

# Clustering via Nonparametric Density Estimation: the R Package pdfCluster

Adelchi Azzalini  
Università di Padova

Giovanna Menardi  
Università di Padova

---

## Abstract

The R package **pdfCluster** performs cluster analysis based on a nonparametric estimate of the density of the observed variables. Functions are provided to encompass the whole process of clustering, from kernel density estimation, to clustering itself and subsequent graphical diagnostics. After summarizing the main aspects of the methodology, we describe the features and the usage of the package, and finally illustrate its application with the aid of two datasets.

*Keywords:* cluster analysis, graph, kernel methods, nonparametric density estimation, R, unsupervised classification, Delaunay triangulation.

---

## 1. Clustering and density estimation

The traditional approach to the clustering problem, also called ‘unsupervised classification’ in the machine learning literature, hinges on some notion of distance or dissimilarity between objects. Once one such notion has been adopted among the many existing alternatives, the clustering process aims at grouping together objects having small dissimilarity, and placing instead those with large dissimilarity in different groups. The classical reference for this approach is [Hartigan \(1975\)](#); a current standard account is [Kaufman and Rousseeuw \(1990\)](#).

In more recent years, substantial work has been directed to the idea that the objects to be clustered represent a sample from some  $d$ -dimensional random variable  $X$  and the clustering process can consequently be carried out on the basis of the distribution of  $X$ , to be estimated from the data themselves. It is usually assumed that  $X$  is of continuous type; denote its density function by  $f(x)$ , for  $x \in \mathbb{R}^d$ . We shall discuss later how the assumption that  $X$  is continuous can be circumvented.

The above broad scheme can be developed in at least two distinct directions. The first one regards  $X$  as a mixture of, say,  $J$  subpopulations, so that its density function takes the form

$$f(x) = \sum_{j=1}^J \pi_j f_j(x) , \quad (1)$$

where  $f_j$  denotes the density function of the  $j$ -th subpopulation and  $\pi_j$  represents its relative weight; here  $\pi_j > 0$  and  $\sum_j \pi_j = 1$ . In this logic, a cluster is associated with each component  $f_j$  and any given observation  $x'$  is allocated to the cluster with maximal density  $f_j(x')$  among the  $J$  components.

To render the estimation problem identifiable, some restriction must be placed on the  $f_j$ 's. This is typically achieved by assuming some parametric form for the component densities, whence the term ‘model-based clustering’ for this formulation, which largely overlaps with the notion of finite mixture modelling of a distribution. Quite naturally, the most common assumption for the  $f_j$ 's is of Gaussian type, but other families have also been considered. The clustering problem now involves estimation of the  $\pi_j$ 's and the set of parameters which identify each of  $f_1, \dots, f_J$  within the adopted parametric class. An extended treatment of finite mixture models is given by [McLachlan and Peel \(2000\)](#).

There exist some variants of the above mixture approach, but we do not dwell on them, since our main focus of interest is in the alternative direction which places the notion of a density function in a nonparametric context. The chief motivation for this choice is to free the individual clusters from a given density shape, that is, the parametric class adopted for the components  $f_j$  in Equation 1. If the cluster shapes do not match the shapes of the  $f_j$ 's, the mixture approach may face difficulties. This problem can be alleviated by adopting a parametric family of distributions more flexible than the Gaussian family. For instance, [Lin et al. \(2007\)](#) adopt the skew- $t$  distribution for the  $f_j$  components; this family provides better adaptability to the data behaviour, and correspondingly can lead to a reduced number  $J$  of components, compared to the Gaussian assumption. Although variants of this type certainly increase the flexibility of the approach, there is still motivation for considering a nonparametric formulation, completely free from assumptions on the cluster shapes.

The idea on which the nonparametric approach relies is to associate clusters with the regions around the modes of the density underlying the data. An appealing property of this formulation is that clusters are conceptually well defined. Then, their number corresponds to an intrinsic property of the data density and is operationally estimable.

It is appropriate to remark that the above two approaches involve somewhat different notions of a cluster. In the parametric setting (1) clusters are associated with the components  $f_j$  while in the nonparametric context, they are associated with regions with high density. The two concepts are different, even if they often lead effectively to the same outcome. A typical case where they diverge is provided by the mixture of two bivariate normal populations, both with markedly non-spherical distribution, such that where their tails overlap an additional mode, besides the centres of the two normal components, is generated by the superposition of the densities; see for instance Figure 1 of [Ray and Lindsay \(2005\)](#) for a graphical illustration of this situation. In this case, the mixture model (1) declares that two clusters exist, while from the nonparametric viewpoint the three modes translate into three clusters.

Since the nonparametric approach to density-based clustering has only been examined relatively more recently than the parametric one, it is undoubtedly less developed, though growing. It is not our purpose here to provide a systematic review of this literature, especially in the present setting, considering that only a few of the methodologies proposed so far have lead to the construction of a broadly-usable software. We restrict ourselves to mention the works of [Stuetzle \(2003\)](#) and [Stuetzle and Nugent \(2010\)](#). In the supplementary material provided by this latter reference, the R package `gslclust` is also available. Among the few ready-to-use techniques, a quite popular method is ‘dbscan’ by [Ester et al. \(1996\)](#), originated in the machine learning literature and available through the R package `fpc` ([Hennig 2010](#)); the notion of a data density which they adopt is somewhat different from the one of probability theory considered here. For more information on the existing contributions in this stream of literature, we refer the reader to the discussion included in the papers to be summarized in

the following section.

The present paper focuses on the clustering methodology constructed via nonparametric density estimation developed by [Azzalini and Torelli \(2007\)](#) and by [Menardi and Azzalini \(2013\)](#). Of this formulation, we first recall the main components of the methodology and then describe its R implementation ([R Core Team 2012](#)) in the package **pdfCluster** ([Azzalini and Menardi 2013](#)), illustrated with some numerical examples.

## 2. Clustering via nonparametric density estimation

### 2.1. Basic notions

The idea of associating clusters with modes or with regions of high density goes back a long time. Specifically, [Wishart \(1969\)](#) stated that clustering methods should be able to identify “distinct data modes, independently of their shapes and variance.” [Hartigan \(1975, p. 205\)](#) stated that “clusters may be thought of as regions of high density separated from other such regions by regions of low density,” and the subsequent passage expanded somewhat this point by considering ‘density-contour’ clusters formed by regions with density above a given threshold  $c$ , and showing that these regions form a tree as  $c$  varies. However, this direction was left unexplored and the rest of Hartigan’s book, as well as the subsequent mainstream literature, developed cluster analysis methodology in another direction, that is building on the notion of dissimilarity.

Among the few subsequent efforts to build clustering methods based on the idea of density functions in a nonparametric context, we concentrate on the construction by [Azzalini and Torelli \(2007\)](#) and its development by [Menardi and Azzalini \(2013\)](#). In the following, we summarize these contributions up to some minor differences, without attempting a full discussion, for which we refer the reader to the quoted papers.

For a  $d$ -dimensional density function  $f(\cdot)$ , which we assume to be bounded and differentiable, define

$$\begin{aligned} R(c) &= \{x : x \in \mathbb{R}^d, f(x) \geq c\}, & (0 \leq c \leq \max f), \\ p_c &= \int_{R(c)} f(x) \, dx, \end{aligned} \tag{2}$$

which represent the region with density values above a level  $c$  and its probability, respectively. When  $f$  is unimodal,  $R(c)$  is a connected region. Otherwise, it may be connected or not; in the latter case, it is formed by two or more connected components, corresponding to the regions around the modes of  $f$  which are identified by sectioning the density function at the  $c$  level. These notions are illustrated for the case  $d = 1$  by the left panel of [Figure 1](#), where the specific choice of  $c$  identifies a disconnected set  $R(c)$  which is formed by the two disjoint intervals at the basis of the shaded area, with associated probability  $p_c$ . As  $c$  varies, the number of connected regions varies. Since  $c$  and  $p_c$  are monotonically related, we can regard the number of connected regions as a step function  $m(p)$  of  $p$ , for  $0 < p < 1$ . The right panel of [Figure 1](#) displays the function  $m(p)$  corresponding to the density of the left panel; the probability  $p_c$  labeled on the  $x$ -axis corresponds to a value  $m(p_c) = 2$ , meaning that two connected components form the level set  $R(c)$  associated with the probability  $p_c$ .

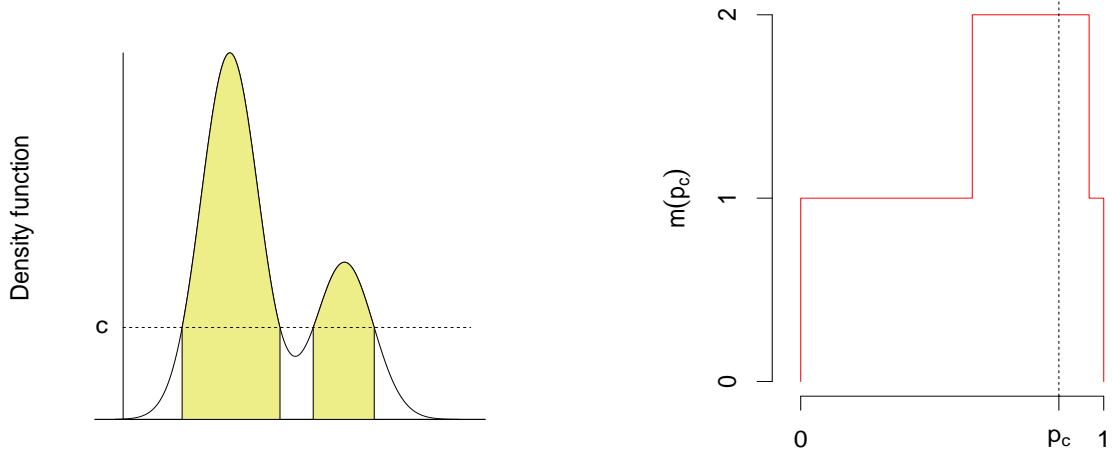


Figure 1: Density function and set  $R(c)$  for a given  $c$  (left panel) and corresponding mode function (right panel).

We shall refer to  $m(p)$  as the ‘mode function’ because it enjoys some useful properties related to the modes of  $f(x)$ . With the further convention that for convenience,  $m(0) = m(1) = 0$ , the most relevant facts are: (a) the total number of increments of  $m(p)$ , counted with their multiplicity, is equal to the number of modes,  $M$ ; (b) a similar statement holds for the number of decrements; (c) the increment of  $m(\cdot)$  at a given point  $p_c$  equals the number of modes whose ordinate is  $c$ . Inspection of the mode function allows us to see, moving along the  $x$ -axis, when a new mode is formed, or when two or more disconnected high-density sets merge into one.

It is worth to point out that, unlike many related methods which focus on a specific choice of  $c$ , we go back to the idea of [Wishart \(1969\)](#) and [Hartigan \(1975\)](#) of letting  $c$  vary from 0 to  $\max(f)$ . In this way, all high density regions can be identified, irrespectively of their level or prominence. Moreover, as established by [Hartigan \(1975, Section 11.13\)](#), the set of regions  $R(c)$  exhibits a hierarchical structure. This tree structure will be illustrated later in the examples to follow.

When a set of observations  $S = \{x_1, \dots, x_n\}$  randomly sampled from the distribution of  $X$  is available, we can compute a nonparametric estimate  $\hat{f}(x)$  of the density. The specific choice of the type of estimate is not crucial at this point, provided it is nonnegative and finite. The sample version  $\hat{R}(c)$  of  $R(c)$  is then obtained replacing  $f(x)$  by  $\hat{f}(x)$  in Equation 2, and a corresponding sample version of the mode function is introduced. Under mild regularity conditions, one can prove ‘strong set consistency’ of  $\hat{R}(c)$  to  $R(c)$ , as  $n \rightarrow \infty$ ; see, for instance, [Wong and Lane \(1983\)](#).

Since we are primarily interested in allocating observations to clusters, far more often than allocating all points of  $\mathbb{R}^d$ , this can be achieved considering

$$\begin{aligned} S(c) &= \{x_i : x_i \in S, \hat{f}(x_i) \geq c\}, & (0 \leq c \leq \max \hat{f}), \\ \hat{p}_c &= |S(c)|/n, \end{aligned} \tag{3}$$

where  $|\cdot|$  denotes the cardinality of a set. It follows from the use of a uniformly strongly consistent estimate of  $f$  that  $\hat{p}_c \rightarrow p_c$  as  $n \rightarrow \infty$ .

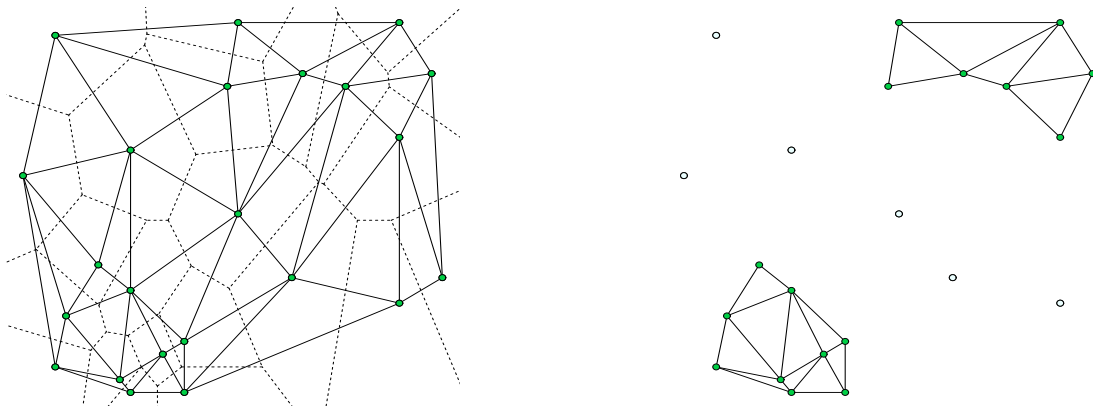


Figure 2: The left plot displays an example of Voronoi tessellation (dashed lines) for a set of points when  $d = 2$ , and superimposed Delaunay triangulation (continuous lines). The right plot removes edges of some points from the original Delaunay triangulation, keeping points with  $\hat{f} \geq c$  for some threshold  $c$ .

The above construction is conceptually simple and clear, but its actual implementation is problematic. While for  $d = 1$  identification of  $R(c)$  and of  $S(c)$  is elementary, as perceivable from Figure 1, the problem complicates substantially for  $d > 1$ , which of course is the really interesting case. More specifically, it is immediate to state whether any given point  $x$  belongs to any given set  $R(c)$ , but it is harder to say how many connected sets comprise  $R(c)$ , and which one they are; a similar problem exists for  $S(c)$ . The next two sections describe two routes to tackle this question.

## 2.2. Spatial tessellation

To establish whether a set  $S(c)$  is formed by points belonging to one or more connected regions which comprise  $\hat{R}(c)$ , [Azzalini and Torelli \(2007\)](#) make use of some concepts from computational geometry. The first of these is the Voronoi tessellation which partitions the Euclidean space into  $n$  polyhedra, possibly unbounded, each formed by those points of  $\mathbb{R}^d$  which are closer to one given point in  $S$  than to the others. Conceptually, from here one derives the Delaunay triangulation which is formed by joining points of  $S$  which share a facet in the Voronoi tessellation. From a computational viewpoint, the Delaunay triangulation is the only one required for the subsequent steps and can be obtained directly, without forming the Voronoi tessellation first. The elements of the Delaunay triangulation are simplices in  $\mathbb{R}^d$ ; since for  $d = 2$  they reduce to triangles, this explains the term ‘Delaunay triangulation.’ For a detailed treatment, see [Okabe \*et al.\* \(1992\)](#).

These notions are illustrated in the left panel of Figure 2 which shows the Voronoi tessellation and the Delaunay triangulation, for a set of points in  $\mathbb{R}^2$ .

The procedure of [Azzalini and Torelli \(2007\)](#) consists of two main phases. The first one comprises itself a few steps, as follows. First, we construct the Delaunay triangulation of the sample  $S$ , and compute the nonparametric estimate  $\hat{f}(x_i)$  for each  $x_i \in S$ . Then, for any given value  $p_c \in (0, 1)$ , we eliminate all points  $x_i$  such that  $\hat{f}(x_i) < c$  and determine the connected sets of the remaining points. This step is illustrated graphically in the right panel

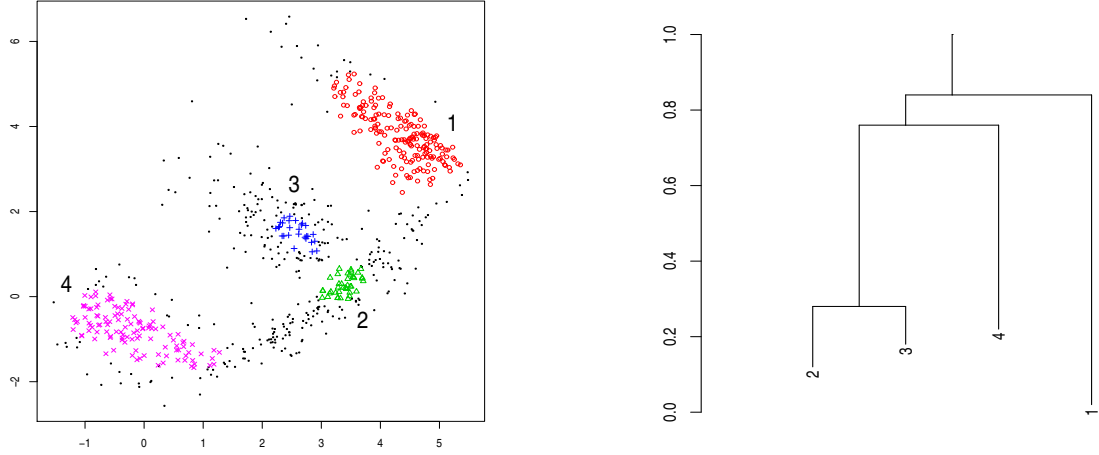


Figure 3: Cluster cores and cluster tree for a set of simulated data with  $d = 2$ . In the left panel the points belonging to the cluster cores are marked with different symbols. The unlabelled points are marked by dots. The right panel shows the corresponding cluster tree.

of Figure 2, where two connected sets are visible after removing the points of low density and the associated arcs from the triangulation on the left panel. In principle, this operation is repeated for all possible values  $p_c \in (0, 1)$ , while in practice for a fine grid of such points. Provided that we have kept track of the group membership of the sample components as  $p_c$  ranges from 0 to 1, at the end of this process, we can construct a tree of the connected sets. The leaves of this tree correspond to the modes of the density function, and for each level of  $p_c$  the tree provides the number of associated connected components of  $R(c)$ . In the same process, say,  $M$  groups of points have been singled out, formed by the connected sets so identified; we call them ‘cluster cores.’ Essentially, each cluster core is formed by the points that unquestionably pertain to one mode. The density function in the left panel of Figure 1, for instance, would give rise to two cluster cores, formed by those points with density above the lowest level which identifies a level set made of two connected components. It is then a quite distinctive feature of this method to pinpoint a number  $M$  of groups, while most methods require that  $M$  is specified on input or it is left undetermined, like in hierarchical distance-based methods.

The outcome of the first phase is illustrated in Figure 3 for a set of simulated data with  $d = 2$ . The left panel displays the observation points, marked with different symbols for the four cluster cores; the unlabelled points are denoted by a simple dot. The right panel shows the cluster tree of the four groups. Notice that the tree is upside-down with respect to the density: the root of the tree corresponds to zero density and the top mode, marked by 1, is at the bottom. This confirms that it would make more sense if mathematical trees had their roots at the bottom, like real trees.

In the second phase of the procedure, the data points still unlabelled must be allocated among the  $M$  cluster cores. This operation is a sort of classification problem, although of a peculiar type, since the unlabelled points are not randomly sampled from  $X$ , but inevitably in the outskirts of the cluster cores. [Azzalini and Torelli \(2007\)](#) propose a block-sequential criterion of allocation, based on density estimation and density ratios. Fixed a number  $K$  of stages, at

each stage compute an estimate  $\hat{f}_m(x_u)$  of the density of all the unallocated data  $x_u$ , based on the only data already allocated to group  $m$  (for  $m = 1, 2, \dots, M$ ). Then, for each unallocated  $x_u$ , compute the log-ratios

$$r_m(x_u) = \log \frac{\hat{f}_m(x_u)}{\max_{l \neq m} \hat{f}_l(x_u)} \quad m = 1, 2, \dots, M, \quad (4)$$

where the use of the log-ratios, instead of the simple ratios, is employed only for conventional reasons. Then, for each  $x_u$ , retain the highest  $r_m(x_u)$  value, say  $r_{m_0}(x_u)$ , and sort the set of these  $r_{m_0}(x_u)$  values in decreasing order; finally, allocate those  $x_u$ 's for which the corresponding  $r_{m_0}(x_u)$  belongs to the top  $(1 - K/100)$ -th fraction of the distribution of these values. At the next stage the process is repeated with  $\hat{f}_m(x_u)$  and  $r_m(x_u)$  updated with the new points which have been allocated to the groups.

A variant of this allocation rule, not examined by [Azzalini and Torelli \(2007\)](#), but which we have found preferable on the whole as it takes into account the precision of the estimates, weights the log-ratios in Equation 4 inversely with their variability while determining the order of allocation of the low density data. In practice, computation of the standard error of Equation 4 is quite complicated, even by employing asymptotic expressions of variances. We take a rough approximation of those quantities, instead, by first identifying the index  $m'$  such that  $\hat{f}_{m'}(x_0) = \max_{k \neq m} \hat{f}_k(x_0)$  and then considering the  $m'$  index as given. Next, we apply standard approximations for transformation of variables; see [Bowman and Azzalini \(1997, page 29\)](#), for the specific case of transforming  $\hat{f}(\cdot)$ .

As a diagnostic device to evaluate the quality of the clusters so obtained, the density-based silhouette (*dbs*) proposed by [Menardi \(2011\)](#) fits naturally in this framework. This tool is the analogue of the classical silhouette information ([Rousseeuw 1987](#)), when the distances among points are replaced by probability log-ratios. Specifically, define for observation  $x_i$ ,

$$\hat{\tau}_m(x_i) = \frac{\pi_m \hat{f}_m(x_i)}{\sum_{k=1}^M \pi_k \hat{f}_k(x_i)}, \quad m = 1, \dots, M, \quad (5)$$

where  $\pi_m$  is a ‘prior’ probability for the  $m$ -th group. Its specification may depend on the prior knowledge about the composition of the clusters and a lack of information would imply the choice of a uniform distribution of the  $\pi_m$  over the groups. Alternatively, information derived from the detected partition can also be used. For instance, these probabilities can be chosen proportional to the cardinalities of the cluster cores.

The *dbs* for  $x_i$  is defined as

$$dbs_i = \frac{\log \left( \frac{\hat{\tau}_{m_0}(x_i)}{\hat{\tau}_{m_1}(x_i)} \right)}{\max_{x_i} \left| \log \left( \frac{\hat{\tau}_{m_0}(x_i)}{\hat{\tau}_{m_1}(x_i)} \right) \right|}, \quad (6)$$

where  $m_0$  is the group to which  $x_i$  has been allocated and  $m_1$  is the group for which  $\hat{\tau}_m(x_i)$  takes the second largest value. The *dbs* of one point is then proportional to the log-ratio between the posterior probability that it belongs to the group to which it has been allocated and the maximum posterior probability that it belongs to another group. Then, its interpretation is the same as for the classical ‘silhouette,’ namely large values of *dbs* are evidence of a well clustered data point while small values of *dbs* mean a low confidence in the classification.

We close this section with some remarks on computational aspects. From the point of view of memory usage, while for methods based on dissimilarities it grows quadratically, the requirement of this procedure grows linearly with  $n$ , as it depends on the number of points on which the density is evaluated and on the number of arcs from each point in the Delaunay triangulation (that is, number of shared Voronoi facets). From a computational point of view, the most critical aspect here is the construction of the Delaunay triangulation. This can be produced by the Quickhull algorithm by [Barber \*et al.\* \(1996\)](#), whose implementation is publicly available at <http://www.qhull.org><sup>1</sup>. This algorithm works efficiently for increasing  $n$  when  $d$  is small to moderate, but computing time grows exponentially when  $d$  increases.

### 2.3. Pairwise connections

The final remarks of the previous section motivate the development of a variant procedure to build a connection network of the elements of  $S$  by using a different criterion instead of the Delaunay triangulation, leaving unchanged the rest of the above-described process.

The proposal of [Menardi and Azzalini \(2013\)](#) starts by reconsidering the notion of connected sets for  $d = 1$  and then extending this view to the case  $d > 1$ . The basic idea is to examine the behaviour of  $\hat{f}(x)$  when we move along a segment  $[x_1, x_2]$ , since it depends on whether the sample values  $x_1$  and  $x_2$  belong to the same connected set of  $\hat{R}(c)$  or not. To visualize the process, it is convenient to refer to the left panel of Figure 1, regarding the density there as the estimate  $\hat{f}$ , and consider the set  $\hat{R}(c)$  formed by the union of the two intervals of the shaded area. If  $x_1$  and  $x_2$  belong to the same interval, then the corresponding portion of density along the segment  $(x_1, x_2)$  has no local minimum. On the contrary, if  $x_1$  and  $x_2$  belong to different subsets of  $\hat{R}(c)$ , then at some point along  $[x_1, x_2]$  the density exhibits a local minimum, which we shall refer to as ‘presence of a valley.’

When  $d > 1$ , the same idea can be carried over by examining the behaviour of the section of  $\hat{f}(x)$  along the segment joining  $x_1$  and  $x_2$  where  $x_1, x_2 \in \mathbb{R}^d$  and applying the same principle as above. A graph is then created whose vertices are the sample points and an arc is set between any pair of points such that there is no valley in the section of  $\hat{f}(x)$  between them.

In practical terms, the claim of presence of a valley must allow some tolerance for the inevitable variability of  $\hat{f}$ . Given the above premises, we must take into account the amplitude of the valley detected along the stated section of  $\hat{f}$ , and declare that  $x_1$  and  $x_2$  are connected points if this amplitude is below a certain threshold. Clearly, if no valley exists, the connection of  $x_1$  and  $x_2$  is unquestioned.

The broad idea of tolerance must be given a specific form to be operational. For simplicity, we describe here only informally the criterion adopted by [Menardi and Azzalini \(2013\)](#), as a general definition of the quantities involved, allowing for multiple and nested valleys, would be quite tricky. We refer the reader to the original paper for the full specification. When a valley is detected, we introduce an auxiliary function  $\varphi$  derived from the original one by increasing it by the minimum amount required to fill the valley; one can think that water is poured into the valley until it starts to overflow. To visualize this process, consider Figure 4 where, in the first panel, a sample of points in  $\mathbb{R}^2$  is depicted and the segments joining two pairs of them,  $(x_1, x_2)$  and  $(x_3, x_4)$ , are highlighted. The second panel displays the section of  $\hat{f}(x)$  along the segment joining  $x_1$  and  $x_2$ , represented by the smooth curve, and the auxiliary function  $\varphi(x)$  which fills the valley. Denote by  $A_1$  the dark-shaded area, that is the area filled

<sup>1</sup>Conditions for copying, modifying, and redistributing the source code are available on the same web site.

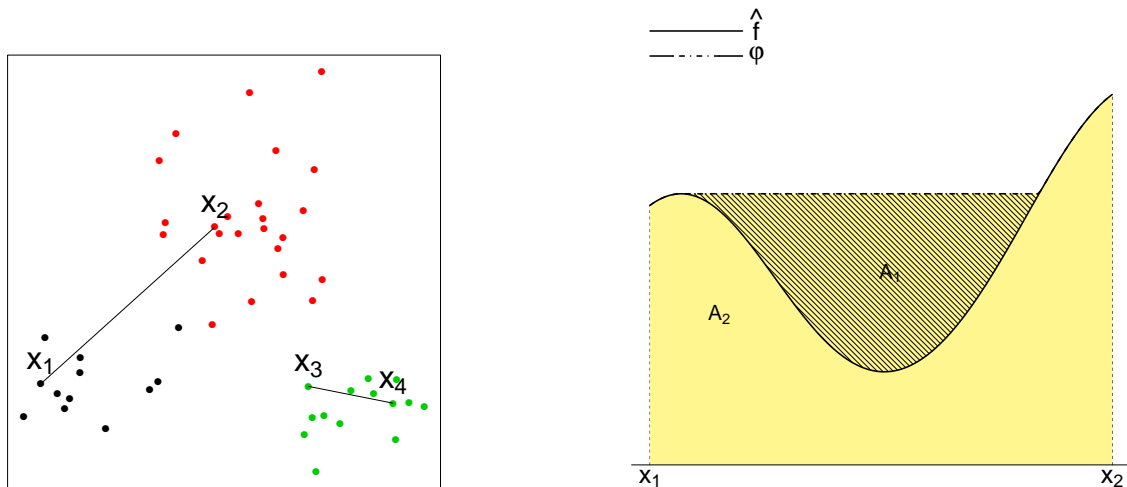


Figure 4: The left panel displays a set of points in  $\mathbb{R}^2$ , of which two pairs,  $(x_1, x_2)$  and  $(x_3, x_4)$ , are highlighted by drawing the segment joining them. The right plot displays the section of  $\hat{f}(x)$  along the segment joining  $x_1$  and  $x_2$  (solid curve delimiting the light-shaded area  $A_2$ ) and the associated non-decreasing function  $\varphi(x)$  delimiting the whole filled area  $A_1 + A_2$ .

by “pouring water”, and by  $A_2$  the light-shaded area, that is the area under  $\hat{f}$ . The amplitude of the valley in this case is quantified by

$$V = \frac{A_1}{A_1 + A_2} \in [0, 1). \quad (7)$$

For a given tolerance parameter  $\lambda \in (0, 1)$ , if  $V < \lambda$  then  $x_1$  and  $x_2$  are considered connected points, and an edge is set in the connection graph. For the pair  $(x_3, x_4)$  in Figure 4, the section of  $\hat{f}(x)$  along the segment joining them is concave and  $\varphi(x)$  coincides with  $\hat{f}(x)$  because there is no valley (plot not shown); hence in this case, as well as in those situations where  $\hat{f}(x)$  is monotone along the segment,  $V = 0$  and the points are declared to be connected.

Once the connection of all pairs of sample values has been examined, the rest of the process is carried out exactly as in the previous section. Since now we are always working in a one-dimensional world, higher values of  $d$  can be tackled. However, for large  $d$  we are facing another problem: the degrade of performance of nonparametric density estimates, known as ‘curse of dimensionality.’ On the other hand, it can be argued that for the clustering problem we need to identify only the main features of a distribution, especially its modes, not the fine details. This consideration indicates that the method can be considered also for a broader set of cases than those where the density is the focus of interest. See [Menardi and Azzalini \(2013\)](#) for an extended discussion of this issue.

### 3. The R package pdfCluster

#### 3.1. Package overview

The R package **pdfCluster** performs, as its main task, cluster analysis by implementing the

	clustering	density estimation	diagnostics
S4 class	"pdfCluster-class"	"kepdf-class"	"dbs-class"
functions to produce the class	pdfCluster() pdfClassification()	kepdf()	dbs()
related methods	pdfCluster plot show summary	plot show summary	dbs plot show summary
other functions	groups()	h.norm() hprop2f()	adj.rand.index()

Table 1: Summary of the **pdfCluster** package. Each of the three main tasks of the package is associated with a set of functions and methods, and results in an object of a specific **S4** class.

methodology described in the previous sections. Furthermore, two other tasks are accomplished: density estimation and clustering diagnostics.

Each of these tasks is associated with a set of functions and methods, and results in an object of a specific class, as summarized in Table 1. The package is built by making use of classes and methods of the **S4** system (Chambers 1998). The **S4** system requires to declare and explicitly define the nature of all classes and methods involved in programming, thus guaranteeing a stricter validation and a greater consistency of the programming.

To ease presentation, an overview of the package is first provided, aiming at an introductory usage of its features. The next section is devoted to a more in-depth examination of computational and technical details.

### *Clustering via nonparametric density estimation*

The package is built around the main routine `pdfCluster()`, which actually performs the clustering process. `pdfCluster()` is defined as a generic function and dispatches different methods depending on the class of its first argument. For a simple use of the function, the user is only required to provide the data  $\mathbf{x}$  to be clustered, in the form of a vector, "`matrix`", or "`data.frame`" of "`numeric`" elements. A further method dispatched by the `pdfCluster()` generic function applies to objects of class "`pdfCluster`" itself. This option will be described in Section 3.2, item (c).

Further arguments include `graphtype`, which defines the procedure to identify the connected components of the density level sets. The elementary case  $d = 1$  is handled by the "`unidimensional`" option. When data are multidimensional, instead, this argument may be set to "`deLaunay`" or "`pairs`", to run the procedures described in Section 2.2 and 2.3, respectively. When not specified, the latter option is selected if  $d > 6$ . When "`pairs`" is selected, explicitly or implicitly, the user may wish to specify the parameter `lambda`, which defines the tolerance threshold to claim the presence of valleys in the section of  $\hat{f}(x)$  between pairs of points. Default value is `lambda = 0.10`, a choice that has been empirically proven to be usually effective; see also Section 3.2, item (b).

After that the connected components associated with the density level sets are identified,

`pdfCluster()` builds the cluster tree and detects the cluster cores. An internal call to function `pdfClassification()` follows by default, to carry on the second phase of the clustering procedure, that is, the allocation of the lower density data points not belonging to the cluster cores. The user can regulate the process of classification by setting some optional parameters to be passed to `pdfClassification()`. Details are discussed in Section 3.2, item (f).

Further optional arguments may be given to `pdfCluster()` in order to regulate density estimation. These arguments are passed internally to function `kepdf()`, which is described below.

The outcome of the clustering process is encapsulated in an object of class `"pdfCluster"`, whose slots include, among others, a list `"nc"` providing details about the connected components identified for each considered density threshold  $c$  at which the level sets have been evaluated, an object of class `"dendrogram"`, providing information about the cluster tree, and the number `"noc"` of detected groups. Furthermore, when the classification procedure is carried on, the slot `"stages"` is a list with elements corresponding to the data allocation to groups at the different stages of the classification procedure. Clusters may be extracted from objects of class `"pdfCluster"` by the application of function `groups()`. The optional argument `stage` may be specified to extract the cluster cores (set `stage = 0`) or the cluster labels at the different stages of the classification procedure.

Methods to `show`, to provide a `summary` and to `plot` objects of class `"pdfCluster"` are available. Four types of plot are selectable, by setting the argument `which`. Argument `which = 1` plots the mode function; `which = 2` plots the cluster tree; a scatterplot of the data or of pairs of selected coordinates reporting the label group is provided when `which = 3` and `which = 4` plots the density-based silhouette information as described below. Multiple choices are possible.

### *Density estimation*

Density estimation is performed by the kernel method throughout the `kepdf()` function. Estimates are computed by a product estimator of the form:

$$\hat{f}(y) = \sum_{i=1}^n \frac{1}{nh_{i,1} \cdots h_{i,d}} \prod_{j=1}^d K\left(\frac{y_j - x_{i,j}}{h_{i,j}}\right). \quad (8)$$

The kernel function  $K$  is an argument of function `kepdf()` and can either be a Gaussian density (if `kernel = "gaussian"`) or a Student's  $t_\nu$  density, with  $\nu = 7$  degrees of freedom (when `kernel = "t7"`). The choice of the kernel function is known not to be critical in density estimation, thus a Gaussian kernel is, in general, adequate. The uncommon option of selecting a  $t_7$  distribution is motivated only by computational reasons, as explained in Section 3.2, item (h).

The user may choose to estimate the density with a fixed or an adaptive bandwidth  $h_i = (h_{i,1} \cdots h_{i,d})^\top$ , by setting the argument `bwtype` accordingly. Leaving the argument unspecified entails the use of a fixed bandwidth estimator. When `bwtype = "fixed"`, that is  $h_i = h$ , a constant smoothing vector is used for all the observations  $x_i$ . Default values are set as asymptotically optimal for a multivariate Normal distribution (see, e.g., [Bowman and Azzalini 1997](#), page 32). Alternatively, `bwtype = "adaptive"`, which corresponds to specifying a vector of bandwidths  $h_i$  for each observation  $x_i$ . Default values are selected according to the

approach described by Silverman (1986, Section 5.3.1), implemented in the package through the function `hprop2f()`.

The outcome of the application of the `kepdf()` function is encapsulated in an object of class `"kepdf"`, whose slots include the `"estimate"` at the evaluation points and the parameters used to obtain that estimate.

Methods which `show`, provide a `summary`, and `plot` objects of class `"kepdf"` are also available. When the density estimate is based on two or higher dimensional data, these functions make use of functions `contour()`, `image()` and `persp()`, depending on how the argument `method` is set. When  $d > 2$ , the pairwise marginal estimated densities are plotted for all pairs of coordinates, or for a subset of coordinates specified by the user via the argument `indcol`.

### *Diagnostics of clustering*

As a third feature, the package provides diagnostics of the clustering outcome. Density-based silhouette is computed by the generic function `dbs()`, which dispatches two methods. One method applies to objects of class `"pdfCluster"` directly; a second method is thought to compute the density-based silhouette information on partitions produced by a possibly different density-based clustering technique. The latter method applies to two arguments: the matrix of clustered data and a numeric vector of cluster labels.

Computation of the density-based silhouette requires the density function to be estimated, conditional to the group membership. Hence, further arguments of function `kepdf()` can be given to `dbs()` to set parameters of density estimation. Moreover, some `prior` probability may be specified for each cluster, as will be clarified in the next section.

Results of the application of function `dbs()` are provided in objects of class `"dbs"`. An S4 method for plotting objects of class `"dbs"` is available: data are partitioned into the clusters, sorted in a decreasing order with respect to their `dbs` value and displayed on a bar graph.

As a further diagnostic tool, the package provides the function `adj.rand.index()` which evaluates the agreement between two partitions, through the adjusted Rand criterion (Hubert and Arabie 1985).

## 3.2. Further details

- (a) As already mentioned, `pdfCluster` automatically selects the procedure to be used for detecting connected components of the density level sets, depending on the data dimensionality. While the user is allowed to change this choice by setting argument `graphtype`, we warn against setting the argument `graphtype = "delaunay"` for large dimensions. The number of operations required to compute the Delaunay triangulation grows with  $n^{\lfloor d/2 \rfloor}$  while the computational complexity due to running the pairwise connection criterion grows quadratically with the sample size. Hence, at the present state of computing resources, running the Delaunay triangulation when  $d > 6$  appears hardly feasible for values of  $n$  greater than about 200. Instead, data with any dimensionality may be handled by the pairwise connection criterion, although the computational speed slows down, and memory requirement increases for very large  $n$ .
- (b) The higher computational feasibility of the pairwise connection criterion is paid for by the need of setting the tolerance threshold  $\lambda$ . According to our experience, the default

value of `lambda = 0.10` is usually a sensible choice for moderate to high dimension (say, for  $d > 6$ ) while a lower `lambda` is sufficient in low dimensions, when the density estimate is more accurate. A larger value can be useful when the procedure detects a number of small spurious groups, because this choice results in aggregating clusters. While there is not an automatic procedure to set the value of  $\lambda$ , we recall that the introduction of this parameter has the only purpose of disregarding valleys of negligible extent, presumably due to small sampling errors. Hence, while in principle  $\lambda$  is defined in  $[0, 1]$ , large values (say, `lambda > 0.3`) would be meaningless. We refer the reader to the paper of [Menardi and Azzalini \(2013\)](#) for further discussion.

- (c) Since running the procedure several times with different choices of `lambda` may be time consuming, the package allows for a more efficient route, implemented by an additional method dispatched by function `pdfCluster()`. Once an object of class "pdfCluster" is created by function `pdfCluster()` with argument `graphtype = "pairs"`, `pdfCluster()` can be called again by setting the same object of class "pdfCluster" as a first argument `x` and a different value of `lambda`. Slot "graph" of the object with class "pdfCluster" contains the amplitude of the valleys detected by the evaluation of the density along the segments connecting observations. Then, the pairwise evaluation does not need to be run again to check results for different values of `lambda` and the procedure speeds up considerably. An example will be presented in the next section.
- (d) Both the Delaunay and the pairwise connection criteria to build a graph on the observed data are implemented by some specifically designed foreign functions. In the former case, the `delaunayn` function in package **geometry** ([Barber et al. 2012](#)) is the R interface to the Quickhull algorithm. Pairwise connection is, instead, implemented in the C language.
- (e) After building the selected connection network among the observations, `pdfCluster` determines the level set  $S(c)$  for a grid of values of  $p_c$ . The number of points of such grid may be chosen by the user through the `n.grid` argument, with the proviso that a fine grid is selected. The default value for `n.grid` is the set to the minimum between  $n$  and  $\lfloor (5 + \sqrt{n})4 + 0.5 \rfloor$ , which is just an empirically-based rule. For each value of  $p_c$  the identification of connected components of  $S(c)$  is carried out by means of a C function borrowed from the R package **spdep** ([Bivand 2013](#)), which implements a depth first search algorithm.
- (f) The procedure to allocate the low density data to the cluster cores is block-sequential and the user is allowed to select the number `n.stage` of such blocks as an optional parameter of `pdfCluster()`, to be passed internally to `pdfClassification()` (default value is `n.stage = 5`). When this argument is set to 0, the clustering process stops when the cluster cores are identified. Otherwise, further arguments can be passed from `pdfCluster()` to `pdfClassification()`. Among them, `se` takes "logical" values, and `se = TRUE` accounts for the standard error of the log-likelihood ratios in Equation 4. Argument `hcores` declares if the densities in Equation 4 have to be estimated by selecting the same bandwidths as the ones used to form the cluster cores. Default value is set to `FALSE`, in which case a vector of bandwidths specific for the clusters is used.
- (g) `pdfCluster()` makes an internal call to function `kepdf()` both to estimate the density underlying the data and to build the connection network when the pairwise connec-

tion criterion is selected. By default, a kernel density estimation with fixed kernel is built, with vector of smoothing parameters set to the one asymptotically optimal under the assumption of multivariate normality. Although arguably sub-optimal, this choice produces sensible results in most applications. When the dimensionality of the data is low-to-moderate, it is often advantageous to shrink the smoothing parameter slightly towards zero; we adopt a shrinkage factor `hmult = 3/4` when  $d \leq 6$ , as recommended by Azzalini and Torelli (2007), but the default value may be optionally modified by the user. For higher-dimensional data, instead, we suggest the use of a kernel estimator with adaptive bandwidth, which can be obtained by setting argument `bwtype` to `"adaptive"`.

- (h) Function `kepdf` represents the R interface of two C routines which allow to speed computations. Each of these routines is designed to perform kernel density estimation with a specific kernel function. It is worth to remind that, when connected sets are identified by the pairwise connection criterion, computation of the  $V$  measure in Equation 7 requires the density function to be evaluated along the segments joining each pair of observations. In practice, a grid of `grid.pairs` points is considered for each segment, so that the number of operation required grows with  $n^2 \text{grid.pairs}$ .

When the sample size is very large, any saving in the arithmetic computations of the kernel can make a noticeable difference. In particular, the use of a Gaussian kernel requires a call to the exponential function at each evaluation, which is computationally more expensive than sum, multiplication and power function. This explains the non-standard option to select a  $t_\nu$  kernel, with  $\nu = 7$  degrees of freedom. The relatively more critical computation involves an 4th degree power, which can be coded efficiently and still the kernel has unbounded support, which is more appropriate for the classification stage than alternatives like bisquare or similar kernels. The use of this option is then recommended when the sample size is huge.

- (i) To compute `dbs()`, it is possible to specify some prior probability for each cluster, through the argument `prior`. To this end, consider that clusters are labelled according to the maximal value of the density, so that cluster 1 includes the observations with maximal density, cluster 2 includes observations pertaining to the second highest mode and so on. As already mentioned, the choice of these probabilities may depend on the prior knowledge about the composition of the clusters. If no prior information is available, it makes sense to assign equal prior probability to all clusters. This is the default choice for the `dbs` method which applies to objects of class different from `"pdfCluster"`. Alternatively, information derived from the detected partition can also be used. When diagnostics are computed on an object of class `"pdfCluster"`, 'prior' probabilities can be chosen as proportional to the cardinalities of the cluster cores. This is then the default choice when applying the `dbs` method specific for objects of class `"pdfCluster"`. Conversely, when using a model-based clustering method, the mixing proportions could be a natural choice.

## 4. Some illustrative examples

### 4.1. Quantitative variables

The `wine` data set was introduced by [Forina \*et al.\* \(1986\)](#). It originally included the results of 27 chemical measurements on 178 wines grown in the same region in Italy but derived from three different cultivars: Barolo, Grignolino and Barbera. The **pdfCluster** package provides a selection of 13 variables. The data set is here used to illustrate the main features of the package.

As a first simple example, let us suppose to have some knowledge about which variables are relevant to the aim of reconstructing the cultivar of origin of each wine. We then perform cluster analysis on a small subset of variables.

```
> library("pdfCluster")
> data("wine", package = "pdfCluster")
> winesub <- wine[, c(2, 5, 8)]
```

As the number of considered variables is very small, visual exploration of the density estimate of the data may already give us some indication about the clustering structure.

```
> pdf <- kepdf(winesub)
> plot(pdf)
```

The resulting plot is displayed in [Figure 5](#). A three-cluster structure is clearly evident from the marginal density estimate of the variables “Alcohol” and “Flavanoids”. The main content of the object `pdf` having class “`kepdf`” may be printed by the associated `show` method.

```
> pdf
```

```
An S4 object of class "kepdf"
```

```
Call: kepdf(x = winesub)
```

```
Kernel:
```

```
[1] "gaussian"
```

```
Estimator type: fixed bandwidth
```

```
Diagonal elements of the smoothing matrix: 0.3750856 1.542968 0.4614995
```

```
Density estimate at evaluation points:
```

```
[1] 0.015211471 0.001994922 0.009822658 0.010526400 0.009014892 0.013104296
[7] 0.005910667 0.013900582 ...
```

Clustering is performed by a call to `pdfCluster()`. Note that its usage does not require a preliminary call to `kepdf()`. A summary of the resulting object provides the cardinalities of both the cluster cores and the clusters, along with the main structure of the cluster tree.

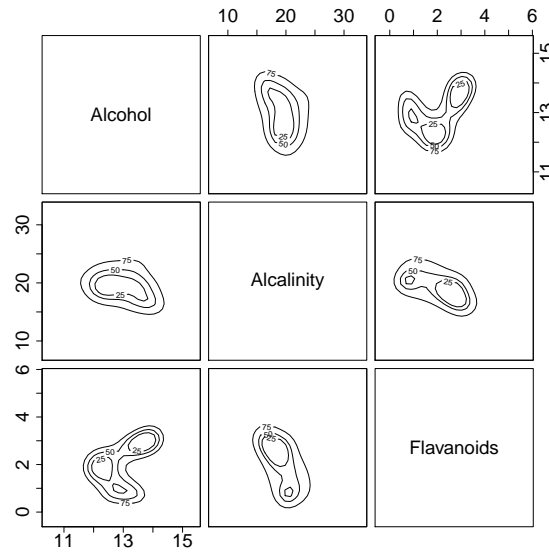


Figure 5: Plot of the pairwise marginal density estimates of three variables of the *wine* data, as given by function `kepdf`.

```
> cl.winesub <- pdfCluster(winesub)
> summary(cl.winesub)
```

An S4 object of class "pdfCluster"

Call: pdfCluster(x = winesub)

Initial groupings:

label	1	2	3	NA
count	29	15	17	117

Final groupings:

label	1	2	3
count	62	63	53

Groups tree (here 'h' denotes 'height'):

```
--[dendrogram w/ 1 branches and 3 members at h = 1]
  |--[dendrogram w/ 2 branches and 3 members at h = 0.361]
    |--leaf "1 "
    `--[dendrogram w/ 2 branches and 2 members at h = 0.333]
      |--leaf "2 " (h = 0.0556 )
      `--leaf "3 " (h = 0.0694 )
```

The group assignment can be accessed through function `groups()`:

```
> groups(cl.winesub)
```

```
[1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1
[38] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 1 2 2 2 1 2 3 2 3 3 ...
```

The optional argument `stage` of function `groups()` allows to extract groups at different stages of allocation. To extract the cluster cores write:

```
> groups(cl.winesub, stage = 0)
```

```
[1] 1 NA NA NA NA 1 NA 1 NA 1 NA NA 1 NA NA 1 NA NA NA 1 1 NA 1 NA 1
[26] NA 1 NA 1 NA NA NA 1 1 1 NA NA 1 NA NA 1 1 NA 1 NA NA 1 1 ...
```

or to see the group labels assigned up to the selected `stage` of the classification procedure:

```
> groups(cl.winesub, stage = 2)
```

```
[1] 1 NA 1 NA 1 1 NA 1 NA 1 NA NA 1 NA NA 1 NA 1 NA 1 1 NA 1 NA 1
[26] NA 1 1 1 NA NA 1 1 1 1 1 1 1 1 NA 1 1 NA 1 1 NA 1 1 ...
```

Here, the NA values label observations that are still unallocated at the selected stage.

Objects of class `pdfCluster` may be further inspected by accessing its slots. Slot `"pdf"`, for instance, contains main information about the estimated density:

```
> cl.winesub@pdf
```

```
$kernel
```

```
[1] "gaussian"
```

```
$bwtype
```

```
[1] "fixed"
```

```
$par
```

```
$par$h
```

```
Alcohol Alcalinity Flavanoids
0.2813142 1.1572259 0.3461246
```

```
$par$hx
```

```
NULL
```

```
$estimate
```

```
[1] 0.021153490 0.003723019 0.009561598 0.013346244 0.011821547 ...
```

Note that the vector of smoothing parameters used to estimate density, during the process of clustering, differs from the one produced by the previous call to function `kepdf()`, whose default value is asymptotically optimal for Normal data, as given by function `h.norm()`. As already mentioned, when low-dimensional data are clustered (in this example  $d = 3$ ), this vector is multiplied by a shrinkage factor  $3/4$  by default. To change the shrinkage factor, select the optional argument `hmult`, as illustrated below.

The user may be also interested about details regarding the procedure used to find the connected components associated with the level set, available through the slot `"graph"`.

```
> cl.winesub@graph
```

```
$type
[1] "delaunay"
```

The user may explicitly select this option, by setting the argument `graphtype` from "delaunay" to "pairs", as follows:

```
> cl.winesub.pairs <- pdfCluster(winesub, graphtype = "pairs")
> summary(cl.winesub.pairs)
```

An S4 object of class "pdfCluster"

```
Call: pdfCluster(x = winesub, graphtype = "pairs")
```

```
Initial groupings:
label    1    2    3  NA
count   12    3   10 153
```

```
Final groupings:
label    1    2    3
count   58   69   51
```

```
Groups tree (here 'h' denotes 'height'):
--[dendrogram w/ 1 branches and 3 members at h = 1]
  |--[dendrogram w/ 2 branches and 3 members at h = 0.167]
    |--[dendrogram w/ 2 branches and 2 members at h = 0.0972]
      | |--leaf "1 "
      | `--leaf "2 " (h = 0.0556 )
      `--leaf "3 " (h = 0.0694 )
```

As previously mentioned, computational arguments explain the preference for adopting in low dimension the Delaunay graph and in higher dimension the procedure to identify connected components described in Section 2.3. Beyond this, there is no particular reason to give preference to one of the two procedures, since they often produce comparable results:

```
> table(groups(cl.winesub), groups(cl.winesub.pairs))

      1  2  3
1 58  4  0
2  0 63  0
3  0  2 51
```

Additional information about the detected partition may be further visualized through a call to the associated `plot` methods. If argument `which` is not selected, the four available types of plots are displayed one at a time.

```
> plot(cl.winesub)
```

The resulting plots are reported in Figure 6. In particular, the diagnostic plot presents values of the *db*s appreciably larger than zero for almost all the observations throughout the three groups, suggesting the soundness of the detected partition. This is confirmed by the cross-classification frequencies of the obtained clusters and the actual cultivar of origin of the wines (indicated in the first column of the `wine` data).

```
> table(wine[, 1], groups(cl.winesub))
```

	1	2	3
Barolo	58	1	0
Grignolino	4	62	5
Barbera	0	0	48

Consider now the entire data set `wine`, only removing the true label class of wines.

```
> wineall <- wine[, -1]
```

When high-dimensional data are clustered, some caution should be used to deal with the curse of dimensionality and the increased variability of the density estimate. Indeed, the use of a fixed bandwidth in high dimension may cause some troubles:

```
> cl.wineall.fixedbwd <- pdfCluster(wineall)
> summary(cl.wineall.fixedbwd)
```

An S4 object of class "pdfCluster"

Call: pdfCluster(x = wineall)

Initial groupings:

label	1	2	3	4	5	6	7	8	9	10	NA
count	31	26	4	7	7	5	2	2	2	2	90

Final groupings:

label	1	2	3	4	5	6	7	8	9	10
count	31	29	17	17	16	12	13	24	17	2

...

The use of a fixed bandwidth has led to a large number of small clusters, that raises doubts about the accuracy of the estimate, due to the quite high dimensionality of data. For moderate to high dimensions, [Menardi and Azzalini \(2013\)](#) suggest to allow for different amount of smoothing by using an adaptive bandwidth:

```
> cl.wineall <- pdfCluster(wineall, bwtype = "adaptive")
> summary(cl.wineall)
```

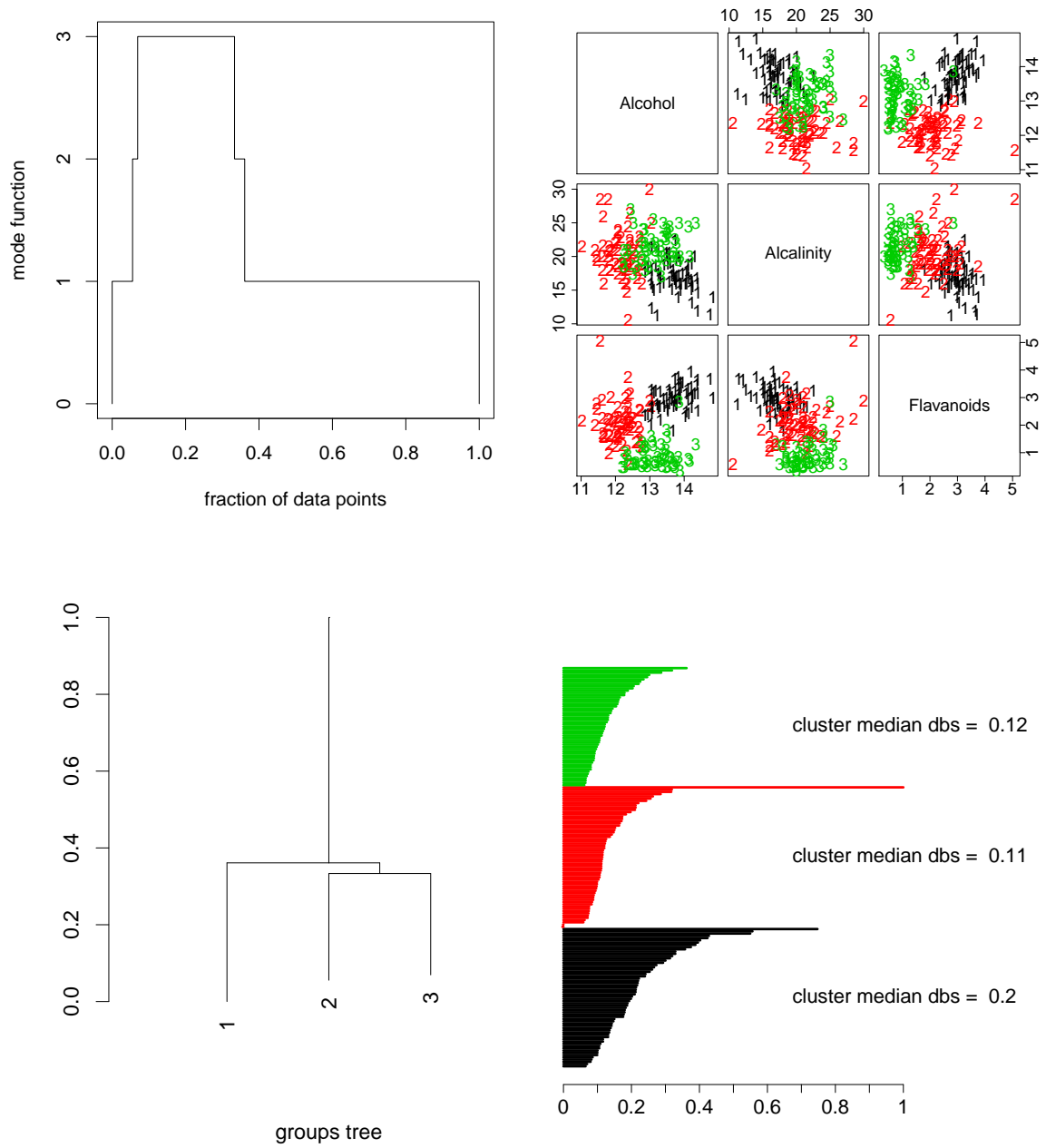


Figure 6: Result of the application of plot methods on objects of class `pdfCluster`. From the left: the mode function, the pairwise scatterplot, the cluster tree, and the *dbs* plot.

An S4 object of class "pdfCluster"

Call: pdfCluster(x = wineall, bwtype = "adaptive")

Initial groupings:

label	1	2	3	4	5	6	NA
count	5	2	10	4	5	3	149

Final groupings:

label	1	2	3	4	5	6
count	35	32	36	40	20	15

...

Note that six groups have now been produced. The identification of more than three clusters, when all thirteen variables of the `wine` data are used, is consistent with results of the application of other clustering methods; see, e.g., [McNicholas and Murphy \(2008\)](#). However, the pairwise aggregation of the six groups as  $\{1, 2\}$ ,  $\{3, 6\}$ ,  $\{4, 5\}$  essentially leads to the three actual cultivars, with a few misclassified points left:

```
> table(wine[, 1], groups(cl.wineall))
```

	1	2	3	4	5	6
Barolo	32	27	0	0	0	0
Grignolino	3	5	0	40	19	4
Barbera	0	0	36	0	1	11

A more accurate classification may be obtained by additionally increasing the value of `hmult`, an optional argument of function `pdfCluster()` that is multiplied by the bandwidths, thus controlling the level of smoothing of the density. When the data dimensionality is greater than 6, its default value is 1. However, for high-dimensional data, a larger value may be convenient, entailing to oversmooth the density function and then possibly aggregating the less prominent modes.

```
> cl.wineall.os <- pdfCluster(wineall, bwtype = "adaptive", hmult = 1.2)
```

```
> table(wine[, 1], groups(cl.wineall.os))
```

	1	2	3
Barolo	59	0	0
Grignolino	5	4	62
Barbera	0	47	1

A similar effect may be caused by relaxing the condition to set connections among points, while building the pairwise connection graph among the observations. In the current example, data dimensionality equal to 13 has prompted the use of the pairwise connection criterion, as may be seen from:

```
> cl.wineall@graph[1:2]
```

```
$type
[1] "pairs"
```

```
$lambda
[1] 0.1
```

In addition to the criterion adopted to build the network connection, the slot "graph" reports the tolerance threshold `lambda`. The third, not displayed, element of the slot contains the amplitude of the valleys detected by the evaluation of the density along the segments connecting all the pairs of observations. As already mentioned, setting  $\lambda$  to some 'large' value corresponds to the choice of being more tolerant in declaring the disconnection between pairs of observations, that possibly results in a smaller number of detected connected regions. To set a different value of `lambda`, either re-run the whole procedure:

```
> cl.wineall.10.25 <- pdfCluster(wineall, bwtype = "adaptive", lambda = 0.25)
```

or apply the method which the `pdfCluster()` function dispatches to objects of class "pdfCluster":

```
> cl.wineall.10.25 <- pdfCluster(cl.wineall, lambda = 0.25)
```

The latter choice does not re-run pairwise evaluation and it exploits, instead, computations saved in the slot "graph" of the object with class "pdfCluster". This speeds up the clustering procedure considerably, especially when the sample size is large. It can be seen that the larger value of  $\lambda$  results in aggregating clusters:

```
> table(groups(cl.wineall.10.25), groups(cl.wineall))
```

	1	2	3	4	5	6
1	34	32	0	1	0	0
2	0	0	35	0	1	11
3	1	0	0	38	1	0
4	0	0	1	1	18	4

## 4.2. Mixed variables

As they stand, density-based clustering methods may be applied to continuous data only. However, we illustrate here one possible way to circumvent this restriction.

Consider the data set `plantTraits` from package **cluster** (Maechler *et al.* 2005).

```
> library("cluster")
> data("plantTraits", package = "cluster")
```

These data refer to 136 plant species, described according to 31 morphological and reproductive attributes. Among them, only columns 1 to 3 correspond to quantitative variables; columns 4 to 11 are coded as ordered factors, while the remaining columns correspond to binary variables, which Maechler *et al.* (2005) further classify into symmetric (columns 12 and 13) and asymmetric (columns 14 to 31), thus considered when the two outcomes are not

equally important (e.g., describing the presence of a relatively rare attribute). This distinction was first advocated by [Kaufman and Rousseeuw \(1990, Chapter 1\)](#).

Applying standard methods for transforming categorical data into dichotomic variables would lead to an augmented computational complexity because of the increased number of variables, and still the hypothesis of continuity required by density-based clustering would not be fulfilled. Thus, in order to use all available information to cluster the data, we propose to proceed as follows: first, a dissimilarity matrix among the points is created, using criteria commonly employed in classical distance-based methods; next, from this matrix, a configuration of points in the  $d'$ -dimensional Euclidean space is sought by means of multidimensional scaling, for some given  $d'$ . In this way,  $d'$  numerical coordinates are obtained, to be passed to the clustering procedure.

It is inevitable that any recoding procedure of this sort involves some arbitrariness and the one just described is no exception. However, of the two steps involved, the first one is exactly that considered by classical hierarchical clustering techniques. The second step requires a few additional choices, especially the number  $d'$  of principal coordinates. A detailed exploration of these aspects is definitely beyond the scope of this paper and will be pursued elsewhere.

For the present illustrative purposes, we refer to the example reported in the documentation of the **cluster** package, that is we choose the Gower coefficient to measure the dissimilarities among the points, implemented in R by function `daisy`.

```
> dai.b <- daisy(plantTraits,
+               type = list(ordratio = 4:11, symm = 12:13, asymm = 14:31))
```

We then make use of classical multidimensional scaling ([Gower 1966](#)), which is computed by function `cmdscale()`, and we take the simple option of setting  $d' = 6$  as in [Maechler et al. \(2005\)](#).

```
> cmdsdai.b <- cmdscale(dai.b, k = 6)
```

With the new set of data, we may run the clustering procedure.

```
> cl.plntTr <- pdfCluster(cmdsdai.b)
> summary(cl.plntTr)
```

An S4 object of class "pdfCluster"

Call: pdfCluster(x = cmdsdai.b)

Initial groupings:

label	1	2	NA
count	23	2	111

Final groupings:

label	1	2
count	128	8

Groups tree (here 'h' denotes 'height'):

```
--[dendrogram w/ 1 branches and 2 members at h = 1]
  `--[dendrogram w/ 2 branches and 2 members at h = 0.171]
    |--leaf "1 "
    `--leaf "2 " (h = 0.167 )
```

The outcome indicates the presence of two clusters. However, the number of data points assigned to one of the two clusters is very small, and the associated cluster core is formed by two observations only. In these circumstances, selecting a global bandwidth to classify the lower density data seems to be more appropriate than using cluster specific bandwidths. This can be pursued with the following command:

```
> cl.plntTr.hc <- pdfCluster(cmdsdai.b, hcores = TRUE)
```

While the `plantTraits` data set does not include information about a true label class of the cases, it is interesting to note that the two identified groups roughly correspond to the aggregation of the clusters  $\{1, 3, 5, 6\}$  and  $\{2, 4\}$  identified by running a distance-based method and then cutting the dendrogram at six clusters, as suggested by Maechler *et al.* (2005).

```
> agn.trts <- agnes(dai.b, method = "ward")
> cutree6 <- cutree(agn.trts, k = 6)
> table(groups(cl.plntTr.hc), cutree6)
```

```
cutree6
  1  2  3  4  5  6
1 10 11 21  4 18 35
2  0 20  0 15  1  1
```

Note that selecting six principal coordinates places us at the threshold we defined to choose both the type procedure to find the connected level sets and the way of smoothing the density function. Since the threshold  $d = 6$  is merely indicative, in these circumstances it makes sense to check results deriving from a different setting of the arguments as, for example,

```
> cl.plntTr.hc.newset <- pdfCluster(cmdsdai.b, hcores = TRUE,
+                                   graphtype = "pairs", bwtype = "adaptive")
```

A variant procedure to handle mixed data would recode the categorical variables only, and run `pdfCluster()` on the set of data obtained by merging the principal coordinates so constructed and the original quantitative variables.

## References

- Azzalini A, Menardi G (2013). *R package pdfCluster: Cluster Analysis via Nonparametric Density Estimation*. R Package Version 1.0-1, with contribution of T. Rosolin up to version 0.0-13, URL <http://cran.r-project.org/package=pdfCluster>.
- Azzalini A, Torelli N (2007). “Clustering via Nonparametric Density Estimation.” *Statistics and Computing*, **17**(1), 71–80.

- Barber C, Habel K, Grasman R, Gramacy RB, Stahel A, Sterratt DC (2012). *Geometry: Mesh Generation and Surface Tesselation*. R Package Version 0.3-1, URL <http://cran.r-project.org/package=geometry>.
- Barber CB, Dobkin DP, Huhdanpaa H (1996). “The Quickhull Algorithm for Convex Hulls.” *ACM Transactions on Mathematical Software*, **22**(4), 469–483.
- Bivand R (2013). *spdep: Spatial Dependence: Weighting Schemes, Statistics and Models*. R Package Version 0.5-60, URL <http://cran.r-project.org/package=spdep>.
- Bowman AW, Azzalini A (1997). *Applied Smoothing Techniques for Data Analysis: the Kernel Approach with S-Plus Illustrations*. Oxford University Press, Oxford.
- Chambers JM (1998). *Programming with Data*. Springer-Verlag, New York.
- Ester M, Kriegel H, Sander J, Xu X (1996). “A Density Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise.” In *Proceedings of the 2Nd International Conference On Knowledge Discovery and Data Mining (KDD-96)*. Aaai Press, Portland.
- Forina M, Armanino C, Castino M, Ubigli M (1986). “Multivariate Data Analysis as a Discriminating Method of the Origin of Wines.” *Vitis*, **25**, 189–201.
- Gower JC (1966). “Some Distance Properties of Latent Root and Vector Methods Used in Multivariate Analysis.” *Biometrika*, **53**(3-4), 325–328.
- Hartigan JA (1975). *Clustering Algorithms*. John Wiley & Sons, New York.
- Hennig C (2010). *Fpc: Flexible Procedures for Clustering*. R Package Version 2.0-3, URL <http://cran.r-project.org/package=fpc>.
- Hubert L, Arabie P (1985). “Comparing Partitions.” *Journal of Classification*, **2**(1), 193–218.
- Kaufman L, Rousseeuw PJ (1990). *Finding Groups in Data: an Introduction to Cluster Analysis*. John Wiley & Sons, New York.
- Lin TI, Lee JC, Hsieh WJ (2007). “Robust Mixture Modeling Using the Skew  $t$  Distribution.” *Statistics and Computing*, **17**(2), 81–92.
- Maechler M, Rousseeuw P, Struyf A, Hubert M (2005). “**Cluster**: Cluster Analysis Basics and Extensions.” R Package Version 1.14-2, URL <http://cran.r-project.org/package=cluster>.
- McLachlan GJ, Peel D (2000). *Finite Mixture Models*. John Wiley & Sons, New York.
- McNicholas PD, Murphy TB (2008). “Parsimonious Gaussian Mixture Models.” *Statistics and Computing*, **18**(3), 285–296.
- Menardi G (2011). “Density-based Silhouette Diagnostics for Clustering Methods.” *Statistics and Computing*, **21**(3), 295–308.
- Menardi G, Azzalini A (2013). “An Advancement in Clustering via Nonparametric Density Estimation.” *Statistics and Computing*. doi:10.1007/s11222-013-9400-x. To appear.

- Okabe A, Boots BN, Sugihara K (1992). *Spatial Tessellations: Concepts and Applications of Voronoi Diagrams*. John Wiley & Sons, New York.
- R Core Team (2012). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org/>.
- Ray S, Lindsay BG (2005). “The Topography of Multivariate Normal Mixtures.” *The Annals of Statistics*, **33**(5), 2042–2065.
- Rousseeuw PJ (1987). “Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis.” *Journal of Computational Applied Mathematics*, **20**(1), 53–65.
- Silverman BW (1986). *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, London.
- Stuetzle W (2003). “Estimating the Cluster Tree of a Density by Analyzing the Minimal Spanning Tree of a Sample.” *Journal of Classification*, **20**(1), 25–47.
- Stuetzle W, Nugent R (2010). “A Generalized Single Linkage Method for Estimating the Cluster Tree of a Density.” *Journal of Computational and Graphical Statistics*, **19**(2), 397–418.
- Wishart D (1969). *Mode Analysis: a Generalization of Nearest Neighbor which Reduces Chaining Effects*, pp. 282–308. Numerical Taxonomy. Cole A. J., Ed. Academic Press, London.
- Wong AM, Lane T (1983). “The  $k$ -th Nearest Neighbour Clustering Procedure.” *Journal of the Royal Statistical Society B*, **45**(3), 362–368.

**Affiliation:**

Adelchi Azzalini  
Dipartimento di Scienze Statistiche  
Università di Padova, Italia  
E-mail: [azzalini@stat.unipd.it](mailto:azzalini@stat.unipd.it)  
URL: <http://azzalini.stat.unipd.it>

Giovanna Menardi  
Dipartimento di Scienze Statistiche  
Università di Padova, Italia  
E-mail: [menardi@stat.unipd.it](mailto:menardi@stat.unipd.it)  
URL: <http://homes.stat.unipd.it/menardi>