

not optimally, according to the adopted famous algorithm of this category is [Lloy 82, Duda 73].

gorithms, follow Bayesian classification is assigned to the cluster C_i for which (probability) is maximum. These probabilities are defined on the basis of an appropriately defined optimization task. where a vector belongs to a specific

gorithms. In this case we measure the x to belong to a cluster C_i .

Instead of determining the clusters in advance, these algorithms adjust iteratively where clusters lie. These algorithms, based on a cost function optimization philosophy, are different from the above algorithms. All the algorithms use cluster representatives and the goal is to find an optimal way. In contrast, boundary placement algorithms place boundaries before the decision to treat these algorithms with algorithms to be discussed next. Some special clustering techniques that belong to these categories. These include:

gorithms. These algorithms provide clustering without having to consider a fixed number m of clusters, and for this reason, their computational burden is

These algorithms use an initial population and iteratively generate new populations in better clusterings than those of the previous population according to a prespecified criterion.

These are methods that guarantee convergence in probability to the globally optimal solution with respect to a prespecified criterion, at the expense of a large number of iterations.

gorithms. These algorithms treat the feature vector x as a (multidimensional) random variable x . The only accepted assumption that regions of high probability density correspond to regions of increased probability density function (pdf) of x . The pdf may highlight the regions where

- Competitive learning algorithms.* These are iterative schemes that do not employ cost functions. They produce several clusterings and they converge to the most “sensible” one, according to a distance metric. Typical representatives of this category are the *basic competitive learning scheme* and the *leaky learning algorithm*.
- Algorithms based on morphological transformation techniques.* These algorithms use morphological transformations in order to achieve better separation of the involved clusters.

12.3 SEQUENTIAL CLUSTERING ALGORITHMS

In this section we describe a basic sequential algorithmic scheme (BSAS) (which is a generalization of that discussed in [Hall 67]) and we also give some variants of it. First, we consider the case where all the vectors are presented to the algorithm only once. The number of clusters is not known a priori in this case. In fact, new clusters are created as the algorithm evolves.

Let $d(x, C)$ denote the distance (or dissimilarity) between a feature vector x and a cluster C . This may be defined by taking into account either all vectors of C or a representative vector of it (see Chapter 11). The user-defined parameters required by the algorithmic scheme are the threshold of dissimilarity Θ and the maximum allowable number of clusters, q . The basic idea of the algorithm is the following: As each new vector is considered, it is either assigned to an existing cluster or assigned to a newly created cluster, depending on its distance from the already formed ones. Let m be the number of clusters that the algorithm has created up to now. Then the algorithmic scheme may be stated as:

Basic Sequential Algorithmic Scheme (BSAS)

- $m = 1$
- $C_m = \{x_1\}$
- For $i = 2$ to N
 - Find C_k : $d(x_i, C_k) = \min_{1 \leq j \leq m} d(x_i, C_j)$
 - If $(d(x_i, C_k) > \Theta)$ AND $(m < q)$ then
 - * $m = m + 1$
 - * $C_m = \{x_i\}$
 - Else
 - * $C_k = C_k \cup \{x_i\}$
 - * Where necessary, update representatives.²
 - End { if }
- End { For }

²This statement is activated in the cases where each cluster is represented by a single vector. For example, if each cluster is represented by its mean vector, this must be updated each time a new vector becomes a member of the cluster.

Different choices of $d(x, C)$ lead to different algorithms and any of the measures introduced in Chapter 11 can be employed. When C is represented by a single vector, $d(x, C)$ becomes

$$d(x, C) = d(x, m_C) \quad (12.4)$$

where m_C is the representative of C . In the case in which the mean vector is used as a representative, the updating may take place in an iterative fashion, that is,

$$m_{C_k}^{new} = \frac{(n_{C_k}^{new} - 1)m_{C_k}^{old} + x}{n_{C_k}^{new}} \quad (12.5)$$

where $n_{C_k}^{new}$ is the cardinality of C_k after the assignment of x to it and $m_{C_k}^{new}$ ($m_{C_k}^{old}$) is the representative of C_k after (before) the assignment of x to it (Problem 12.2).

Algorithms where each cluster is represented by a single vector are said to be based on *global clustering criteria* [Jain 88] and algorithms where all vectors are used for its representation are said to be based on *local clustering criteria*.

It is not difficult to realize that the order in which the vectors are presented to the BSAS plays an important role in the clustering results. Different presentation ordering may lead to totally different clustering results, in terms of the number of clusters as well as the clusters themselves (see Problem 12.3).

Another important factor affecting the result of the clustering algorithm is the choice of the threshold Θ . This value directly affects the number of clusters formed by BSAS. If Θ is too small, unnecessary clusters will be created. On the other hand, if Θ is too large a smaller than appropriate number of clusters will be created. In both cases, the number of clusters that best fits the data set is missed.

If the number q of the maximum allowable number of clusters is not constrained, we leave to the algorithm to "decide" about the appropriate number of clusters. Consider for example Figure 12.1, where three compact and well-separated clusters are formed by the points of X . If the maximum allowable number of clusters is set equal to two, the BSAS algorithm will be unable to discover three clusters. Probably, in this case the two rightmost groups of points will form a single cluster. On the other hand, if q is unconstrained the BSAS algorithm will probably form three clusters (with an appropriate choice of Θ), at least for the case in which the mean vector is used as a representative. However, constraining q becomes necessary when dealing with implementations where the available computational resources are limited. In the next subsection a simple technique is given for determining the number of clusters.³

³This problem is also treated in Chapter 16.

ms and any of the measures
is represented by a single

(12.4)

which the mean vector is
in an iterative fashion, that

(12.5)

of x to it and $m_{C_k}^{new}$ ($m_{C_k}^{old}$)
of x to it (Problem 12.2).
single vector are said to be
thms where all vectors are
clustering criteria.
e vectors are presented to
its. Different presentation
in terms of the number of
12.3).

clustering algorithm is the
number of clusters formed
be created. On the other
of clusters will be created.
ata set is missed.
clusters is not constrained,
prie number of clusters.
and well-separated clusters
able number of clusters is
to discover three clusters.
s will form a single cluster.
m will probably form three
he case in which the mean
ing q becomes necessary
e computational resources
given for determining the



FIGURE 12.1: Three clusters are formed by the feature vectors. When q is constrained to a value less than 3, the BSAS algorithm will not be able to reveal them.

Remarks

- The BSAS scheme may be used with similarity instead of dissimilarity measures with appropriate modification; that is, the *min* operator is replaced by *max*.
- It turns out that BSAS, with point cluster representatives, favors compact clusters. Thus, it is not recommended if there is strong evidence that other types of clusters are present.
- The preceding algorithm is closely related to the algorithm implemented by the ART2 (adaptive resonance theory) neural architecture [Carp 87, Burk 91].

12.3.1 Estimation of the Number of Clusters

In this subsection, a simple method is described for determining the number of clusters (see also Chapter 16). The method is suitable for BSAS as well as other algorithms, for which the number of clusters is not required as an input parameter. In what follows, BSAS(Θ) denotes the BSAS algorithm with a specific threshold of dissimilarity Θ .

- For $\Theta = a$ to b step c
 - Run s times the algorithm BSAS(Θ), each time presenting the data in a different order.
 - Estimate the number of clusters, m_{Θ} , as the most frequent number resulting from the s runs of BSAS(Θ).
- Next Θ

The values a and b are the minimum and maximum dissimilarity levels among all pairs of vectors in X , that is, $a = \min_{i,j=1,\dots,N} d(x_i, x_j)$ and $b = \max_{i,j=1,\dots,N} d(x_i, x_j)$. The choice of c is directly influenced by the choice of $d(x, C)$. As far as the value of s is concerned, the greater the s , the larger the statistical sample

and, thus, the higher the accuracy of the results. In the sequel, we plot the number of clusters m_Θ versus Θ . This plot has a number of flat regions. We estimate the number of clusters as the number that corresponds to the largest flat region. It is expected that at least for the case in which the vectors form well-separated compact clusters, this is the desired number. Let us explain this argument intuitively. Suppose that the data form two compact and well-separated clusters C_1 and C_2 . Let the minimum distance between two vectors in C_1 (C_2) be r_1 (r_2) and suppose that $r_1 < r_2$. Also let r ($> r_2$) be the minimum among all distances $d(x_i, x_j)$, with $x_i \in C_1$ and $x_j \in C_2$. It is clear that for $\Theta \in [r_2, r - r_2]$, the number of clusters created by BSAS is 2. In addition, if $r \gg r_2$, the interval has a big range, and thus it corresponds to a large flat region in the plot of m_Θ versus Θ . Example 12.2 illustrates the idea.

Example 12.2. Consider two 2-dimensional Gaussian distributions with means $[0, 0]^T$ and $[20, 20]^T$, respectively. The covariance matrices are $\Sigma = 0.5I$ for both distributions, where I is the 2×2 identity matrix. Generate 50 points from each distribution (Figure 12.2a). The number of underlying clusters is 2. The plot resulting from the application of the previously described procedure is shown in Figure 12.2b, with $a = \min_{x_i, x_j \in X} d_2(x_i, x_j)$, $b = \max_{x_i, x_j \in X} d_2(x_i, x_j)$, and $c \simeq 0.3$. It can be seen that the biggest flat region corresponds to the number 2, which is the number of underlying clusters.

In the foregoing procedure, we have implicitly assumed that the feature vectors do form clusters. If this is not the case, the method is useless. Methods that deal with the problem of discovering whether any clusters exist are discussed in Chapter 16. Moreover, if the vectors form compact clusters, which are not well separated, the

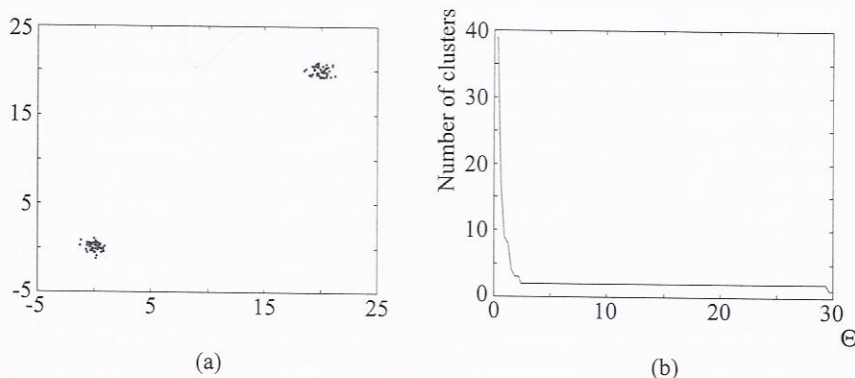


FIGURE 12.2: (a) The data set. (b) The plot of the number of clusters versus Θ . It can be seen that for a wide range of values of Θ , the number of clusters, m , is 2.

procedure may g
 Θ to contain wi

In some cases
 that correspond
 If, for example,
 other and away
 second flattest f
 the three-cluster

12.4 A MOD

As has already
 x is assigned t
 a decision for t
 is determined a
 BSAS, which v
 The cost we p
 the algorithm.
 involves the det
 of X to them. I
 a second time t
 MBSAS may b

Modified Basi

Cluster Dete

- $m = 1$
 - $C_m = \{x\}$
 - For $i = 2$
 - Find
 - If (d
 - *
 - *
 - End
 - End { For
- Pattern Clas
- For $i = 1$
 - If x_i
 - *
 - *
 - *
 - End
 - End { For

of the number of clusters versus
es of Θ , the number of clusters, m .

- For $i = 1$ to N
 - If \mathbf{x}_i has not been assigned to a cluster, then
 - * Find $C_k: d(\mathbf{x}_i, C_k) = \min_{1 \leq j \leq m} d(\mathbf{x}_i, C_j)$
 - * $C_k = C_k \cup \{\mathbf{x}_i\}$
 - * Where necessary, update representatives.
 - End { if }
- End { For }

The number of clusters is determined in the first phase and then it is frozen. Thus, the decision taken during the second phase for each vector takes into account all clusters.

When the mean vector of a cluster is used as its representative, the appropriate cluster representative has to be adjusted using Eq. (12.5), after the assignment of each vector in a cluster.

Also, as it was the case with BSAS, MBSAS is sensitive to the order in which the vectors are presented.

Finally, it must be stated that, after minor modifications, MBSAS may be used when a similarity measure is employed (see Problem 12.7).

12.5 A TWO-THRESHOLD SEQUENTIAL SCHEME

As already pointed out, the results of BSAS and MBSAS are strongly dependent on the order in which the vectors are presented to the algorithm, as well as on the value of Θ . Improper choice of Θ may lead to meaningless clustering results. One way to overcome these difficulties is to define a "gray" region (see [Trah 89]). This is achieved by employing two thresholds, Θ_1 and $\Theta_2 (> \Theta_1)$. If the dissimilarity level $d(x, C)$ of a vector x from its closest cluster C is less than Θ_1 , x is assigned to C . If $d(x, C) > \Theta_2$, a new cluster is formed and x is placed in it. Otherwise, if $\Theta_1 \leq d(x, C) \leq \Theta_2$, there exists uncertainty and the assignment of x to a cluster will take place at a later stage. Let $clas(x)$ be a flag that indicates whether x has been classified (1) or not (0). Again, we denote by m the number of clusters that have been formed up to now. In the following, we assume no bounds to the number of clusters (i.e., $q = N$). The algorithmic scheme is:

The Two-Threshold Sequential Algorithmic Scheme (TTSAS)

```

m = 0
clas(x) = 0,  ∀x ∈ X
prev_change = 0
cur_change = 0
exists_change = 0

```

While (there exists at least one feature vector x with $clas(x) = 0$) do

- For $i = 1$ to N
 - if $clas(x_i) = 0$ AND it is the first in the new while loop AND $exists_change = 0$ then
 - * $m = m + 1$
 - * $C_m = \{x_i\}$

```

* clas
* cur
—Else if cl
* Find
* if d(
.
.
.
* else i
.
.
.
* End
—Else if cl
* cur
—End {If}

```

- End { For }
- exists_change
- prev_change
- cur_change =

End {While}

The *exists_change* is classified at the current pass. This is achieved by checking if any vector was classified up to the previous pass that is, no vector has been unclassified vector in the current pass.

The first if condition forces the first unassigned vector during the last pass to be assigned at a new cluster.

However, in practice, it should be pointed out that the previous two schemes are not efficient. Moreover, since the

first phase and then it is frozen. Thus, for each vector takes into account all

as its representative, the appropriate Eq. (12.5), after the assignment of

AS is sensitive to the order in which

modifications, MBSAS may be used Problem 12.7).

IAL SCHEME

and MBSAS are strongly dependent to the algorithm, as well as on the meaningless clustering results. One a "gray" region (see [Trah 89]). This and $\Theta_2(> \Theta_1)$. If the dissimilarity cluster C is less than Θ_1 , x is assigned and x is placed in it. Otherwise, if and the assignment of x to a cluster is a flag that indicates whether x has been by m the number of clusters that we assume no bounds to the number is:

eme (TTSAS)

x with $clas(x) = 0$ do

in the new while loop AND

```

* clas( $x_i$ ) = 1
* cur_change = cur_change + 1
—Else if clas( $x_i$ ) = 0 then
* Find  $d(x_i, C_k) = \min_{1 \leq j \leq m} d(x_i, C_j)$ 
* if  $d(x_i, C_k) < \Theta_1$  then
    ·  $C_k = C_k \cup \{x_i\}$ 
    · clas( $x_i$ ) = 1
    · cur_change = cur_change + 1
* else if  $d(x_i, C_k) > \Theta_2$  then
    ·  $m = m + 1$ 
    ·  $C_m = \{x_i\}$ 
    · clas( $x_i$ ) = 1
    · cur_change = cur_change + 1
* End {If}
—Else if clas( $x_i$ ) = 1 then
    * cur_change = cur_change + 1
—End {If}

• End {For}
• exists_change = |cur_change - prev_change|
• prev_change = cur_change
• cur_change = 0
    
```

End {While}

The *exists_change* checks whether there exists at least one vector that has been classified at the current pass on X (i.e., the current iteration of the while loop). This is achieved by comparing the number of vectors that have been classified up to the current pass on X , *cur_change*, with the number of vectors that have been classified up to the previous pass on X , *prev_change*. If *exists_change* = 0, that is, no vector has been assigned to a cluster during the last pass on X , the first unclassified vector is used for the formation of a new cluster.

The first *if* condition in the *For* loop ensures that the algorithm terminates after N passes on X (N executions of the while loop) at the most. Indeed, this condition forces the first unassigned vector to a new cluster when no vector has been assigned during the last pass on X . This gives a way out to the case in which no vector has been assigned at a given circle.

However, in practice, the number of required passes is much less than N . It should be pointed out that this scheme is almost always at least as expensive as the previous two schemes, because in general it requires at least two passes on X . Moreover, since the assignment of a vector is postponed until enough information