf{x}_j)\|^2,$$

hence optimising the correlation between the original and projected squared distances. More generally we can see minimisation as aiming to put large distances between points with small inner products and small distances between points having large inner products. The constraints ensure equal scaling in all dimensions centred around the origin, while the different dimensions are required to be mutually orthogonal to ensure no structure is unnecessarily reproduced.

Our next theorem will characterise the solution of the above optimisation using the eigenvectors of the so-called Laplacian matrix. This matrix can also be used in clustering as it frequently possesses more balanced properties than the kernel matrix. It is defined as follows.

Definition 8.31 [Laplacian matrix] The *Laplacian matrix* $\mathbf{L}(\mathbf{K})$ of a kernel matrix \mathbf{K} is defined by

$$\mathbf{L}(\mathbf{K}) = \mathbf{D} - \mathbf{K},$$

where \mathbf{D} is the diagonal matrix with entries

$$D_{ii} = \sum_{j=1}^{\ell} K_{ij}.$$

■

Observe the following simple property of the Laplacian matrix. Given any real vector $\mathbf{v} = (v_1, \dots, v_{\ell}) \in \mathbb{R}^{\ell}$

$$\begin{aligned} \sum_{i,j=1}^{\ell} K_{ij} (v_i - v_j)^2 &= 2 \sum_{i,j=1}^{\ell} K_{ij} v_i^2 - 2 \sum_{i,j=1}^{\ell} K_{ij} v_i v_j \\ &= 2\mathbf{v}'\mathbf{D}\mathbf{v} - 2\mathbf{v}'\mathbf{K}\mathbf{v} = 2\mathbf{v}'\mathbf{L}(\mathbf{K})\mathbf{v}. \end{aligned}$$

It follows that the all 1s vector \mathbf{j} is an eigenvector of $\mathbf{L}(\mathbf{K})$ with eigenvalue 0 since the sum is zero if $v_i = v_j = 1$. It also implies that if the kernel matrix has positive entries then $\mathbf{L}(\mathbf{K})$ is positive semi-definite. In the statement of the following theorem we separate out λ_1 as the eigenvalue 0, while ordering the remaining eigenvalues in ascending order.

Theorem 8.32 *Let*

$$S = \{\mathbf{x}_1, \dots, \mathbf{x}_{\ell}\}$$

be a set of points with kernel matrix \mathbf{K} . The visualisation problem given in Computation 8.30 is solved by computing the eigenvectors $\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^\ell$ with corresponding eigenvalues $0 = \lambda_1, \lambda_2 \leq \dots \leq \lambda_\ell$ of the Laplacian matrix $\mathbf{L}(\mathbf{K})$. An optimal embedding $\boldsymbol{\tau}$ is given by $\tau_i = \mathbf{v}^{i+1}$, $i = 1, \dots, k$ and the minimal value of $E(\boldsymbol{\tau})$ is

$$2 \sum_{\ell=2}^{k+1} \lambda_\ell.$$

If $\lambda_{k+1} < \lambda_{k+2}$ then the optimal embedding is unique up to orthonormal transformations in \mathbb{R}^k .

Proof The criterion to be minimised is

$$\begin{aligned} \sum_{i,j=1}^{\ell} \kappa(\mathbf{x}_i, \mathbf{x}_j) \|\boldsymbol{\tau}(\mathbf{x}_i) - \boldsymbol{\tau}(\mathbf{x}_j)\|^2 &= \sum_{s=1}^k \sum_{i,j=1}^{\ell} \kappa(\mathbf{x}_i, \mathbf{x}_j) (\tau_s(\mathbf{x}_i) - \tau_s(\mathbf{x}_j))^2 \\ &= 2 \sum_{s=1}^k \boldsymbol{\tau}_s' \mathbf{L}(\mathbf{K}) \boldsymbol{\tau}_s. \end{aligned}$$

Taking into account the normalisation and orthogonality constraints gives the solution as the eigenvectors of $\mathbf{L}(\mathbf{K})$ by the usual characterisation of the Raleigh quotients. The uniqueness again follows from the invariance under orthonormal transformations together with the need to restrict to the subspace spanned by the first k eigenvectors. \square

The implementation of this visualisation technique is very straightforward.

Algorithm 8.33 [Data visualisation] Matlab code for the data visualisation algorithm is given in Code Fragment 8.4. \blacksquare

```
% original kernel matrix stored in variable K
% tau gives the embedding in k dimensions
D = diag(sum(K));
L = D - K;
[V,Lambda] = eig(L);
Lambda = diag(Lambda);
I = find(abs(Lambda) > 0.00001);
objective = 2*sum(Lambda(I(1:k)))
Tau = V(:,I(1:k));
plot(Tau(:,1), Tau(:,2), 'x')
```

Code Fragment 8.4. Matlab code to implementing low-dimensional visualisation.