INEB-PSI Technical Report 2005-1


# Human Clustering on Bi-dimensional Data: An Assessment

by

Jorge M. Santos, J. P. Marques de Sá
jms@isep.ipp.pt

# Contents

# 1    Introduction

Data clustering performed by humans is characterized by a high variability of solutions for non-trivial data sets. The complexity and subjectivity involved in the clustering process are highly related to the personal experience and sometimes to knowledge about the problem domain. Clustering solutions may depend on a variety of features perceived in the data set. Figure 1 illustrates some of the features that seem to have a main role in guiding human solutions to clustering. They are as follows:

**Connectedness** – This is probably the most basic feature leading us to join points into clusters whenever connecting paths are perceived. This feature is valued in the data set of Figure 1a when a human "sees" one cluster instead of two.
**Structuring direction** – This feature leads us to "see" the two arms of the cross in Figure 1b instead of only one cluster. Humans are good at perceiving structuring directions in data set graphs, independently of those directions being straight or curved lines.
**Structuring density** - This feature leads us to "see" two clusters in Figure 1c instead of only one.
**Structuring morphology** - This feature leads us to "see" two clusters in Figure 1d instead of only one, deciding differently of the similar figure 1a. The reason is that, contrary to Figure 1a, we now identify the bulging out wart of Figure 1d with a known form.
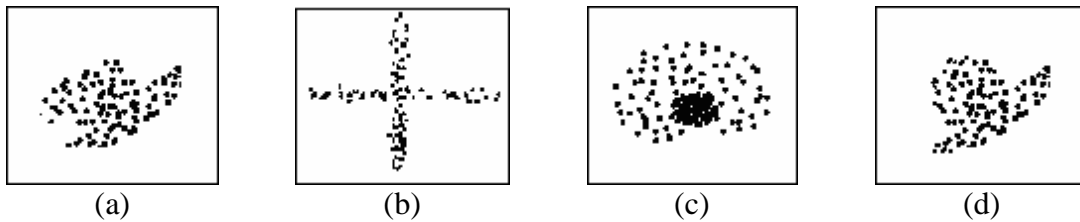


Figure 1. Clustering features: a) connectedness; b) structuring direction; c) structuring density; d) structuring morphology.

How much influence have these features in the clustering process? How do they interplay? In order to obtain some knowledge about these issues we performed a variety of 2D data clustering experiments involving children and adults. The reason to involve children in clustering experiments is related to the fact that we expected in this way to discriminate (and characterize) basic clustering skills present in children from more advanced skills present in adults. Based on the experimental results we were able to extract a few guidelines on the human approach to data clustering.

# 2    Clusters experiments

We performed tests involving several individuals (including children) in order to grasp, based on the results, the mental process of data clustering. We made the experiment with 37 individuals, 17 of them children (6-7 years old), 15 adults with no knowledge about clustering and 5 adults with some knowledge of clustering problems. The experiments were performed with the bi-dimensional data sets shown in Figure 2 and Figure 3. All data sets were manually drawn and we tried to create different situations using examples similar to those usually seen in clustering-related works and others created by us.

We have presented to the individuals all the data sets in the same order as in Figures 2 and 3 and they were asked to circle the possible groups of points in each data set. We haven't given any other explanation or made any comment on the way they should perform the experiment. We just said that in each figure some groups of points could exist, or not, and if they thought they existed they should circle them with a line.

A few similar data sets with small differences among them were deliberately included in order to appreciate how small differences influence the clustering solutions. Examples of such data sets are the pairs (b-f) and (p-aa).

(a) Data set "a"

(b) Data set "b"

(c) Data set "c"

(d) Data set "d"

(e) Data set "e"

(f) Data set "f"

(g) Data set "g"

(h) Data set "h"

(i) Data set "i"

(j) Data set "j"

(k) Data set "k"

(l) Data set "l"

(m) Data set "m"

(n) Data set "n"

(o) Data set "o"

**Figure 2: Data sets I.**

(a) Data set "p"

(b) Data set "q"

(c) Data set "r"

(d) Data set "s"

(e) Data set "t"

(f) Data set "u"

(g) Data set "v"

(h) Data set "w"

(i) Data set "x"

(j) Data set "y"

(k) Data set "z"

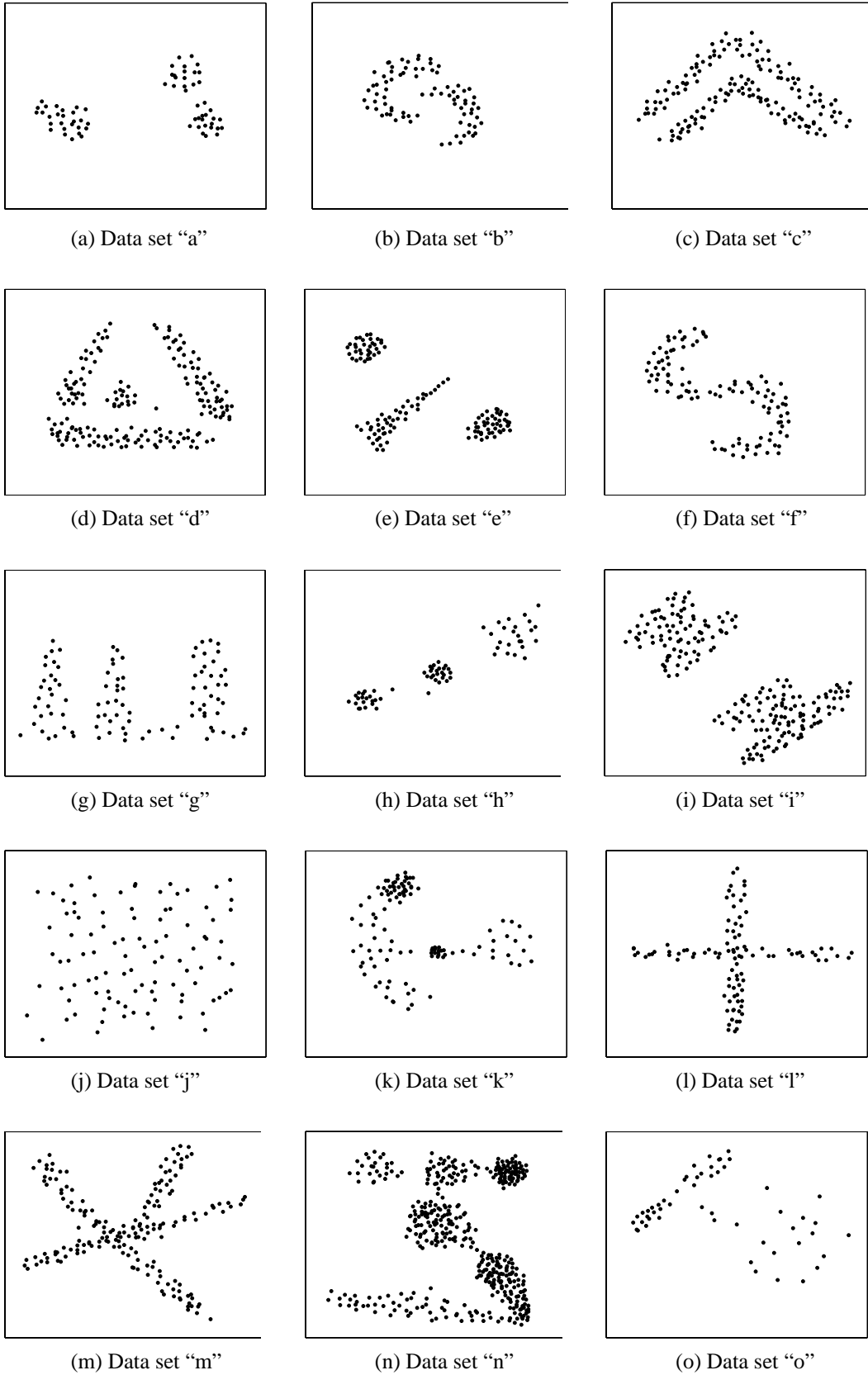(l) Data set "aa"

(m) Data set "bb"
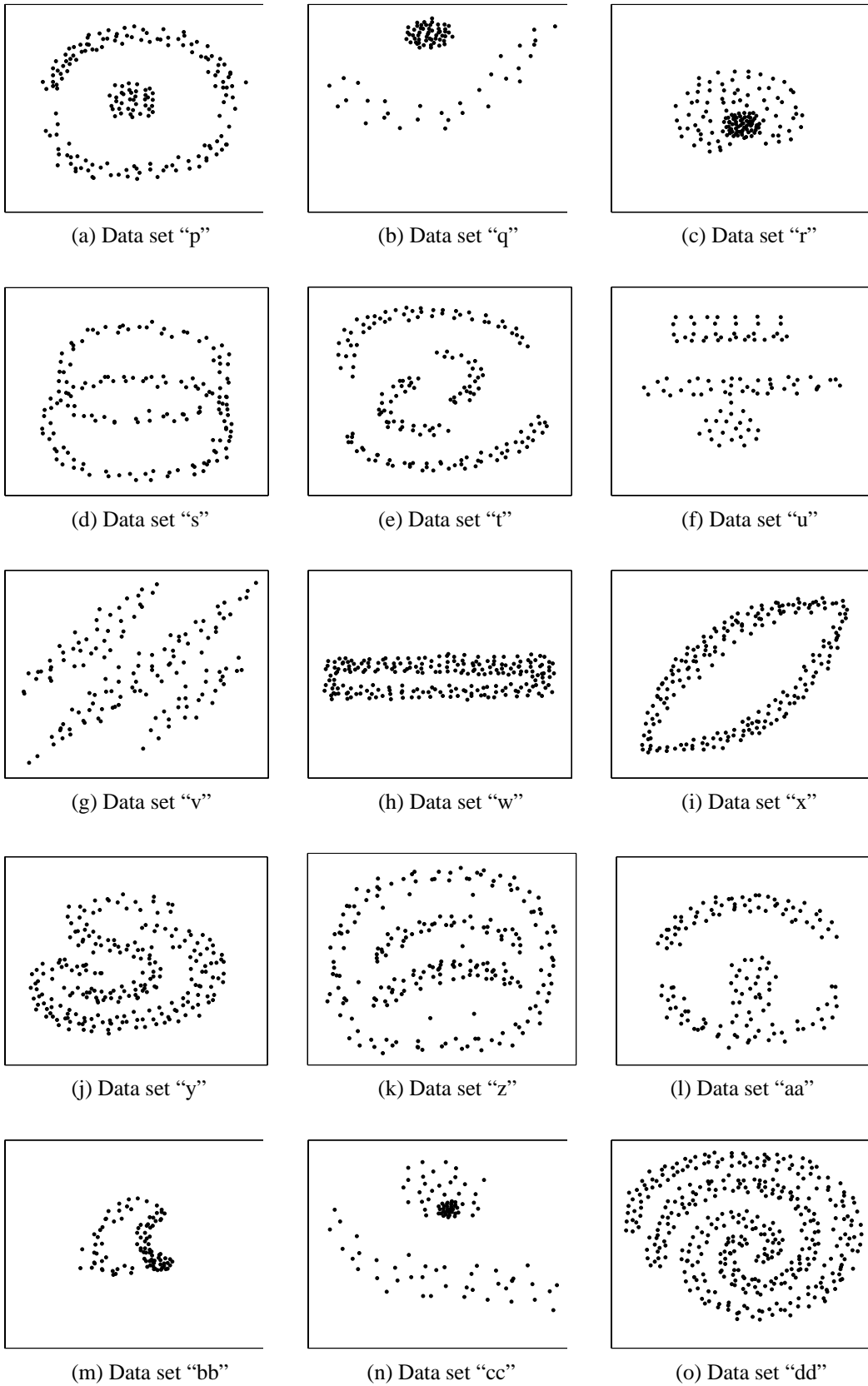
(n) Data set "cc"

(o) Data set "dd"

Figure 3: Data sets II.

5

# 3 Results

## 3.1 Global View

In this section, we present the results of the experiments in a global perspective.

The clustering solutions proposed by the adults are summarized in Table 1. The clustering solutions proposed by the children are summarized in Table 2 following the same labelling as for the adults. In the labelling of the solutions, we used the label "Others" to designate a group of various solutions different from the most occurring ones, labelled with numerals.

**Table 1: Experimental results with adults.**

| | Data sets | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Solutions** | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z | aa | bb | cc | dd |
| 1 | 19 | 20 | 17 | 9 | 19 | 15 | 13 | 20 | 16 | 19 | 2 | 10 | 6 | 14 | 5 | 12 | 20 | 14 | 4 | 19 | 10 | 19 | 9 | 8 | 16 | 16 | 11 | 16 | 10 | 17 |
| 2 | 1 | | 2 | 9 | | 5 | 6 | | 3 | | 4 | 5 | 4 | | 11 | 7 | | 6 | 14 | | 7 | | 5 | 6 | 3 | | 8 | 4 | 8 | 3 |
| 3 | | | | | | | | | | | 4 | | 6 | | | | | | | | | | 5 | 6 | | | | | | |
| 4 | | | | | | | | | | | 5 | | | | | | | | | | | | | | | | | | | |
| Others | | 1 | 1 | 1 | | 1 | | 1 | 1 | | 5 | 5 | 4 | 6 | 4 | 1 | | 2 | 1 | 3 | 1 | 1 | | 1 | 4 | 1 | | 2 | | |

**Table 2: Experimental results with children.**

| | Data sets | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Solutions** | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z | aa | bb | cc | dd |
| 1 | 14 | 12 | 13 | 4 | 14 | 2 | 6 | 14 | 14 | 9 | 1 | 2 | 9 | 5 | 5 | 10 | 13 | 9 | 4 | 12 | 9 | 12 | 2 | 3 | 8 | 10 | 8 | 5 | 1 | 9 |
| 2 | 1 | | 1 | 10 | | 13 | 7 | | | | 4 | 6 | 1 | | 3 | 3 | | | | 2 | | 2 | 3 | 3 | | 2 | 6 | 11 | 3 | |
| 3 | | | | | | | | | | | | | | 1 | | | | | | | | 9 | 9 | | | | | | | |
| 4 | | | | | | | | | | | | 1 | | | | | | | | | | | | | | | | | | |
| Others | 2 | 5 | 3 | 3 | 3 | 2 | 3 | 2 | 2 | 7 | 10 | 9 | 6 | 10 | 8 | 3 | 3 | 7 | 12 | 4 | 5 | 4 | 1 | 1 | 5 | 6 | 5 | 5 | 4 | 4 |

A glance at Tables 1 and 2 immediately shows that the solutions proposed by the adults are more consistent, exhibiting fewer solutions for each data set than the ones proposed by the children (6-7 years). Detailed observation of the children solutions revealed that a large percentage of children build clusters based on a small number of points. It seems that they focus on more local regions giving particular attention to small groups. An example of such behavior is shown in Figure 4.

During the labelling process we only considered "well-grown" clusters proposed in the solution, disregarding very small clusters (up to 2 points). This often happened with solutions proposed by children. An example of this situation is the one depicted in Figure 4a. In this case, we considered the proposed 3-cluster solution like the one shown in Figure 20c.
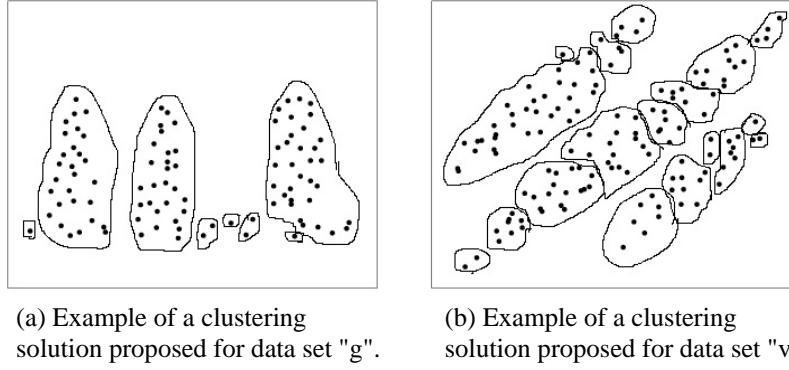
(a) Example of a clustering
solution proposed for data set "g".

(b) Example of a clustering
solution proposed for data set "v".

**Figure 4: Children usually consider the existence of small clusters.**

## 3.2 Detailed View

In this section, we present a detailed view of the results together with statistical assessment and some comments.

In order to understand in detail the clustering process, we have divided the data sets into several types. Type A: data sets with well-separated clusters; Type B: data sets with different point densities; Type C: data sets with crossing clusters; Type D: data sets with nested clusters; Type F: data sets with spiral-shaped clusters; Type E: other data sets.

In the next subsections, we take a closer view to each group of data sets and make some comments about the proposed solutions. We also present analyses of the clustering results with the following statistical tests: $X^2$ test for goodness of fit to a postulated distribution; $X^2$ test for independence between the Age variable (two categories: adults and children) and Solution variable (categories to be presented in the subsections). The independence test is complemented with Cramer's V measure of association for nominal variables. The level of significance of the tests was set at 5%. The usual conditions of validity of the $X^2$ tests were taken into consideration: for one degree of freedom no expected value below 5; for more than one degree of freedom no expected value below 1 and no more than 20% of the expected values below 5. When these conditions were not met the tests were not applied.

### 3.2.1 Type A: Data sets with well-separated clusters

In this subsection, we analyze the group of data sets with well-separated clusters. This group is constituted by the set of data sets {a, b, c, d, e, h, i, q, t, v}.

For these data sets there is basically a unique solution shown in Figure 5 proposed by a large majority of adults and children. Connectedness and sometimes structuring direction (data sets b, c and d) are the main features valued in this unique clustering solution. The results for these data sets are shown in Table 3.

**Table 3: Experimental results with adults and children for well-separated clusters.**

| *adults* Solutions | a | b | c | d | e | h | i | q | t | v |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 19 | 20 | 17 | 18 | 19 | 20 | 16 | 20 | 19 | 19 |
| 2 | 1 | | 2 | | | | 3 | | | |
| Others | | | 1 | 1 | 1 | | 1 | | 1 | 1 |

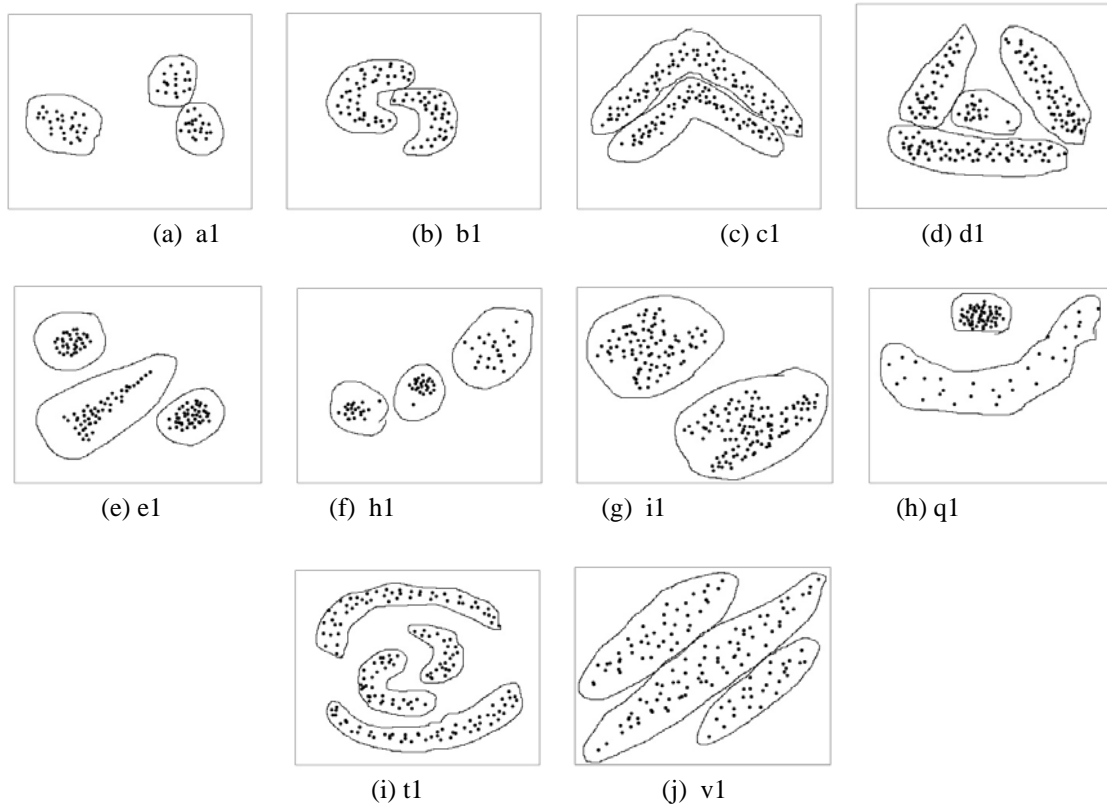| *children* Solutions | a | b | c | d | e | h | i | q | t | v |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 14 | 12 | 13 | 14 | 14 | 14 | 14 | 13 | 12 | 12 |
| 2 | 1 | | 1 | | | | | | | |
| Others | 2 | 5 | 3 | 3 | 3 | 2 | 2 | 3 | 4 | 4 |

7

**Figure 5: The solutions proposed for data sets a, b, c, e, h, i, q, t and v.**

The $X^2$ test for independence was performed for a Solution variable with two categories: major solutions; minor solutions. Thus, the following 2×2 table was used:

|  | Adults | Children |
|---|---|---|
| Major solutions | 187 | 132 |
| Minor solutions | 13 | 33 |

As expected, the independence hypothesis was rejected with p≈0. The Cramer V of the association is low (V=0.2).

*3.2.2*   Type B: Data sets with different point densities

In this subsection we analyze the data sets exhibiting clusters with different point densities. This group is constituted by the set of data sets {k, n, o, r, bb, cc}. For data set "n" there is basically a unique proposed solution, shown in Figure 6.
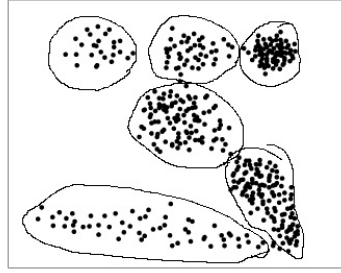
**Figure 6: The solution proposed for data set "n".**

The other Type B data sets are discussed in the following subsections.

### 3.2.2.1 Data set "k"

This data set was probably the one with the largest number of different proposed solutions (see Figure 7). Apart from the 4 considered solutions (k1 to k4) the adults proposed 5 more different solutions. The reasons for this variability can be attributed to the existence of different density regions and the peculiar structure of the data.
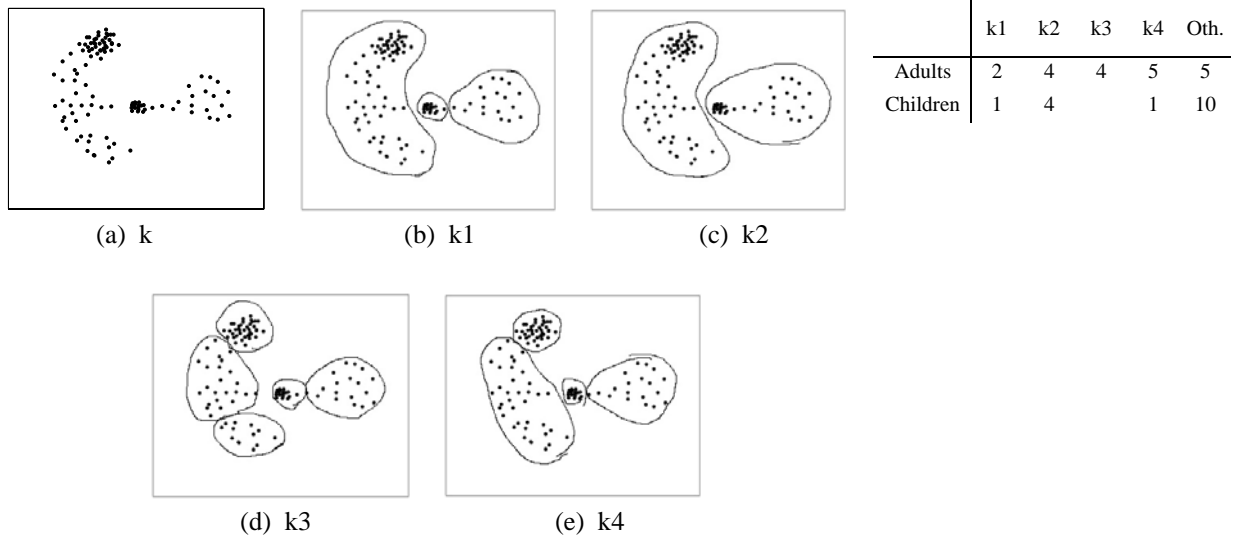


|          | k1 | k2 | k3 | k4 | Oth. |
|----------|----|----|----|----|------|
| Adults   | 2  | 4  | 4  | 5  | 5    |
| Children | 1  | 4  |    | 1  | 10   |

(a) k  (b) k1  (c) k2

(d) k3  (e) k4

**Figure 7: The solutions proposed for data set "k".**

Solution "k4" is the most significant for adults and solution "k2" for children and adults. We think that solution "k3" was suggested by adults based on the symmetry of the data set. We can see that solution "k2" gives more importance to the global structure and that solution "k4" gives relevance to the local structure of the data. Therefore, this data set suggests that children do not value the density feature to the point of sacrificing local connectedness

The $X^2$ test for goodness of fit lead us to accept the uniformity hypothesis (equiprobability of the solutions) for the adults (p=0.66). The $X^2$ test for independence, for a Solution variable with two categories ("regular clusters", "other clusters"≡"non-regular clusters"), lead us to reject the independence hypothesis (p=0.05). The Cramer V

is moderate (V=0.38). The rejection of the independence hypothesis is related to the fact that there is a regular vs. non-regular balance for children which is the opposite for adults.
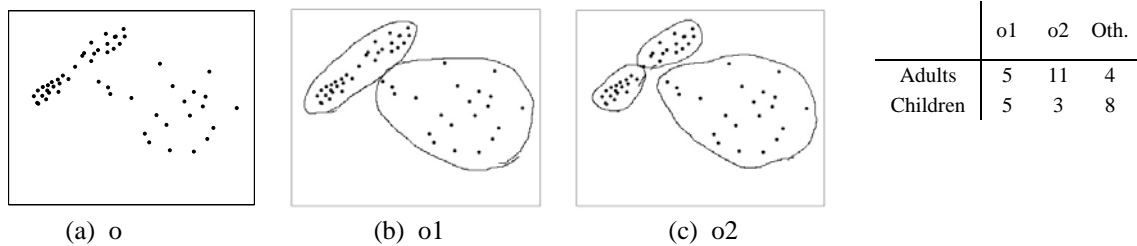
### 3.2.2.2 Data set "o"



| | o1 | o2 | Oth. |
|---|---|---|---|
| Adults | 5 | 11 | 4 |
| Children | 5 | 3 | 8 |

(a) o  (b) o1  (c) o2

**Figure 8: The solutions proposed for data set "o".**

For the data set "o" the solution "o2" was proposed by the majority of the adults. However, the $X^2$ test for goodness of fit lead us to accept the uniformity hypothesis for the adults (p=0.13) and for the children (p=0.26). Therefore, the behaviour of adults and children was quite similar in this case. The $X^2$ test for independence, for a Solution variable with the three categories as above, lead us to reject the independence hypothesis (p=0.056). The Cramer V is moderate (V=0.39). These findings further support the idea of identical behaviour of adults and children when connectedness prevails over light differences of point density.

### 3.2.2.3 Data set "r"

This data set was produced in order to try to percept the influence of a high density region situated inside a low density region. The performed tests indicate that this high density region is considered by the majority of the individuals, both adults (70%) and children (56%), as a separate cluster.
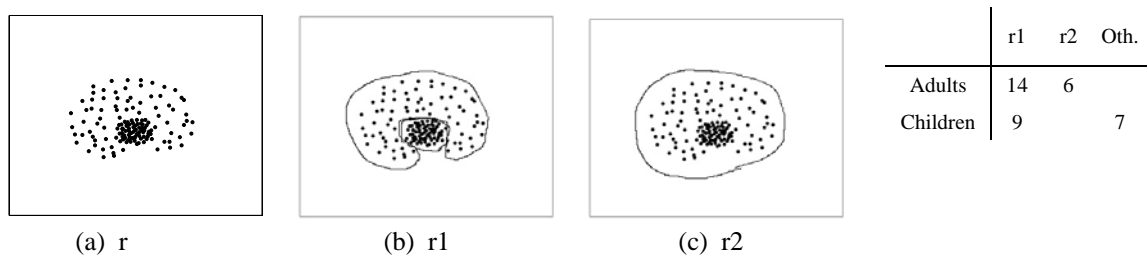


| | r1 | r2 | Oth. |
|---|---|---|---|
| Adults | 14 | 6 | |
| Children | 9 | | 7 |

(a) r  (b) r1  (c) r2

**Figure 9: The solutions proposed for data set "r".**

The $X^2$ test for goodness of fit lead us to reject the uniformity hypothesis for the adults (p=0.05) and accept it for the children (p=0.6) for the "regular"-"non-regular" categories.

### 3.2.2.4  Data set "bb"



(a) bb  (b) bb1  (c) bb2

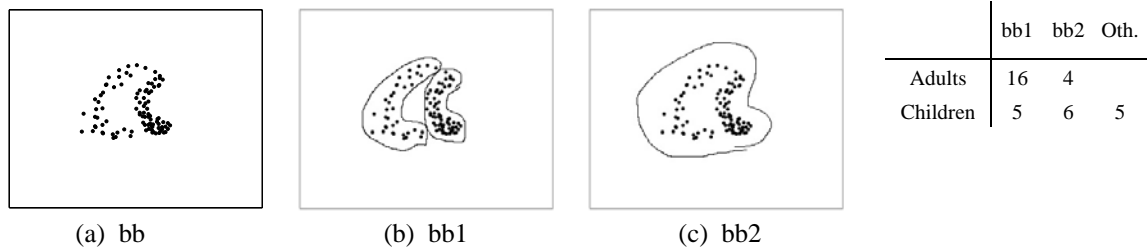|  | bb1 | bb2 | Oth. |
|---|---|---|---|
| Adults | 16 | 4 | |
| Children | 5 | 6 | 5 |

**Figure 10: The solutions proposed for data set "bb".**

The results show that solution "bb1" was overwhelmingly chosen by the adults. For the children the two solutions "bb1" and "bb2" are almost equally suggested, confirming what we noted previously: children are less inclined to sacrifice connectedness to point density differences.

The $X^2$ test for goodness of fit rejects the uniformity hypothesis for the adults and accepts it for the children (p=0.94), confirming the different behaviour of children and adults. The $X^2$ test for independence, for a Solution variable with the three categories as above, lead us to reject the independence hypothesis (p=0.004). The Cramer V is high (V=0.594). This can be attributed to the lack of "Others" in the adult solutions.

### 3.2.2.5  Data set "cc"



(a) cc  (b) cc1  (c) cc2

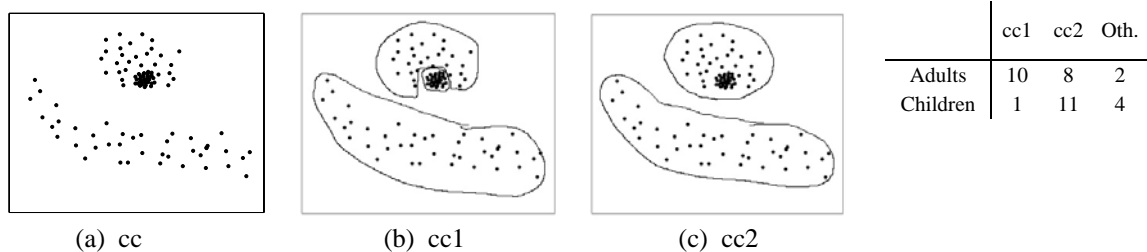|  | cc1 | cc2 | Oth. |
|---|---|---|---|
| Adults | 10 | 8 | 2 |
| Children | 1 | 11 | 4 |

**Figure 11: The solutions proposed for data set "cc".**

We prepared this data set with the aim of comparing it with data set "r". We have included here a similar region to the one appearing in the data set "r". We were expecting similar solutions in the similar regions. We indeed obtained adult results for this data set very similar to the results for data set "r". For the similar regions, the proposed solutions were also similar. In the children results, this did not happen. Children have just considered the existence of the two most evident clusters (solution "cc2"). It seems that many adults were able to decompose the data set on several levels of clusters (something like hierarchical clustering). First by mentally construct 2 clusters and secondly by separating one of them in 2 clusters. Children, on the contrary, tend to value the most prominent feature: connectedness. We think that only a hierarchical mental process is able to justify the differences between adults and children in this data set.

Disregarding solution "Others", the $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults (p=0.64). The $X^2$ test for independence, for a Solution variable

with the three categories as above, lead us to reject the independence hypothesis (p=0.017). The Cramer V is high (V=0.476).

*3.2.3*  Type C: Data sets with crossing clusters

In this subsection, we analyze the group of data sets with crossing clusters. This group is constituted by the set of data sets {l, m, s}. In the following subsections, we present and comment the different proposed solutions for these data sets.

### 3.2.3.1  Data set "l"

In this data set the tests made on adults show that the preferred solution is the one that considers the 2 arms of the cross.
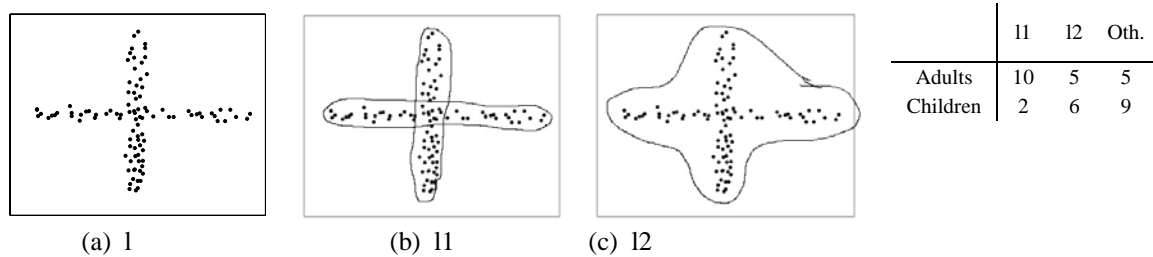


|  | l1 | l2 | Oth. |
|---|---|---|---|
| Adults | 10 | 5 | 5 |
| Children | 2 | 6 | 9 |

(a) l          (b) l1          (c) l2

**Figure 12: The solutions proposed for data set "l".**

Children prefer to consider the cross as a single cluster. Among the other solutions proposed by children, there were a couple of them considering the division of the cross in 4 clusters, one for each branch. These results suggest that adults are able to trade connectedness by structuring direction, a feature not taken into account by children.

The $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults (p=0.33) and rejects it for the children (p=0.018). The $X^2$ test for independence, for a Solution variable with the three categories as above, lead us to reject the independence hypothesis (p=0.04). The Cramer V is high (V=0.415). These results support the different and almost opposite behaviour of adults and children.

### 3.2.3.2  Data set "m"

This data set was the one where there was more reluctance in clustering the data. Adults were divided between the existence of only one cluster and the existence of several clusters.

Among all the solutions, proposed by adults, the most significative was the one that considered the existence of 5 clusters. This solution for data set "m" is very curious when comparing with the solutions proposed for data set "l". In the latter, adults have not considered the hypothesis of dividing the data set in 4 clusters, one for each branch of the cross; however, in data set "m", maybe influenced by the existence of a branch with no correspondence in the other side of the star, adults have decided to consider each branch a single cluster. Children consider this to be a single cluster problem, as they do with data set "l". The same comment made previously on connectedness and structuring direction

applies here.

The $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults (0.64) and rejects it for the children (p≈0). $X^2$ test for independence, for a Solution variable with two categories - "regular clusters" and "non-regular clusters" -, lead us to accept the independence hypothesis (p=0.3). The Cramer V is low (V=0.17).
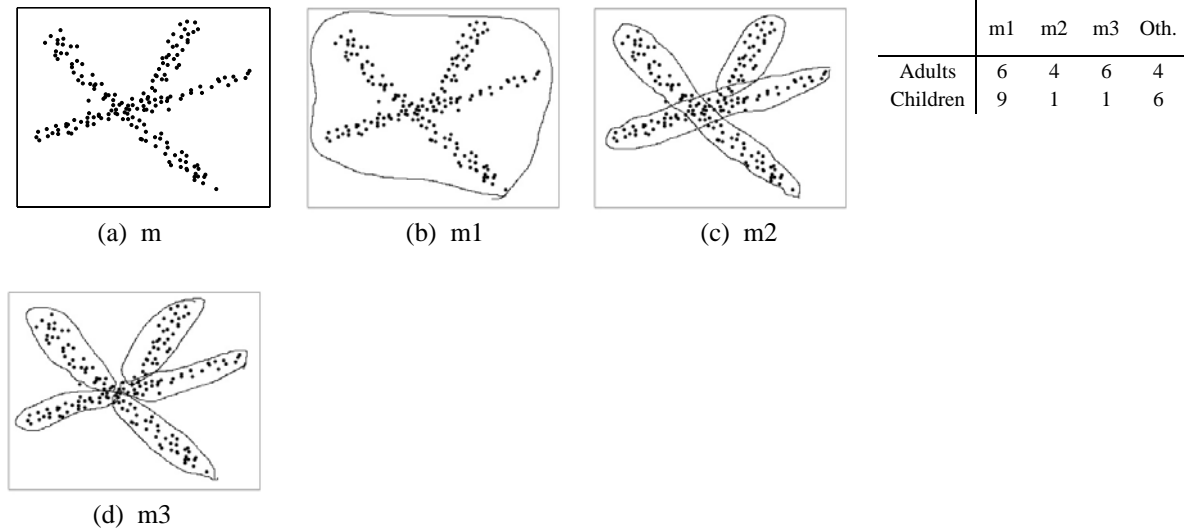


| | m1 | m2 | m3 | Oth. |
|---|---|---|---|---|
| Adults | 6 | 4 | 6 | 4 |
| Children | 9 | 1 | 1 | 6 |

(a) m          (b) m1          (c) m2



(d) m3

**Figure 13: The solutions proposed for data set "m".**

### 3.2.3.3  Data set "s"

In this data set, almost all adults considered the existence of 2 annular clusters as shown in Figure 14c, however children were unable to do the same. We could see on the solutions proposed by the children that, in some cases, they have tried to represent the two clusters without success due to lack of representation skills. This is a data set where the notion of a structuring direction is of primordial importance, explaining the failure of children in "seeing" solution s2.
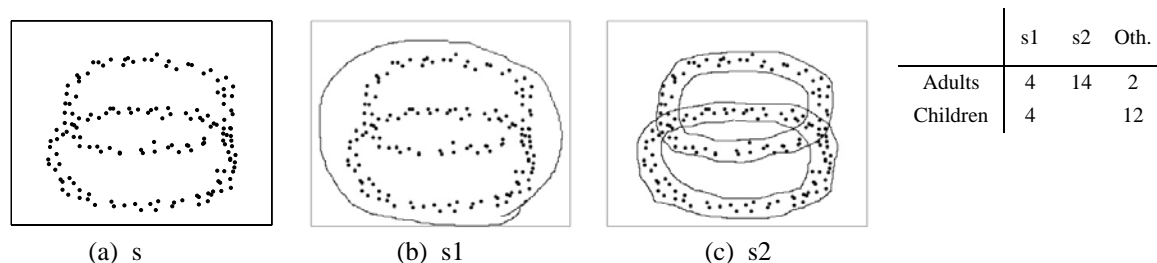


| | s1 | s2 | Oth. |
|---|---|---|---|
| Adults | 4 | 14 | 2 |
| Children | 4 | | 12 |

(a) s          (b) s1          (c) s2

**Figure 14: The solutions proposed for data set "s".**

The $X^2$ test for goodness of fit rejects the uniformity hypothesis for the adults (p<0.014).

*3.2.4*  Type D: Data sets with nested clusters

In this subsection we analyze the group of data sets with nested clusters (clusters inside clusters) not considered in previous types. This group is constituted by the set of data sets {p, z, aa}. For data set "z" there was basically a unique proposed solution, shown in Figure 15.
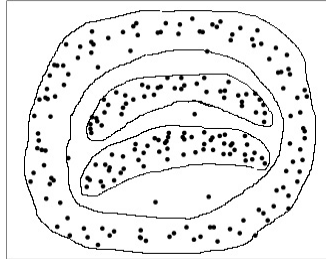


**Figure 15: The solution proposed for data set "z".**

The other Type D data sets are discussed in the following subsections.

### 3.2.4.1  Data set "p"



| | p1 | p2 | Oth. |
|---|---|---|---|
| Adults | 12 | 7 | 1 |
| Children | 10 | 3 | 3 |

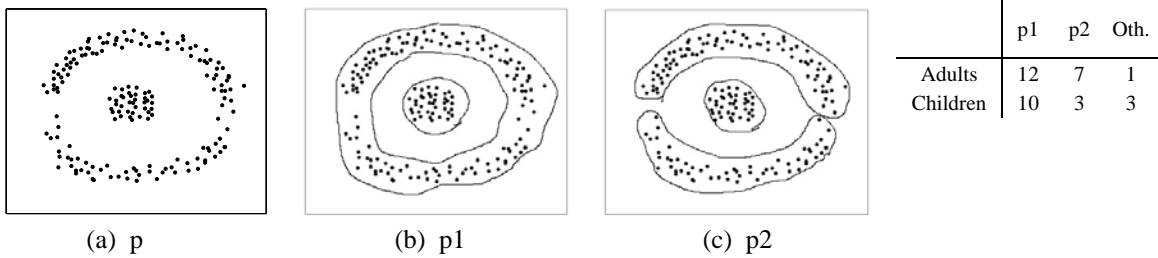(a) p          (b) p1          (c) p2

**Figure 16: The solutions proposed for data set "p".**

Data set "p" has two different proposed solutions. The majority of both children and adults proposed solution "p1".
Disregarding the solution "Others" the $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults (p=0.23) and rejects it for the children (p=0.02. Disregarding the solution "Others" the $X^2$ test for independence lead us to accept the independence hypothesis (p=0.31). The Cramer V is moderate (V=0.256). Thus, although the majority chose "p1", the behaviour of adults and children is different and, in fact, there is a more than chance-explained (at 5% significance level) majority choice of "p1" for the children. This is a strange finding that at first sight could lead us to think that children valued more than adults structuring direction and/or morphology. However, part of the explanation why so many adults chose "p2" may be due to the different point densities of the upper and lower part of the annular cluster; a feature which most of the children didn't see.

### 3.2.4.2  Data set "aa"

As we previously mentioned in section 2, we made this data set similar to data set "p" for comparison purposes. We have separated the annular cluster and we have shifted down

the circular cluster so that it touches the lower section of the annular cluster. By doing that, we tried to percept if the individuals consider the circular cluster as a separate cluster.
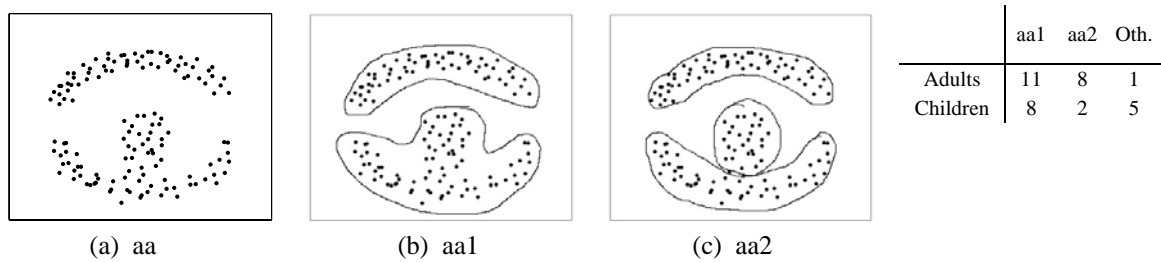


| | aa1 | aa2 | Oth. |
|---|---|---|---|
| Adults | 11 | 8 | 1 |
| Children | 8 | 2 | 5 |

(a) aa  (b) aa1  (c) aa2

**Figure 17: The solutions proposed for data set "aa".**

The results show that this solution ("aa2") was not the preferred solution, especially in the children results, but it almost equals solution "aa1" (only two clusters) in the adult results. Disregarding the solution "Others", the $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults (p=0.49). Uniformity of the three categories is marginally acceptable for the children (p=0.06). The $X^2$ test for independence, for a Solution variable with three categories as above, lead us to reject the independence hypothesis (p=0.038). The Cramer V is high (V=0.42). Therefore, although the "aa2" solution was not the most preferred one by the adults, there is a clear different behaviour of children and adults. Adults were significantly (at 5% significance level) more capable of taking into account the structuring morphological feature present in solution "aa2".

*3.2.5*  Type F: Data sets with spiral-shaped clusters

In this subsection, we analyze the group of data sets with spiral-shaped clusters. This group is constituted by the set of data sets {y, dd}. For both data sets, the individuals basically considered them as 2-cluster data sets (Figure 18), despite the fact that the clusters present a very complex structure compared with the other data sets. We were even surprised by the fact that children also recognized the two spiral clusters presented in data set "dd"; a good illustration of prevalence of a structuring direction over connectedness.
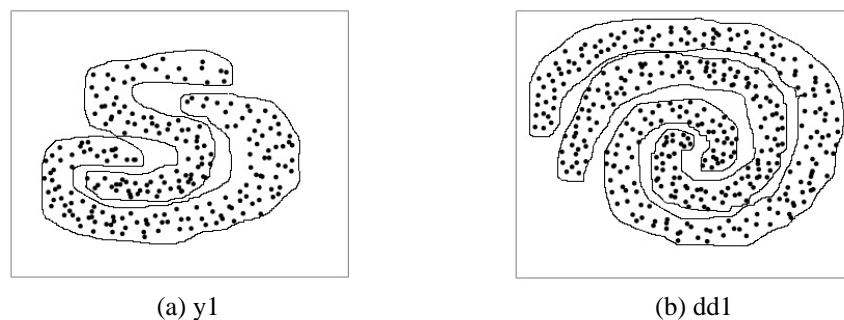


(a) y1  (b) dd1

**Figure 18: The most significative solutions proposed for data sets with spiral shaped clusters.**

*3.2.6*   Type E: Other data sets

In this subsection we analyze the group of data sets not considered in any of the previous groups. This group is constituted by the set of data sets {f, j, u, w, x}. For data set "j" there is basically a unique proposed solution that considers it as a single cluster. The other Type E data sets are discussed in the following subsections.

## 3.2.6.1   Data set "f"

The results suggested by adults for data set "f" were, in our opinion, influenced by the previous solutions given to data set "b". We have already mentioned that these two data sets were intentionally produced with a small difference. In this case, the two clusters of data set "b" were shifted to be almost connected (apparently). We think that this fact, and also the low density in the "connecting" region, was responsible for the predominant 2 clusters solution.
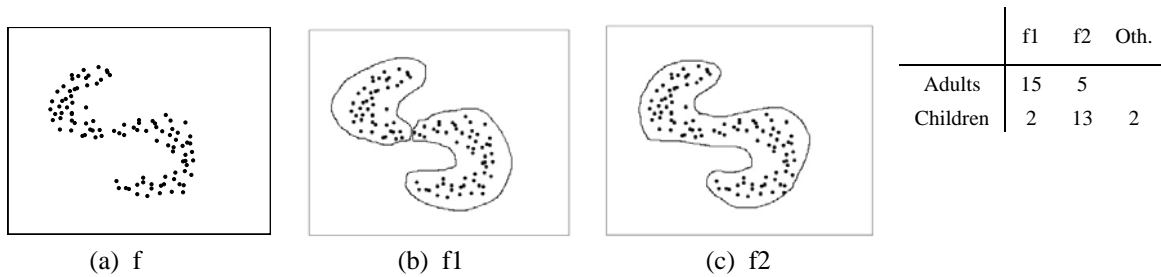


| | f1 | f2 | Oth. |
|---|---|---|---|
| Adults | 15 | 5 | |
| Children | 2 | 13 | 2 |

(a) f          (b) f1          (c) f2

**Figure 19: The solutions proposed for data set "f".**

However, in the children tests, this fact does not happen. It seems that the solutions that they proposed to data set "b" did not affect the proposed solutions for data set "f", confirming the overwhelming value that children attribute to connectedness and/or structuring directions.

Disregarding the solution "Others", the $X^2$ test for goodness of fit rejects (at 5% significance level) the uniformity hypothesis for both adults and children ($p<0.01$). The $X^2$ test for independence, for a Solution variable with the two "regular" categories, lead us to reject the independence hypothesis ($p\approx0$). The Cramer V is quite high ($V=0.61$). Thus statistical analysis confirms that adult and children behaviours are different and opposite of each other.

## 3.2.6.2   Data set "g"

The solutions for this data set are shown in Figure 20.

The majority of the adults have considered this a problem with 3 clusters. The children results are conditioned by the previous mentioned fact that they pay a particular attention to small clusters.

Disregarding solution "Others" the $X^2$ test for goodness of fit lead us to accept the uniformity hypothesis for both children and adults (here, with $p=0.08$). The $X^2$ test for independence, for a Solution variable with two categories (g1, g2), lead us to accept the independence hypothesis ($p=0.21$). The Cramer V is low ($V=0.22$).
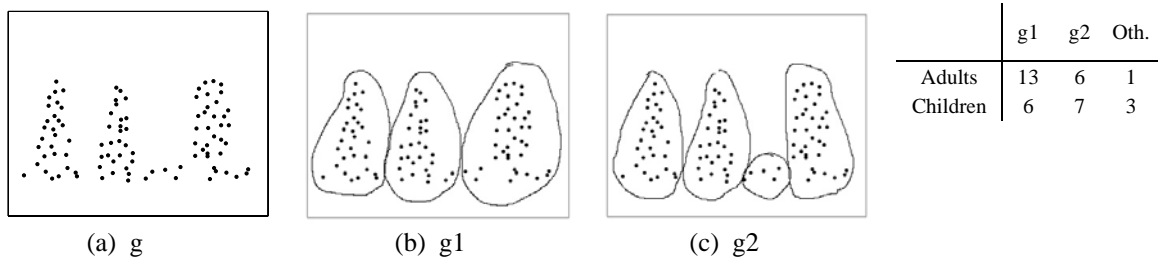
|  | g1 | g2 | Oth. |
|---|---|---|---|
| Adults | 13 | 6 | 1 |
| Children | 6 | 7 | 3 |

(a) g　　　　　(b) g1　　　　　(c) g2

**Figure 20: The solutions proposed for data set "g".**

### 3.2.6.3  Data set "u"

Data set "u" was produced to try to see if differently shaped groups, placed close to each other were considered as one or two clusters. Briefly, the influence of the "structuring morphology" feature.
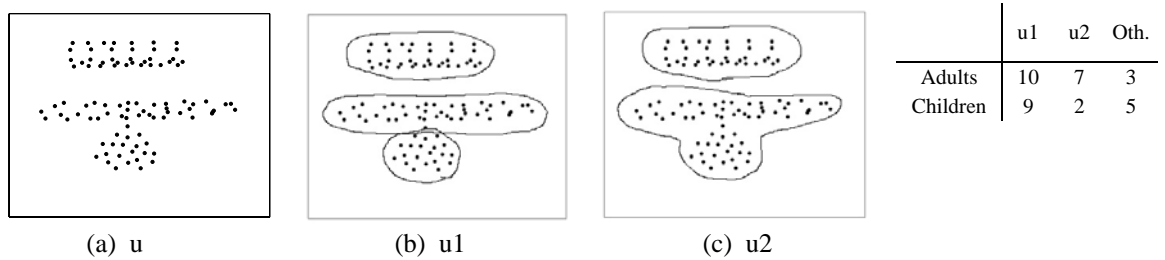


|  | u1 | u2 | Oth. |
|---|---|---|---|
| Adults | 10 | 7 | 3 |
| Children | 9 | 2 | 5 |

(a) u　　　　　(b) u1　　　　　(c) u2

**Figure 21: The solutions proposed for data set "u".**

Regarding the solutions proposed by the adults, we see that surprisingly many adults failed to recognize the existence of three clusters, corresponding to separating the circular cluster from the elongated one. Children, on the contrary, seem to exhibit a definite preference by "u1", valuing the "structuring morphology" feature. They overwhelmingly separate the circular cluster from the elongated one.

The $X^2$ test for goodness of fit marginally accepts the uniformity hypothesis for the adults (p=0.06) and rejects it for the children (p=0.03). The $X^2$ test for independence, for a Solution variable with the three categories as above, lead us to accept the independence hypothesis (p=0.23). The Cramer V is quite moderate (V=0.29).

### 3.2.6.4  Data set "w" and "x"

The solutions proposed for data sets "w" and "x" are very similar. In both cases, there are connections at the ends of the point clouds that influence the different results. Although many two-cluster solutions were proposed, more than 50% of the individuals considered these as one-cluster problems.

Disregarding the solution "Others", the $X^2$ test for goodness of fit accepts the uniformity hypothesis for the adults and for both data sets (p=0.48 and p=0.83 for "w" and "x", respectively). The uniformity hypothesis was only accepted for the children for data set

"x". Also, the $X^2$ test for independence, for a Solution variable with the three regular categories, yielded different results for the datasets: rejection for "w" (p=0.046) and acceptance for "x" (p=0.2).
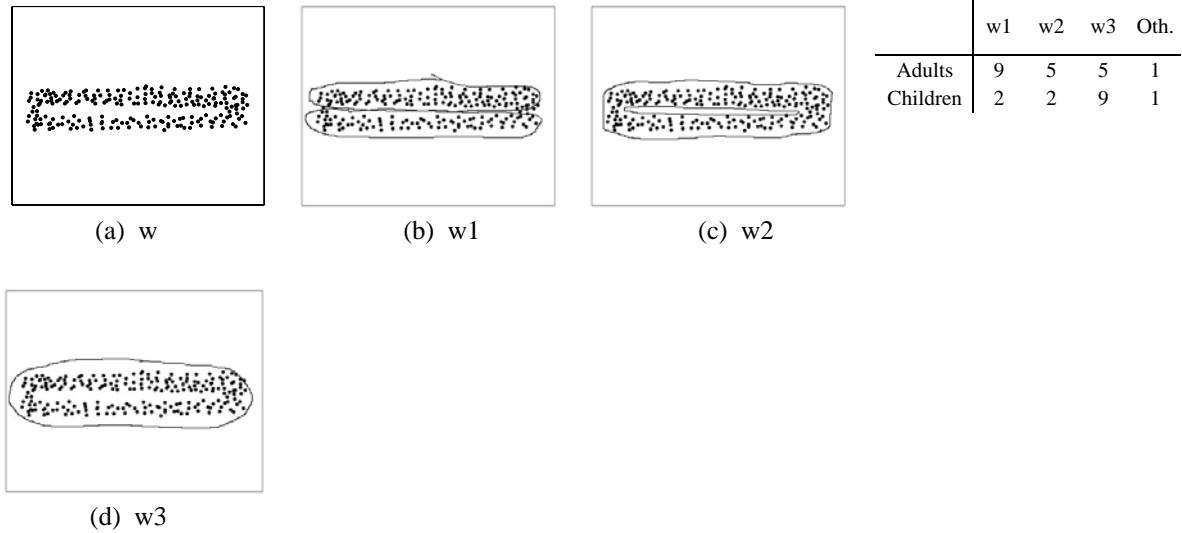


|          | w1 | w2 | w3 | Oth. |
|----------|----|----|----|------|
| Adults   | 9  | 5  | 5  | 1    |
| Children | 2  | 2  | 9  | 1    |

| (a) w | (b) w1 | (c) w2 |

(d) w3

**Figure 22: The solutions proposed for data set "w".**



|          | x1 | x2 | x3 | Oth. |
|----------|----|----|----|------|
| Adults   | 8  | 6  | 6  |      |
| Children | 3  | 3  | 9  | 1    |

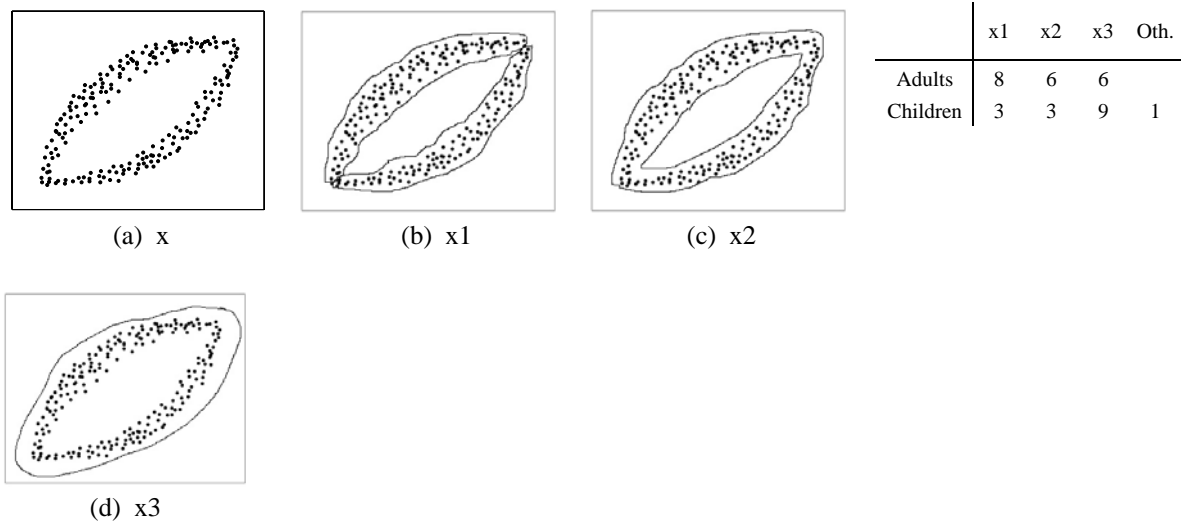| (a) x | (b) x1 | (c) x2 |

(d) x3

**Figure 23: The solutions proposed for data set "x".**

These findings support the different behaviour of adults and children, with the adults valuing more than children the "structuring morphology" feature.

# 4    Conclusions

Clustering solutions proposed by children are quite different from those proposed by adults. We found for several data sets that the $X^2$ test for independence (at 5% significance level) either accepted the independence hypothesis (data sets d, g, m, p, u) or rejected it because of adult and children choices in opposite directions (data sets f, k, l, w, x, aa, bb, cc). Thus, we found statistical evidence supporting the thesis of different cluster behaviour of children and adults in those data sets. Children and adults showed a strong agreement of their clustering preferences for the datasets where clustering is mainly based on the connectedness or structuring direction features (well-separated data sets, nested clusters, spiral-shaped clusters).

Children usually "see" small clusters focusing their attention in small regions, leading to solutions with a large number of clusters. They praise overwhelmingly the connectedness feature to the point of sacrificing other ones.

From the analysis of the different types of data sets we draw the following conclusions:

- Connectedness or structuring-direction features are the easiest features to handle, by both adults and children.
- Children are often unable to sacrifice connectedness by other features. This is especially the case with data sets exhibiting cross-type clusters.
- Point density and morphological structuring are the most difficult clustering features to handle.
- Adults seem capable of performing some sort of hierarchical clustering, using clustering features at different decision levels. This was mainly apparent in the solutions produced for data sets "p", "k", "aa" and "cc".
- A small difference in the data sets, like in the pairs "b"-"f" and "p"-"aa", can lead to very different clustering solutions. This is especially to be expected when the point density and morphological structuring features come into play.