

Maria Luísa de Moraes e Sousa de Castro

**Classificação de Imagens de Satélite por
Aglomerção Hierárquica usando o MatLab**

U. PORTO

**FC FACULDADE DE CIÊNCIAS
UNIVERSIDADE DO PORTO**

Faculdade de Ciências da Universidade do Porto
Junho de 2006

Reg. 509268
Cota TESE N.º 168

Faculdade de Ciências do Porto
MATEMÁTICA

O presidente do júri,
António Vaulianek



FC

Biblioteca
Faculdade de Ciências
Universidade do Porto



0000100842

Maria Luísa de Moraes e Sousa de Castro

**Classificação de Imagens de Satélite por
Aglomerção Hierárquica usando o MatLab**

U. PORTO

**FC FACULDADE DE CIÊNCIAS
UNIVERSIDADE DO PORTO**

*Tese submetida à Faculdade de Ciências da Universidade do Porto
para obtenção do grau de Mestre em Engenharia Matemática*

Dissertação realizada sob a supervisão do
Professor Doutor André Ribeiro da Silva de Almeida Marçal
Departamento de Matemática Aplicada
Faculdade de Ciências da Universidade do Porto
Junho de 2006

Agradecimentos

Ao Professor André Marçal dedico grande parte dos meus agradecimentos, por toda a segurança e motivação com que me contagiou, ainda que à distância, ao longo de toda a orientação do meu trabalho.

À Janete e à Alexandra agradeço a disponibilidade, de tempo e espaço, respectivamente, que, nos momentos certos, sempre demonstraram.

Aos meus pais apresento sincera gratidão pela confiança e paciência demonstradas, em especial em alturas de tensão redobrada.

E finalmente, agradeço ao Martim a compreensão pelo tempo que dediquei a esta dissertação e a força que me transmitiu durante todo o percurso.

Resumo

Neste trabalho é apresentado um método de processamento de imagens multi-espectrais para posterior classificação, utilizando as ferramentas disponíveis no MATLAB. Sendo a classificação um processo computacionalmente pesado e tendo as imagens multi-espectrais um elevado número de padrões, é necessária a aplicação de um tratamento prévio de forma a reduzir o volume de dados. O método descrito compreende ferramentas de processamento de imagem, nomeadamente filtros lineares espaciais, filtros de *sobel* e segmentação *watershed*. O processo de classificação recorre a métodos não supervisionados, tendo sido aplicados 5 métodos distintos de agregação hierárquica. Foram criados conjuntos de dados de teste para avaliar a aplicabilidade dos classificadores e foram determinadas as partições melhor representativas da estrutura intrínseca de cada conjunto de dados. Por fim o processo envolvendo as três fases - filtragem, segmentação e classificação - é aplicado a uma imagem multi-espectral de satélite da zona da Ria de Aveiro (Portugal).

Abstract

This dissertation presents a method to implement classification using hierarchical aggregation tools, available in MATLAB, to multi-spectral images. These classification methods are computationally heavy, and as multi-spectral images contain a large number of patterns, it becomes necessary to apply a previous treatment in order to decrease the amount of data. The developed method includes image processing tools, such as linear spatial filters and watershed segmentation. The classification process employs unsupervised techniques, with 5 distinct hierarchical aggregation methods applied. Synthetic datasets were created to evaluate the applicability of classifiers, while the partitions which best suits the intrinsic structure of the data were found using two internal similarity indices. Finally, the process formed by the three steps – filtering, segmentation and classification – is applied to a multi-spectral satellite image of the region of Ria de Aveiro (Portugal).

REQUISIÇÃO DE MEMBROS DO PORS
MATEMÁTICA

Índice

Agradecimentos	i
Resumo	ii
Abstract	iii
Lista de Figuras	iv
1. Introdução	1
2. Metodologia	5
2.1. Classificação	5
2.1.1. Funcionamento	5
2.1.2. Métodos de Agregação Hierárquica	9
2.2. Índices de Semelhança	11
2.2.1. Índice de Davies e Bouldin	11
2.2.2. Índice de Xu, Kamath e Capson	12
3. Teste com Dados Sintéticos	13
3.1. Construção dos Conjuntos de Dados de Teste	13
3.2. Classificação dos Conjuntos de Dados de Teste	15
3.2.1. Agregação pelo Método <i>Single</i>	16
3.2.2. Agregação pelo Método <i>Complete</i>	22
3.2.3. Agregação pelo Método <i>Average</i>	28
3.2.4. Agregação pelo Método <i>Weighted</i>	34
3.2.5. Agregação pelo Método <i>Ward</i>	39
3.2.6. Agregação pelos Métodos <i>Median</i> e <i>Centroid</i>	45
3.3. Conclusões	46
4. Aplicação a Imagens Multi-espectrais	47
4.1. Método Proposto	47
4.2. Filtragem	48
4.3. Segmentação	52
4.4. Classificação e Transposição dos Resultados	58
5. Teste com Imagem de Satélite	61
5.1. Método de Agregação <i>Single</i>	63
5.2. Método de Agregação <i>Complete</i>	64
5.3. Método de Agregação <i>Average</i>	66
5.4. Método de Agregação <i>Weighted</i>	68
5.5. Método de Agregação <i>Ward</i>	70
6. Conclusão	75
Referências	77

Lista de Figuras

Figura 2.1 Representação gráfica do conjunto de dados A.	6
Figura 2.2 Árvore de agregação hierárquica (dendograma).	7
Figura 2.3 Árvore de agregação hierárquica com 'corte' em $c=4$.	8
Figura 3.1 Representações gráficas dos conjuntos de dados de teste criados: I, II, III, IV e V.	14
Figura 3.2 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.	16
Figura 3.3 Representação dos valores dos índices DB e Xu, para o conjunto I.	17
Figura 3.4 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.	17
Figura 3.5 Representação dos valores dos índices DB e Xu, para o conjunto II.	18
Figura 3.6 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.	18
Figura 3.7 Representação dos valores dos índices DB e Xu, para o conjunto III.	19
Figura 3.8 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.	19
Figura 3.9 Representação dos valores dos índices DB e Xu, para o conjunto IV.	20
Figura 3.10 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.	20
Figura 3.11 Representação dos valores dos índices DB e Xu, para o conjunto V.	21
Figura 3.12 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.	22
Figura 3.13 Representação dos valores dos índices DB e Xu, para o conjunto I.	23
Figura 3.14 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.	23
Figura 3.15 Representação dos valores dos índices DB e Xu, para o conjunto II.	24
Figura 3.16 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.	24
Figura 3.17 Representação dos valores dos índices DB e Xu, para o conjunto III.	25
Figura 3.18 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.	25
Figura 3.19 Representação dos valores dos índices DB e Xu, para o conjunto IV.	26
Figura 3.20 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.	26
Figura 3.21 Representação dos valores dos índices DB e Xu, para o conjunto V.	27
Figura 3.22 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.	28
Figura 3.23 Representação dos valores dos índices DB e Xu, para o conjunto I.	28
Figura 3.24 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.	29
Figura 3.25 Representação dos valores dos índices DB e Xu, para o conjunto II.	30
Figura 3.26 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.	30
Figura 3.27 Representação dos valores dos índices DB e Xu, para o conjunto III.	31
Figura 3.28 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.	31
Figura 3.29 Representação dos valores dos índices DB e Xu, para o conjunto IV.	32
Figura 3.30 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.	32

Figura 3.31. Representação dos valores dos índices DB e Xu, para o conjunto V.	33
Figura 3.32 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.	34
Figura 3.33 Representação dos valores dos índices DB e Xu, para o conjunto I.	34
Figura 3.34 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.	35
Figura 3.35 Representação dos valores dos índices DB e Xu, para o conjunto II.	35
Figura 3.36 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.	36
Figura 3.37 Representação dos valores dos índices DB e Xu, para o conjunto III.	36
Figura 3.38 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.	37
Figura 3.39 Representação dos valores dos índices DB e Xu, para o conjunto IV.	37
Figura 3.40 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.	38
Figura 3.41 Representação dos valores dos índices DB e Xu, para o conjunto V.	38
Figura 3.42 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.	39
Figura 3.43 Representação dos valores dos índices DB e Xu, para o conjunto I.	40
Figura 3.44 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.	40
Figura 3.45 Representação dos valores dos índices DB e Xu, para o conjunto II.	41
Figura 3.46 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.	41
Figura 3.47 Representação dos valores dos índices DB e Xu, para o conjunto III.	42
Figura 3.48 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.	42
Figura 3.49 Representação dos valores dos índices DB e Xu, para o conjunto IV.	43
Figura 3.50 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.	43
Figura 3.51 Representação dos valores dos índices DB e Xu, para o conjunto V.	44
Figura 3.52 Representação de <i>clusters</i> e respectivos centróides.	45
Figura 4.1 Esquema do método proposto para a classificação de imagens multi-espectrais: Imagem Multi-espectral (A), Imagem Segmentada (B), Conjunto de Padrões Reduzido (C), Classificação (D) e Imagem Classificada (E).	48
Figura 4.2 O mecanismo de filtragem espacial.	49
Figura 4.3 Região 3×3 de uma imagem.	51
Figura 4.4 Filtros de <i>Sobel</i> : realce de arestas horizontais (A) e de arestas verticais (B).	52
Figura 4.5. Bacias de captação, <i>watersheds</i> e mínimos regionais.	53
Figura 4.6 A: Duas bacias de captação inundadas na fase <i>n-1</i> da subida do nível da água; B: Fase <i>n</i> da inundação, junção da água das duas bacias; C: Elemento estruturante da dilatação; D: Resultado da dilatação e construção do dique.	56
Figura 4.7 Imagem original (A), imagem de intensidade (B), imagem filtrada (C), imagem gradiente (D), imagem segmentada (E) e segmentação sobreposta com imagem original (F).	57
Figura 4.8 Secção 3×3 da imagem resultante da segmentação.	58
Figura 5.1 Imagem de teste (composição RGB das bandas 1, 2 e 3, com histograma modificado).	61

Figura 5.2 Representação a 3 dimensões dos pontos que representam cada objecto obtido da segmentação.	62
Figura 5.3 Imagem classificada em 40 classes pelo método de agregação <i>single</i> .	63
Figura 5.4 Imagem classificada em 40 classes pelo método de agregação <i>complete</i> .	64
Figura 5.5 Representação dos valores do índice DB para a classificação pelo método <i>complete</i> .	65
Figura 5.6 Representação dos valores do índice Xu para a classificação pelo método <i>complete</i> .	66
Figura 5.7 Imagem classificada em 29, 11, 10 e 2 classes pelo método de agregação <i>complete</i> .	66
Figura 5.8 Imagem classificada em 40 classes pelo método de agregação <i>average</i> .	66
Figura 5.9 Representação dos valores do índice DB para a classificação pelo método <i>average</i> .	67
Figura 5.10 Representação dos valores do índice Xu para a classificação pelo método <i>average</i> .	67
Figura 5.11 Imagem classificada em 26, 5, 4 e 2 classes pelo método de agregação <i>average</i> .	68
Figura 5.12 Imagem classificada em 40 classes pelo método de agregação <i>weighted</i> .	68
Figura 5.13 Representação dos valores do índice Xu para a classificação pelo método <i>weighted</i> .	69
Figura 5.14 Representação dos valores do índice Xu para a classificação pelo método <i>weighted</i> .	69
Figura 5.15 Imagem classificada em 31, 5 e 4 classes pelo método de agregação <i>weighted</i> .	70
Figura 5.16 Imagem classificada em 40 classes pelo método de agregação <i>ward</i> .	70
Figura 5.17 Representação dos valores do índice DB para a classificação pelo método <i>ward</i> .	71
Figura 5.18 Representação dos valores do índice Xu para a classificação pelo método <i>ward</i> .	71
Figura 5.19 Imagem classificada em 16, 7, 4 e 3 classes pelo método de agregação <i>ward</i> .	72
Figura 5.20 Imagem classificada em 40, 37, 34, 31, 28, 25, 22, 19, 16, 13, 10, 8, 6, 4, e 2 classes pelo método de agregação <i>ward</i> .	73

Capítulo 1

Introdução

As imagens da superfície terrestre obtidas a partir de satélites contêm uma grande quantidade de informação, na maioria das vezes impossível de captar recorrendo apenas à visão humana. As imagens de satélite são em geral multi-espectrais, isto é, apresentam mais do que uma banda, normalmente referentes à parte visível e de infra-vermelhos do espectro electromagnético.

O processo de classificação permite a extracção da informação relevante de uma imagem multi-espectral, recorrendo normalmente a meios computacionais, com o objectivo de identificar regiões homogéneas de ocupação do solo, denominadas classes temáticas. Uma amostra de cada classe temática pode ser utilizada para o processo de classificação, neste caso a classificação é supervisionada pois existe informação à priori sobre as propriedades espectrais das zonas que se pretendem identificar. Por vezes tal informação não é disponibilizada e realiza-se uma classificação não supervisionada, contando apenas com a imagem a analisar. Uma classificação não supervisionada pode ser mais vantajosa nos casos em que o levantamento de amostras de classes temáticas seja difícil de realizar, quer por existirem em número elevado quer pela possível inacessibilidade da região a que se reporta a imagem multi-espectral.

Uma imagem pode ser definida como uma função bidimensional $f(x,y)$, onde x e y são coordenadas espaciais e o valor de f para qualquer par de coordenadas (x,y) , representa a intensidade da imagem nesse ponto. Quando x , y e os valores de intensidade são todos finitos e quantidades discretas, a imagem toma a designação de imagem digital. A classificação de uma imagem digital é realizada sobre cada ponto que a compõe, denominado pixel, a que corresponde um vector de intensidade, n -dimensional, sendo n o número de bandas da imagem multi-espectral.

As imagens multi-espectrais de satélite são em geral conjuntos de dados demasiado grandes (número de pixels e/ou bandas) para a utilização de classificadores com exigências computacionais elevadas, como é o caso dos métodos de agregação hierárquica. Neste trabalho apresenta-se um método que permite a utilização de classificadores de aglomeração hierárquica em imagens multi-espectrais de satélite, através de um processo de redução do número de padrões a classificar baseada na segmentação de imagem.

A ideia do classificador aplicado à imagem ASTER é baseada numa abordagem aglomerativa, ou seja, começa por formar cada classe por um padrão distinto (partição inicial), e vai progressivamente unindo as classes que considera mais próximas, obtendo em cada nível uma partição do conjunto de dados inicial até obter a partição formada apenas por uma classe, produzindo uma estrutura ou árvore hierárquica. As medidas utilizadas para determinar as classes mais próximas são várias e a aplicação

de cada uma em separado dá origem a vários métodos de classificação. Em geral, se as classes presentes nos dados forem compactas e bem definidas, todos os métodos produzem os mesmos resultados, caso contrário podem ser obtidos resultados bastante distintos para cada classificador (Duda, 2001). Neste sentido, e com o objectivo de testar o desempenho de cada método de classificação, foram criados conjuntos de dados sintéticos, representativos de diversas situações. No Capítulo 2 são apresentados os métodos de classificação usados neste trabalho. Os resultados da aplicação de cada método de agregação hierárquica a cada conjunto de teste, são apresentados no Capítulo 3. A aplicação de um classificador a um conjunto de dados, produz vários resultados dependendo do número de classes que se pretende encontrar. Cada resultado representa uma partição dos dados e, normalmente, apenas uma partição representa a estrutura intrínseca do conjunto de dados. Para determinar a partição que melhor representa cada conjunto de dados gerados, foram aplicados dois índices de semelhança internos: os índices de Davies e Bouldin e de Xu.

O software utilizado para a classificação foi o MATLAB, no entanto, foi necessário executar um tratamento à imagem a testar, anterior à classificação, de forma a reduzir o volume de dados a analisar. Para tal foi aplicada à imagem digital uma técnica com princípios análogos à da classificação, denominada segmentação, que divide a imagem em diferentes regiões cujos pixels partilham determinadas características (Fu, 1981). A segmentação difere da classificação na medida em que forma conjuntos disjuntos, atribuindo a cada um uma etiqueta diferente, enquanto que a classificação divide a imagem em regiões separadas espacialmente mas que podem pertencer à mesma classe.

Existem essencialmente três tipos de abordagem para realizar uma segmentação: partição de histograma ou *thresholding*, detecção de descontinuidades e extracção de regiões, as quais, por si só, apresentam, vantagens e inconvenientes. A aplicação da segmentação a um conjunto de dados é uma fase de extrema importância no pré-processamento de uma classificação pois os erros possivelmente cometidos nesta etapa irão acentuar-se no processo de classificação (Fu, 1981). A técnica de segmentação utilizada neste trabalho, *watershed*, combina alguns dos conceitos vantajosos das três abordagens, produzindo desse modo resultados mais sólidos e estáveis. O método de classificação *watershed* consiste na construção de linhas divisórias, denominadas *watersheds*, que delimitam as regiões a separar ou objectos. A fase de pré-processamento de imagem, incluindo o método de segmentação *watershed*, será exposta em maior detalhe no Capítulo 4.

A imagem escolhida para este trabalho, foi uma imagem com 3 bandas espectrais, RGB (componentes vermelha, verde e azul) referente à zona da Ria de Aveiro, obtida pelo instrumento ASTER (*Advanced Spaceborne Thermal Emission and Reflection Radiometer*) do satélite Terra.

A imagem que se pretende classificar neste trabalho é uma imagem digital de dimensão 800×800, ou seja, pretendem-se classificar um total de 640.000 pixels. No entanto, a aplicação directa de um processo de segmentação a um número elevado de dados, com a possível existência de ruído, pode provocar uma sobresegmentação, ou seja, a obtenção de um número de regiões, de pequena dimensão, demasiado elevado para ser vantajoso. Para prevenir eventuais situações de sobresegmentação, a

imagem foi suavizada através da aplicação de filtros passa-baixo, ou seja, filtros que atenuam as altas frequências das intensidades dos pixels, e filtros gradiente, que acentuam as descontinuidades.

Após a filtragem e segmentação, obteve-se um conjunto de dados de dimensão mais reduzida que o original, tendo sido efectuada uma classificação não supervisionada sobre este novo conjunto de dados. Neste caso, o número de classes originais é desconhecido e para determinar o número natural de classes nos dados, ou seja, a melhor partição, são aplicados os índices de semelhança de Davies e Bouldin e de Xu. O mecanismo utilizado neste processo pode ser consultado em mais detalhe no Capítulo 4 e os resultados da classificação apresentam-se no Capítulo 5.

Neste trabalho é proposto e desenvolvido (no MATLAB) um método que tem como objectivo a classificação de imagens multi-espectrais e que se divide em três fases: uma fase inicial de filtragem, um processo de segmentação pelo método *watershed* e, finalmente, a classificação por métodos distintos pela forma de agregação hierárquica.

Capítulo 2

Metodologia

Neste capítulo apresenta-se uma análise do classificador `clusterdata`, disponível no MATLAB e pertencente à *toolbox* de estatística (The Mathworks, 2004). A análise foi efectuada pela aplicação deste classificador, sob diferentes configurações, em conjuntos de dados sintéticos.

2.1. Classificação

O `clusterdata` é um classificador de aglomeração que parte do conjunto de elementos dado e os agrupa gradualmente com base numa métrica definida, construindo assim uma estrutura (árvore) hierárquica. Esta ferramenta recebe como *input* o conjunto de dados que se pretende classificar sob a forma de uma matriz, $X (n \times d)$, onde cada linha, n , representa um padrão, x . O conjunto de dados é formado então por n elementos em que cada x_i ($0 < i \leq n$) tem dimensão d e é da forma $x_i = (x_{i1}, x_{i2}, \dots, x_{id})$. Na aplicação da função `clusterdata` é possível definir três parâmetros que vão caracterizar o processo de classificação:

- `pdist`, define o tipo de distância a utilizar entre pares de observações;
- `linkage`, define o método de agregação para a construção da árvore hierárquica;
- `cluster`, define a forma de construção de classes a partir da árvore hierárquica.

2.1.1 Funcionamento

Para o cálculo da distância entre pares de observações, parâmetro `pdist`, foi utilizada a distância Euclidiana. A ferramenta `pdist` recebe como *input* a matriz X e produz um vector Y constituído pelas medidas das distâncias Euclidianas calculadas para todos os pares que é possível formar com as n observações, ou seja, Y será um vector com $n(n-1)/2$ valores.

Considere-se, como exemplo, o conjunto de dados, A , formado por seis elementos a duas dimensões ($n=6; d=2$): $A = \{(2,3); (4,5); (3,7); (8,9); (6,9); (4,6)\}$ (figura 2.1).

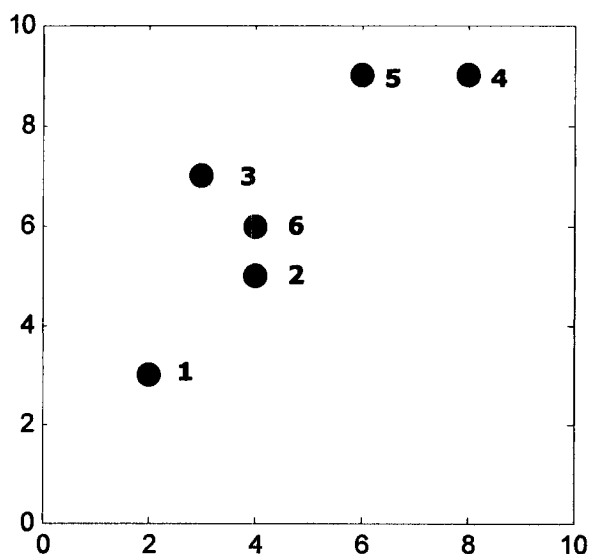


Figura 2.1 Representação gráfica do conjunto de dados A.

O vector obtido, B, aplicando a função `pdist` ao conjunto A, é constituído pelas 15 medidas das distâncias (valores apresentados com 5 algarismos significativos) entre os pares de elementos de A:

$B = (2.8284, 4.1231, 8.4853, 7.2111, 3.6056, 2.2361, 5.6569, 4.4721, 1.000, 5.3852, 3.6056, 1.4142, 2.0000, 5.0000, 3.6056)$.

A primeira medida armazenada em B, b_1 , é a distância entre os elementos $a_1 = (2,3)$ e $a_2 = (4,5)$ (d_{12}); a segunda medida, b_2 , é igual a d_{13} ; $b_3 = d_{14}$ e assim sucessivamente até $b_{15} = d_{56}$.

Na fase que se segue, o vector produzido pela função `pdist`, Y, será processado pela ferramenta `linkage` utilizando uma das sete possibilidades de escolha para o algoritmo de agregação: `single`, `complete`, `average`, `weighted`, `ward`, `centroid` e `median` (cada um destes métodos será descrito com maior detalhe na secção 2.1.2). De acordo com o método seleccionado, a função `linkage` irá construir a estrutura hierárquica que traduz a agregação das observações. O *output* desta função será uma matriz, Z, de dimensão $(n-1) \times 3$ onde as duas primeiras colunas contêm os índices dos pares de elementos/classes que são agregados para formar um novo grupo e a terceira coluna contém a distância entre o par de elementos/classes correspondente. Em relação ao exemplo exposto, a matriz C (5×3) obtida pela aplicação de `linkage`, ao vector B, com o método de agregação `single` é:

$$C = \begin{bmatrix} 2 & 6 & 1.000 \\ 3 & 7 & 1.4142 \\ 4 & 5 & 2.0000 \\ 1 & 8 & 2.8284 \\ 9 & 10 & 3.6056 \end{bmatrix}$$

Pelo método *single*, a agregação é efectuada unindo as duas classes (*clusters*) separadas pela menor distância. Dessa forma, serão agregados em primeiro lugar os *clusters* (constituídos nesta fase por um único elemento) nas posições 2 e 6, já que $d(a_2, a_6) = 1$ é a menor das distâncias armazenadas no vector B. O *cluster* obtido desta união terá a posição $n+1$, ou seja, a posição 7. O processo de agregação continua até à obtenção de apenas uma classe. A figura que se segue representa a árvore hierárquica resultante ou dendograma.

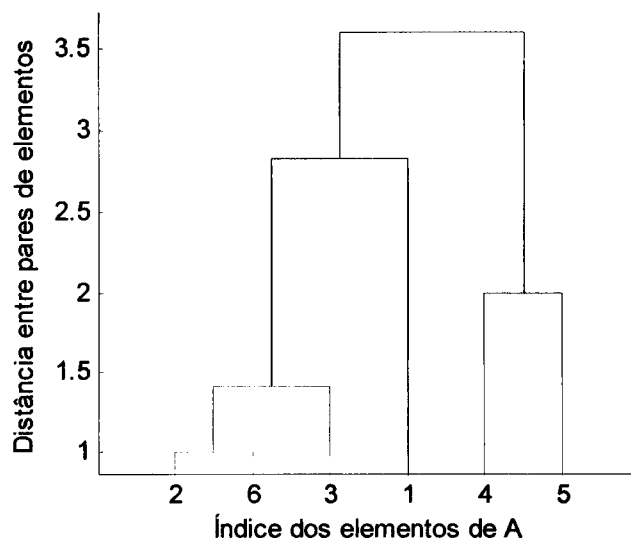


Figura 2.2 Árvore de agregação hierárquica (dendograma).

A árvore de agregação hierárquica obtida mostra que os elementos a_2 e a_6 foram os primeiros a serem agregados, formando o *cluster* 7. Por sua vez, este é unido a a_3 formando o *cluster* 8, e assim sucessivamente até ser formada a última classe, ou seja, aquela que contém todos os elementos.

Finalmente é aplicada a função `cluster` à matriz Z para partir a árvore hierárquica no nível em que esta contenha o número pretendido de classes. Para a utilização desta ferramenta foi seleccionado o parâmetro de quebra `maxclust` o qual, associado a um valor c , produz o resultado da agregação no momento em que existem c classes. O *output* da função `cluster` (que é também o resultado final da aplicação de `clusterdata`) é um vector de dimensão n que em cada posição armazena o número da classe atribuído a cada padrão do conjunto de dados original. Aplicando esta função com $c=4$ à matriz C do exemplo dado, correspondente ao esquema da figura 2.3, obtém-se $D_4 = (3,4,4,1,2,4)$. O resultado obtido corresponde ao 'corte' do dendograma no nível em que existem exactamente 4 classes:

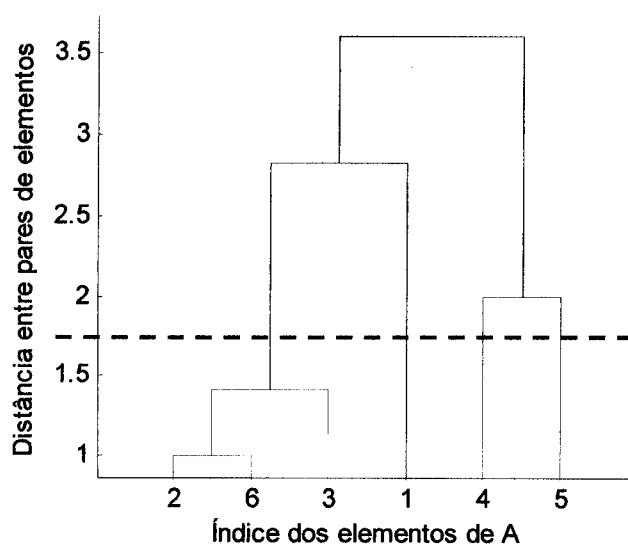


Figura 2.3 Árvore de agregação hierárquica com ‘corte’ em $c=4$.

O *output* da função `clusterdata` atribui a cada elemento dos dados originais uma etiqueta que representa a classe a que esse elemento pertence, ou seja, representa a partição dos dados em c classes.

Num só comando, o processo de classificação aplicado à matriz A do exemplo, traduz-se em:

```
D4 = clusterdata (A, 'maxclust', 4, 'linkage', 'single').
```

2.1.2. Métodos de Agregação Hierárquica

No processo de agregação hierárquica de dois *clusters*, o método `single`, também designado por método do ‘vizinho mais próximo’, determina a menor distância entre todos os pares de elementos dos dois *clusters*. De seguida procede à junção das classes às quais pertencem os elementos do par que apresenta a menor distância, formando um novo *cluster*. Sejam r e s dois *clusters* presentes no conjunto de dados X , sendo n_r e n_s o número de elementos pertencentes a cada um dos *clusters*, r e s , respectivamente. A distância entre estes dois conjuntos, pelo método `single` é dada por:

$$d(r, s) = \min(\text{dist}(x_{ri}, x_{sj})), i \in (1, \dots, n_r), j \in (1, \dots, n_s). \quad (1)$$

Tratando-se de um algoritmo com forte dependência do número de padrões existentes nos dados, a complexidade computacional é $O(cn^2d)$, onde c será o número pretendido de classes a formar (Duda, 2001).

Outro método de agregação é designado por *complete* ou método do ‘vizinho mais afastado’. Neste caso a distância entre dois *clusters*, r e s , é a maior das distâncias entre todos os pares de elementos desses dois *clusters*:

$$d(r, s) = \max(\text{dist}(x_{ri}, x_{sj})), i \in (1, \dots, n_r), j \in (1, \dots, n_s). \quad (2)$$

A adopção da distância máxima e mínima de cada um dos métodos descritos, representa dois extremos na medição da distância entre duas classes. Desta forma estes algoritmos tendem a ser demasiado sensíveis a *outliers*¹ (Duda, 2001). Os métodos que se seguem tendem a diminuir esta susceptibilidade considerando a média das distâncias, o número de elementos e/ou os centros dos *clusters*. A complexidade temporal deste método é da mesma ordem de grandeza do anterior, $O(cn^2d)$.

O método de agregação hierárquica *average* considera a média das distâncias entre todos os pares de elementos nos *clusters* r e s :

$$d_a(r, s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} \text{dist}(x_{ri}, x_{sj}), i \in (1, \dots, n_r), j \in (1, \dots, n_s). \quad (3a)$$

O método de agregação *weighted* em lugar de calcular apenas a média, utiliza a média pesada. A função usada para calcular a distância entre o novo *cluster*, r , recém-formado pela junção dos *clusters* p e q , e um outro *cluster*, s , é dada por (Sung, 2003)²:

$$d(r, s) = \frac{n_p \times d_a(p, s) + n_q \times d_a(q, s)}{n_p + n_q}. \quad (3b)$$

Para a aplicação do método de agregação *centroid*, *median* e *ward* é necessário determinar o centro de cada classe. Para um *cluster* r , o seu centro, m_r , é definido por:

$$m_r = \frac{1}{n_r} \sum_{i=1}^{n_r} x_{ri}. \quad (4)$$

Pelo método de agregação *centroid*, a medida da distância entre dois *clusters* é dada simplesmente pela distância Euclidiana entre os centros dos respectivos *clusters*:

$$d(r, s) = \|m_r - m_s\|.$$

¹ Elementos do conjunto de dados que se desviam significativamente da média dos restantes elementos.

² O manual de apoio do MatLab não refere a fórmula utilizada neste método.

Este método não tem em conta o número de elementos de cada *cluster* nem a sua dispersão em torno do centro.

O método de agregação *median*, utiliza a distância Euclidiana entre os centros pesados das classes em questão:

$$d(r, s) = \|\tilde{x}_r - \tilde{x}_s\|, \text{ onde } \tilde{x}_r \text{ é o centro pesado da classe } r.$$

O método de Ward, consiste num processo iterativo que em cada passo tenta construir a partição que minimiza a soma dos quadrados das distâncias entre os elementos de cada *cluster* e o seu centro, ou seja, a soma dos quadrados *within-cluster*. Este método utiliza a seguinte medida de distância:

$$d(r, s) = \sqrt{\frac{n_r n_s}{n_r + n_s}} \|m_r - m_s\|. \quad (5)$$

Deste modo, ao seleccionar os *clusters* a agregar, o método de Ward tem em conta o número de objectos em cada *cluster* assim como a distância entre *clusters*. Em geral, o uso desta técnica tende a favorecer a agregação de *clusters* formados apenas por um ou por poucos elementos com *clusters* com um número elevado de elementos em vez de unir *clusters* com número médio de elementos (Duda, 2001).

2.2. Índices de Semelhança

Um índice de semelhança é aplicado a partições de n padrões para determinar a que melhor representa o conjunto de dados original. Os índices de semelhança podem ser internos, se comparam a matriz de proximidade dos dados com cada partição, sem utilizarem informação à-priori sobre a natureza dos padrões, ou externos, quando são conhecidas as etiquetas à-priori de cada padrão (Dubes, 1987). A utilização de índices de semelhança é motivada pela frequente necessidade de conhecer o número natural de classes num conjunto de dados. Neste estudo serão aplicados dois índices internos de semelhança: o índice de Davies e Bouldin (Dubes, 1987) e o índice proposto por Xu, Kamath e Capson (Xu et al., 1993).

2.2.1. Índice de Davies e Bouldin

Sejam $\{x_1, x_2, \dots, x_n\}$ os padrões a classificar, sendo cada x_i um vector de dimensão d , tal como definido anteriormente. Uma solução para a classificação da imagem corresponde a uma partição $\{C_1, C_2, \dots, C_k\}$ dos n padrões, tal que $i \in C_k$ se x_i pertencer à classe k . O centro da classe k é um vector m_k de dimensão n , dado pela equação (4), onde n_k é o número de elementos da classe k .

O erro quadrático para a classe k , e_k^2 , é a soma das distâncias quadráticas entre cada um dos seus elementos e o centro da classe, equação (6). A distância entre dois vectores será calculada usando a distância Euclidiana.

$$e_k^2 = \sum_{i \in C_k} \|x_i - m_k\|^2 \quad (6)$$

Tendo sido estabelecida uma partição dos dados em k classes, $R_{i,k}$ dá uma medida da separabilidade *within to between* do par de classes (i,k) , equação (7). Os valores de $R_{j,k}$ serão tanto mais baixos quanto mais concentradas forem as classes j e k , e quanto mais separadas entre si elas estiverem.

$$R_{j,k} = \frac{e_j / \sqrt{n_j} + e_k / \sqrt{n_k}}{\|m_j - m_k\|} \quad (7)$$

O índice de Davies e Bouldin de uma partição K , $DB(K)$, é dado pela média dos valores máximos de R para cada classe, equação (8). Quanto menor o valor de $DB(K)$ melhor é a separabilidade entre as classes, sendo por isso mais adequada a partição K .

$$DB(K) = \frac{1}{K} \sum_{k=1}^K R_k, \text{ onde } R_k = \max_{j \neq k} (R_{j,k}) \quad (8)$$

2.2.2. Índice de Xu, Kamath e Capson

Sejam $\{x_1, x_2, \dots, x_n\}$ os padrões a classificar, $\{C_1, C_2, \dots, C_k\}$ uma partição do conjunto de padrões, m_k o centro da classe k e e_k^2 o erro quadrático para a mesma classe, definidos nas equações (4) e (6).

Considere-se a medida de semelhança “menor distância entre clusters” (Minimum of Between-Cluster Distance) (Xu et al., 1993) definida para uma partição em k classes, $M(k)$:

$M(k) = \min_{j < i} (d_w(j, i))$, $j, i = 1, 2, \dots, k$, onde $d_w(j, i)$ representa a distância de Ward entre as classes i e j , definida em (5).

O índice proposto por Xu, $E(h)$, é uma combinação entre as medidas do erro quadrático e da menor distância entre *clusters* que avalia cada partição K , comparando cada nível h (correspondente à partição K) da estrutura hierárquica com o nível seguinte, equação (9).

$$E(h) = \frac{M(h) - M(h+1)}{\sqrt{J(h)} - \sqrt{J(h+1)}}, \text{ onde } J(h) = \sum_{k=1}^K e_k^2 \quad (9)$$

Na representação gráfica do índice E em função do nível h , espera-se que o valor de h , h_m , a que corresponde o valor máximo de E corresponda ao nível onde se encontra o número natural de classes do conjunto (Xu et al., 1993). Por outras palavras, espera-se que a redução de classes na passagem do nível $h_m + 1$ para o nível h_m , provoque um aumento na distância entre *clusters*, sem ter ocorrido um aumento significativo do erro quadrático médio dos padrões associados a cada *cluster*.

Capítulo 3

Teste com Dados Sintéticos

Com o objectivo de avaliar o desempenho do classificador `clusterdata`, foram produzidos conjuntos de dados sintéticos que posteriormente foram classificados pelos diversos métodos de agregação hierárquica referidos no Capítulo 2.

3.1. Construção dos Conjuntos de Dados de Teste

Os dados sintéticos foram criados utilizando o MATLAB tendo sido produzidas representações gráficas a duas dimensões (figura 3.1).

Cada conjunto de dados de teste foi produzido utilizando um conjunto com 1000 elementos, a duas dimensões, com valores entre 0 e 1 (exclusive), agrupados em gaussianas com características adequadas a cada situação pretendida. Cada conjunto de dados é obtido a partir de um número pré-definido de pontos chave (parâmetro que será o centro de cada gaussiana) em torno dos quais se gera um determinado número de elementos, com uma distribuição gaussiana e uma dada dispersão.

O conjunto I, figura 3.1, foi produzido partindo de quatro pontos, em torno dos quais se geraram gaussianas com um desvio padrão de 0.04, em que os dois conjuntos mais à direita têm o dobro de elementos dos conjuntos da esquerda. Nesta imagem pretende-se avaliar se o classificador tem em consideração tanto a distância entre os centros de cada *cluster* como a sua dimensão (número de elementos).

O conjunto II, figura 3.1, tem como base 18 centros a partir dos quais se geraram grupos distribuídos por uma normal com desvio padrão de 0.025. Neste exemplo há claramente dois conjuntos distintos, em forma de 'C', sendo no entanto relativamente espalhados no espaço bidimensional. Para todos os elementos de cada 'C' há elementos do outro conjunto que distam de si menos do que a distância (Euclidiana) entre si e pelo menos um elemento do seu conjunto. Esta situação é potencialmente difícil de lidar para certos classificadores.

O conjunto III é um exemplo de 5 gaussianas, cada uma com 200 elementos, geradas com um desvio padrão de 0.035. Este conjunto apresenta 5 grupos relativamente densos e separados entre si, não devendo, em princípio, oferecer grandes dificuldades para qualquer um dos classificadores.

O conjunto IV é formado por 20 gaussianas, geradas com um desvio padrão de 0.02, e cujos centros estão dispostos em forma de "lua", com um pequeno conjunto no espaço entre os seus extremos. Este conjunto tem algumas semelhanças com o conjunto II na medida em que um dos grupos tem

elementos muito distantes entre si. Com este exemplo, pretende-se verificar a importância dada pelos classificadores à coesão existente em cada classe.

O conjunto V foi construído com 16 gaussianas, com um desvio padrão de 0.02. Este conjunto representa duas classes dispostas em forma de duas estrelas com cauda, e pretende verificar o comportamento dos classificadores face a uma situação de classes alongadas e pouco uniformes.

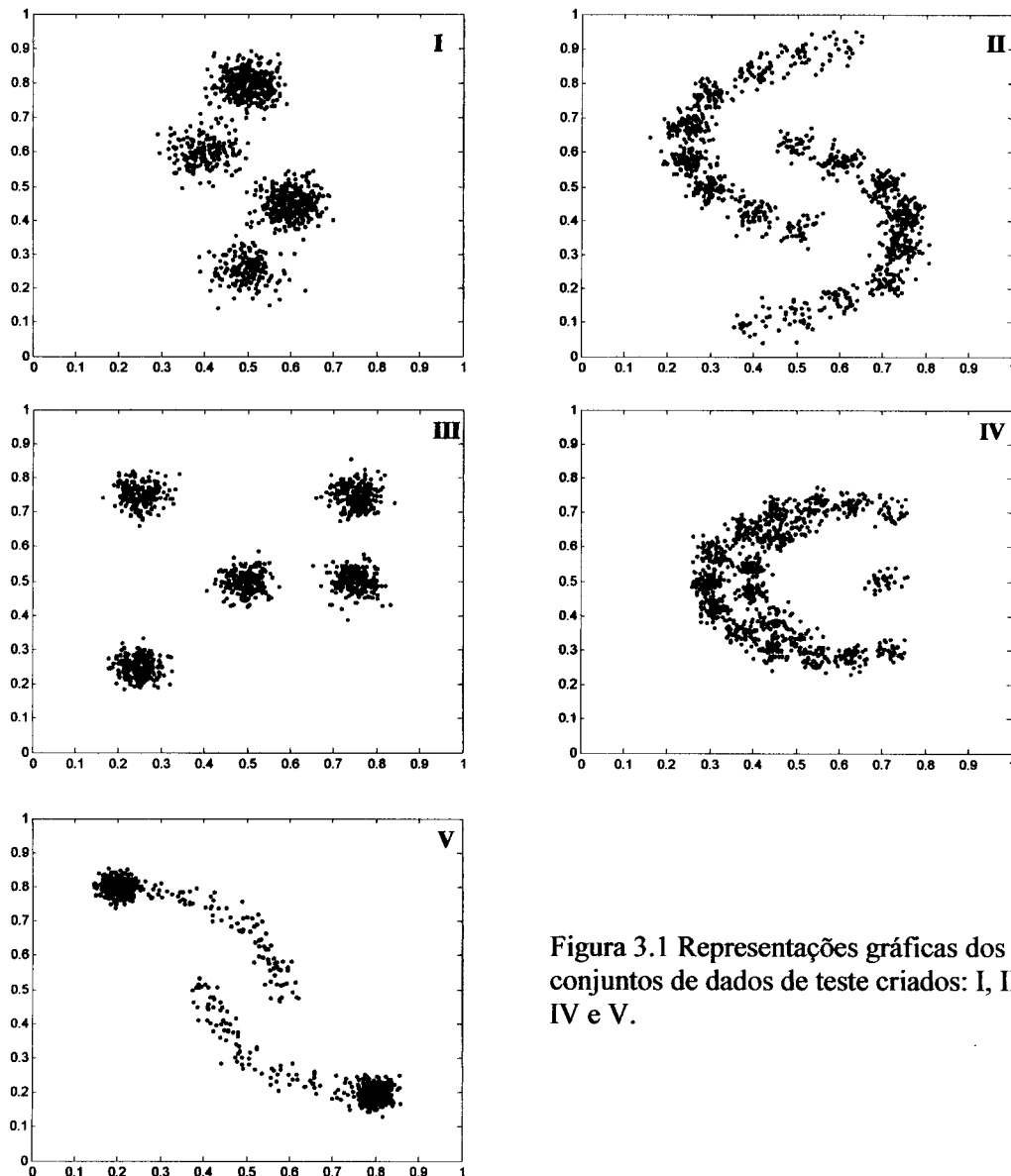


Figura 3.1 Representações gráficas dos conjuntos de dados de teste criados: I, II, III, IV e V.

Num panorama ideal, um classificador que tenha um comportamento adequado num conjunto de dados sintéticos, também o deverá ter para os restantes conjuntos de dados. Como se poderá comprovar na secção que se segue, esta é uma situação em tudo ideal, já que o desempenho de cada classificador varia de situação para situação, não existindo um único classificador adequado a todos os conjuntos de dados.

3.2. Classificação dos Conjuntos de Dados de Teste

Para cada conjunto de dados sintético foi aplicada a função `clusterdata` disponível no MATLAB mantendo fixo o parâmetro (`pdist`) da medida da distância entre os diferentes objectos, distância euclidiana. O parâmetro `maxclust` foi alterado de acordo com o número de classes pretendidas para cada aplicação e alterando o parâmetro (`linkage`) para a construção da árvore de agregação hierárquica. Neste último parâmetro foram estudados todos os métodos disponíveis no MATLAB: `single`, `complete`, `average`, `weighted`, `ward`, `centroid` e `median`. Para cada tipo de agregação, apresentam-se dois gráficos com os valores dos índices DB e Xu. O índice DB apresentado diz respeito à classificação em 2, 3, ..., 12 classes e o índice Xu é aplicado às classificações em 2, 3, ..., 11 classes. Dado que o índice Xu de uma partição é estabelecido comparando-a com a partição seguinte, as representações gráficas deste índice não apresentam valores para a partição em 12 classes. Nos gráficos dos índices DB e Xu o valor do índice para um número c de classes corresponde ao índice atribuído ao valor de k tal que $c=k+1$. Assim, nos gráficos dos índices DB e Xu, por exemplo, o valor 1 na abcissa corresponde a uma partição com 2 classes.

Nesta secção serão apresentados os resultados da aplicação de cada método de classificação a cada um dos 5 conjuntos de dados de teste e os respectivos gráficos dos índices DB e Xu.

3.2.1. Agregação pelo Método *Single*

O método *single*, também conhecido por método do “vizinho mais próximo” utiliza a medida da distância mínima simples entre dois *clusters*. Nas figuras que se seguem apresentam-se os resultados obtidos classificando cada um dos cinco conjuntos de dados teste descritos anteriormente pelo método *single* em 3 números de classes diferentes. Por exemplo, para o primeiro conjunto de dados, onde se espera obter 4 classes, são apresentados os resultados das classificações em 3, 4 e 5 classes (figura 3.2).

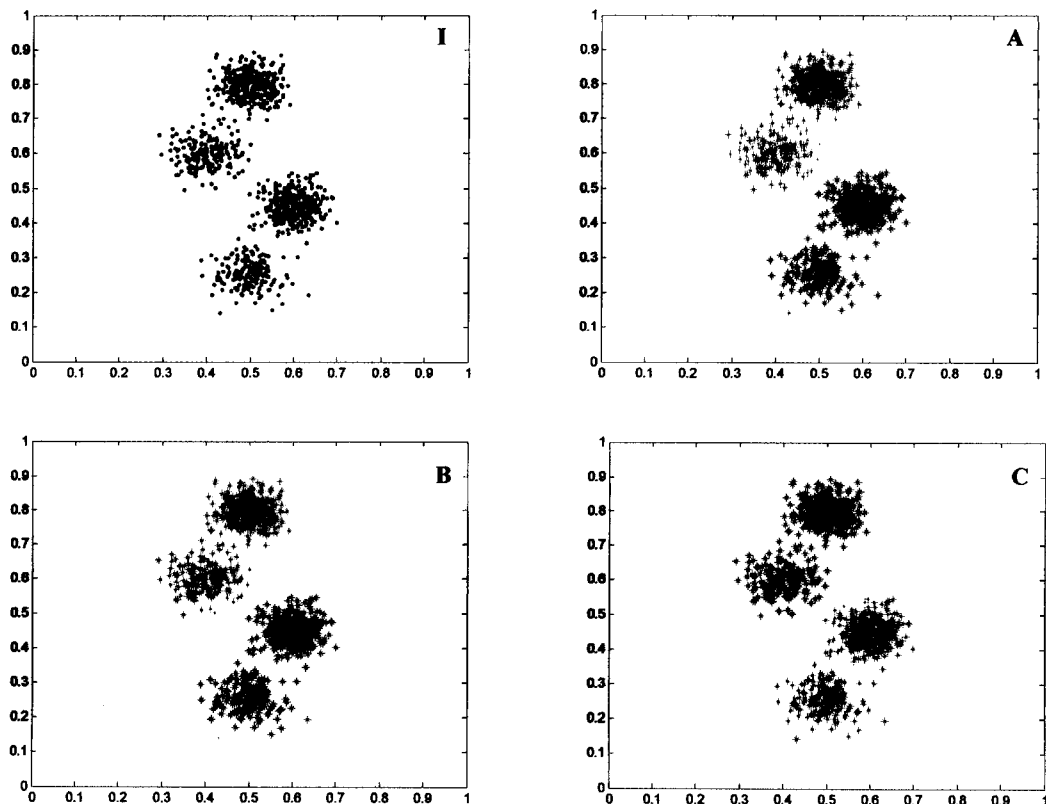


Figura 3.2 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.

Na figura 3.2, os resultados da aplicação do método em estudo, representados em A, B e C, revelam que o facto de a agregação ser baseada nas distâncias mínimas entre *clusters* impede o reconhecimento dos 4 conjuntos esperados. O classificador, nesta situação, opta por separar pontos isolados para formarem *clusters* sem grande relevância.

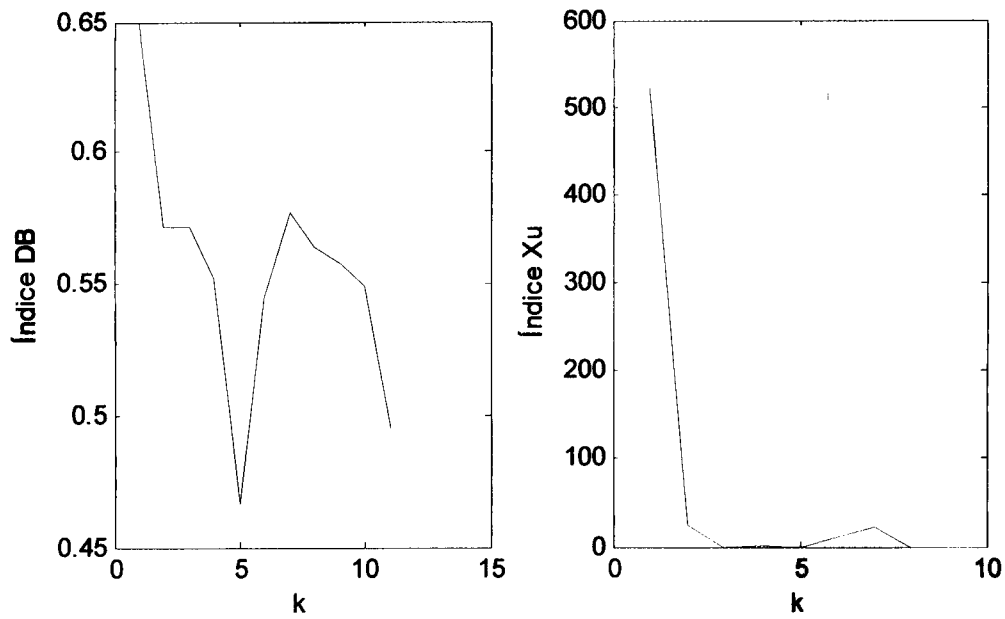


Figura 3.3 Representação dos valores dos índices DB e Xu, para o conjunto I.

Pelo índice DB a melhor classificação será a que compreende 6 classes, enquanto que pelo índice Xu, a imagem possui 2 classes naturais (figura 3.3).

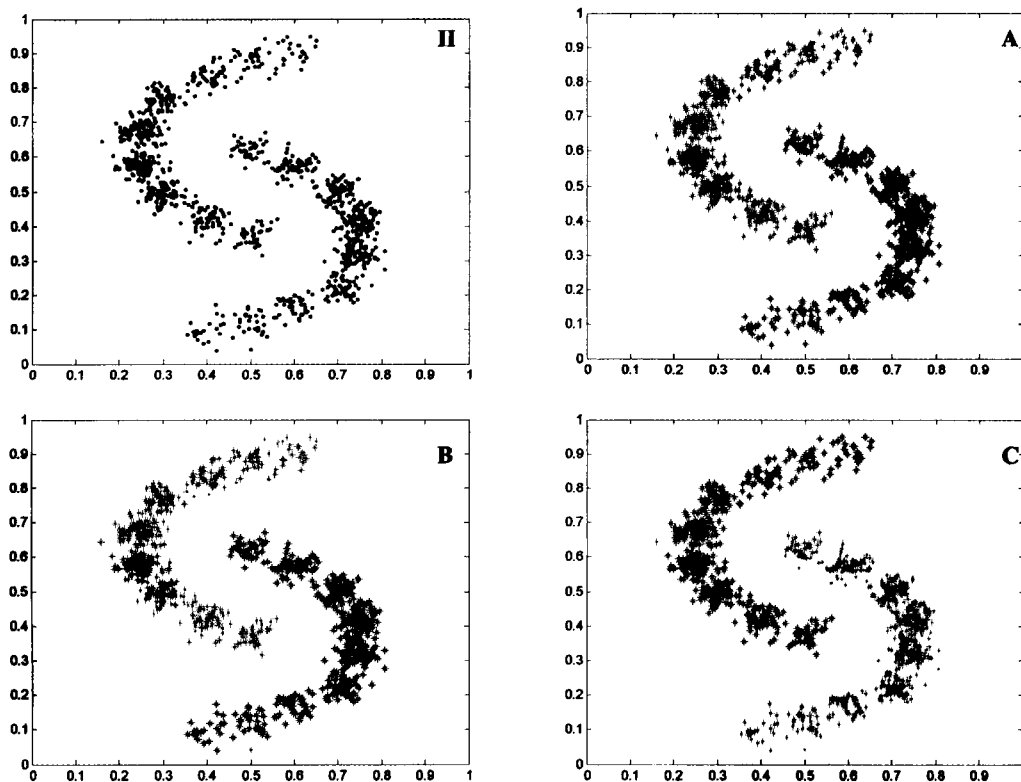


Figura 3.4 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.

No conjunto de dados II o objectivo da classificação foi atingido uma vez que para duas classes o método de agregação *single* reconhece as duas classes esperadas nestes dados. Para mais classes repete-se o sucedido no conjunto anterior, formando classes com apenas um elemento, aquele que estiver mais isolado relativamente aos restantes elementos.

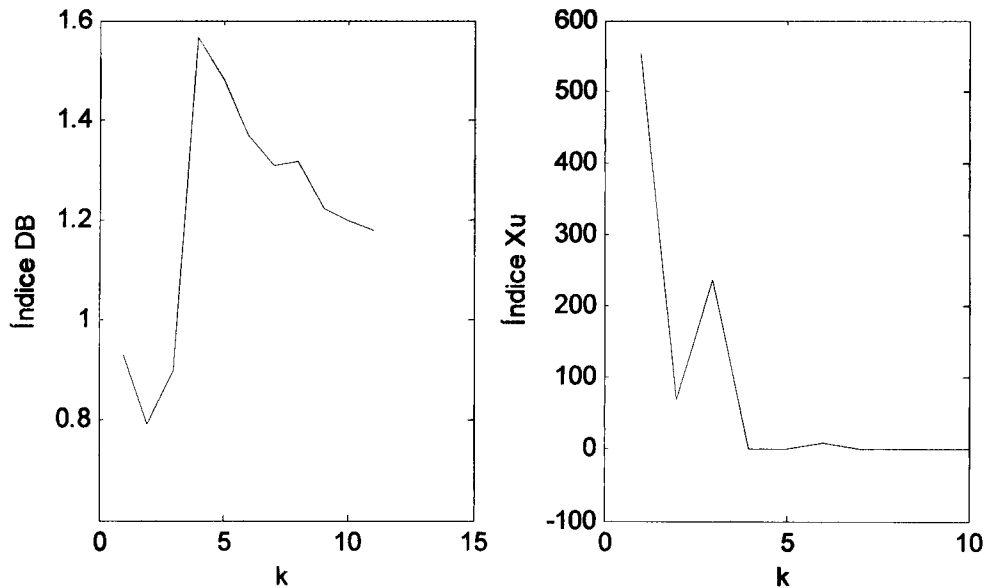


Figura 3.5 Representação dos valores dos índices DB e Xu, para o conjunto II.

Mais uma vez os índices de semelhança produzem resultados divergentes: o índice DB elege a separação em 3 classes enquanto que o índice Xu considera a classificação natural dos dados (2 classes) como a mais eficiente (figura 3.5).

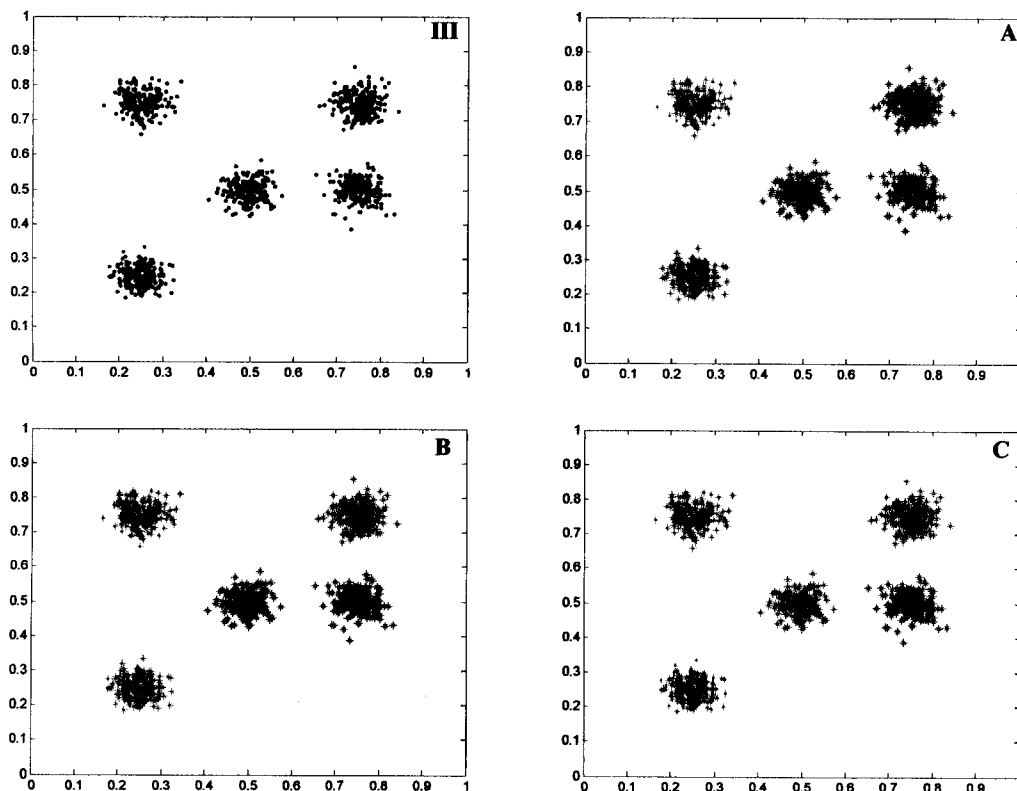


Figura 3.6 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.

Como esperado, a classificação do conjunto de dados III para a partição em 5 classes (figura 3.6 C) capta de forma satisfatória a sua estrutura natural. Os 5 *clusters* produzidos pelas 5 gaussianas são correctamente identificados pelo classificador e as escolhas para 3 e 4 classes são também satisfatórias. Neste caso, o desempenho dos dois índices é semelhante e corresponde ao pretendido, ambos considerando que 5 é o número de classes naturais do conjunto III (figura 3.7).

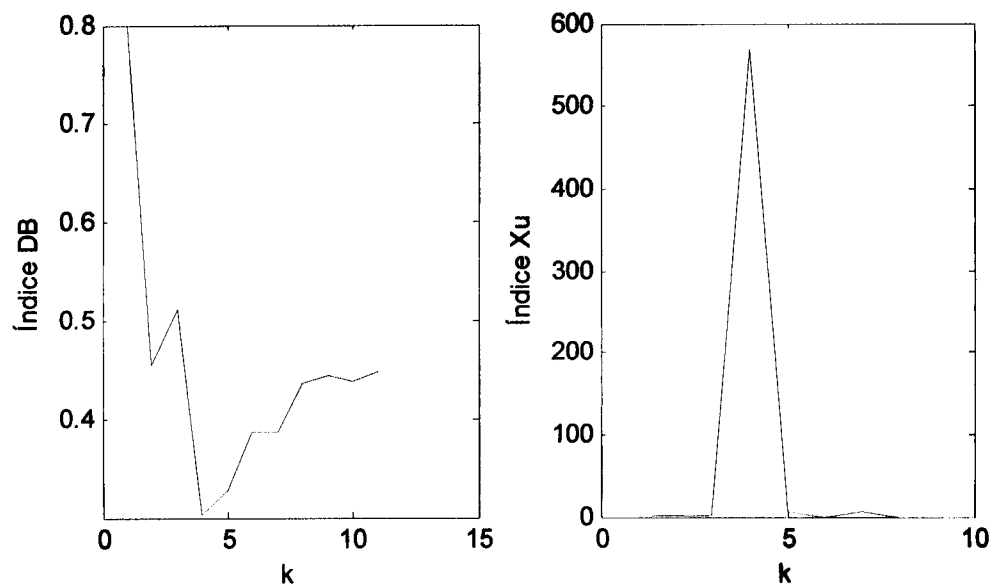


Figura 3.7 Representação dos valores dos índices DB e Xu, para o conjunto III.

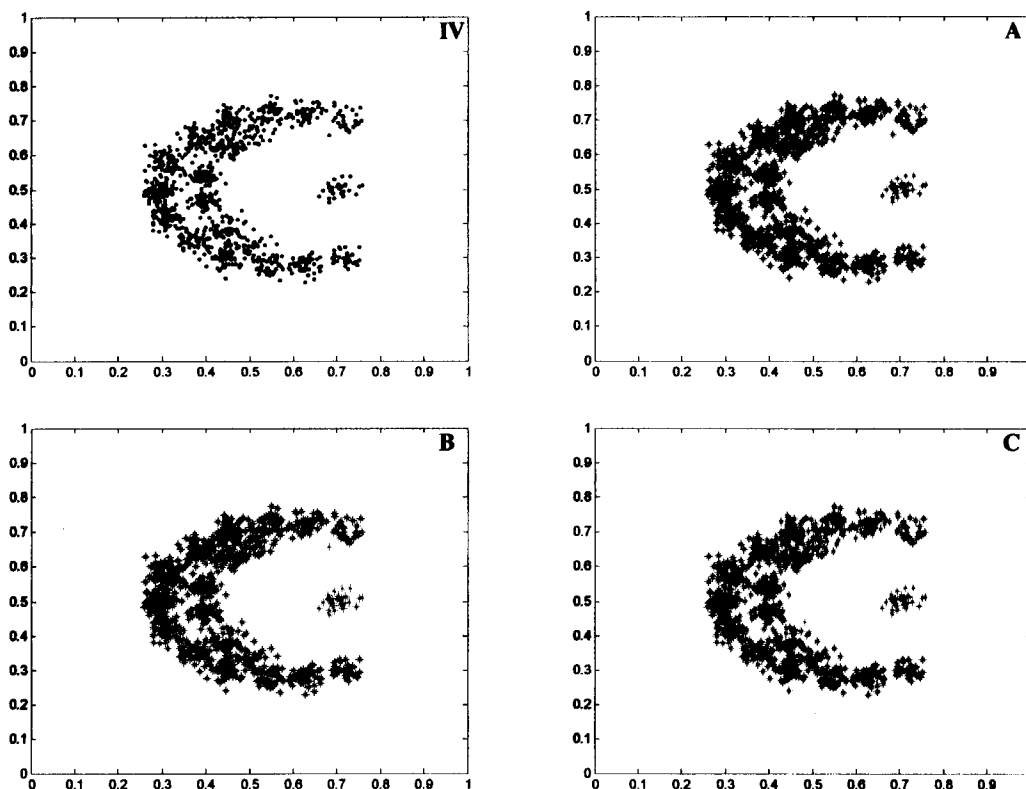


Figura 3.8 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.

Para o conjunto de teste III, o método de agregação *single* reconhece, mais uma vez, a existência das duas classes presentes nos dados. Por considerar a distância mínima entre *clusters*, atinge o resultado pretendido separando a classe em forma de lua da classe mais reduzida que se situa entre os seus extremos.

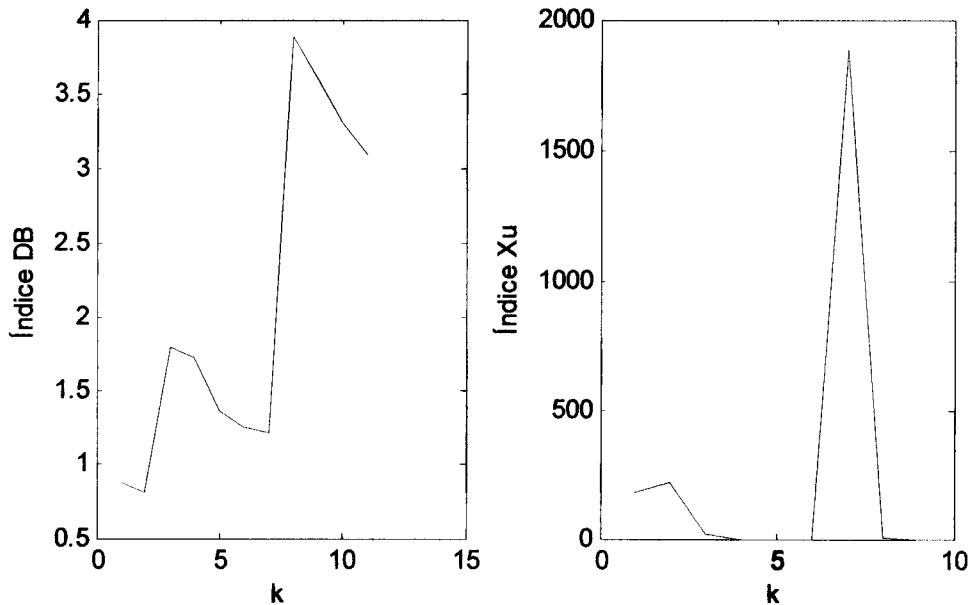


Figura 3.9 Representação dos valores dos índices DB e Xu, para o conjunto IV.

Neste exemplo, os índices voltam a ser concordantes, escolhendo a partição em 8 classes como a mais eficiente. No entanto, o número esperado de classes nos dados é 2, o que não corresponde à escolha feita pelos dois índices (figura 3.9).

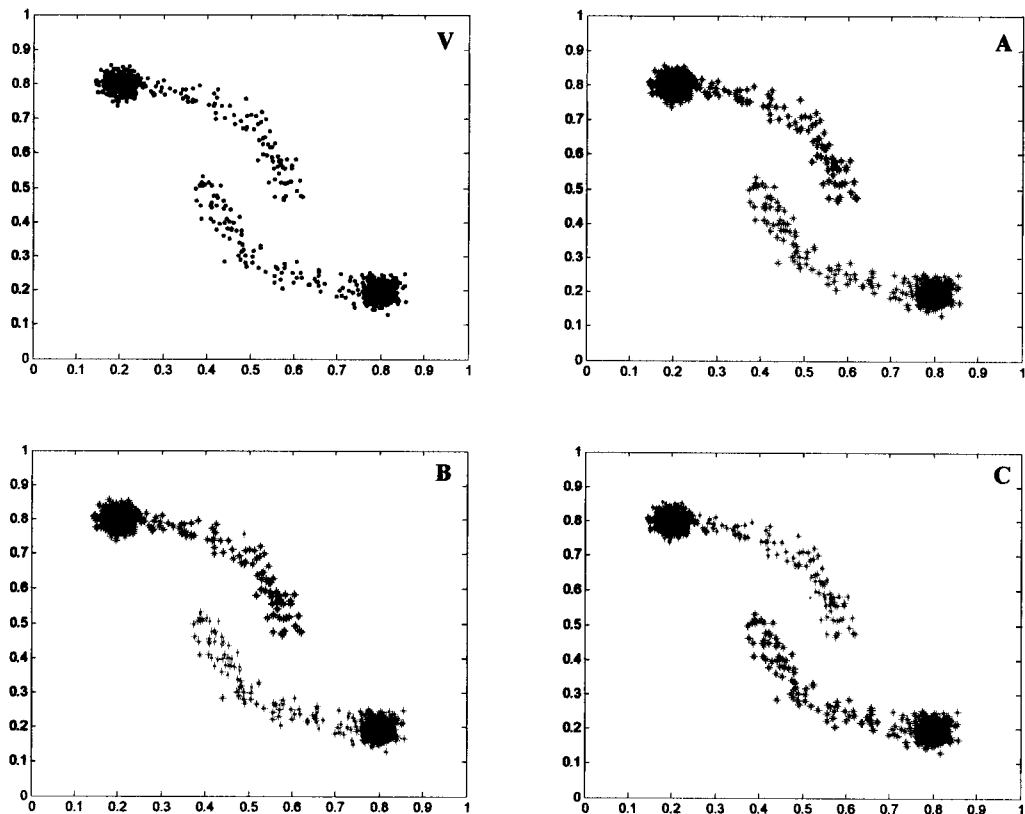


Figura 3.10 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.

O conjunto de dados V da figura 3.10 é classificado correctamente em duas classes (A). O desempenho do método é mais uma vez favorável ao número esperado de classes existentes nos dados. A classificação em 3 classes não é satisfatória, tendo surgido uma classe com apenas 1 elemento, mas a separação que ocorre de 3 para 4 classes parece muito adequada.

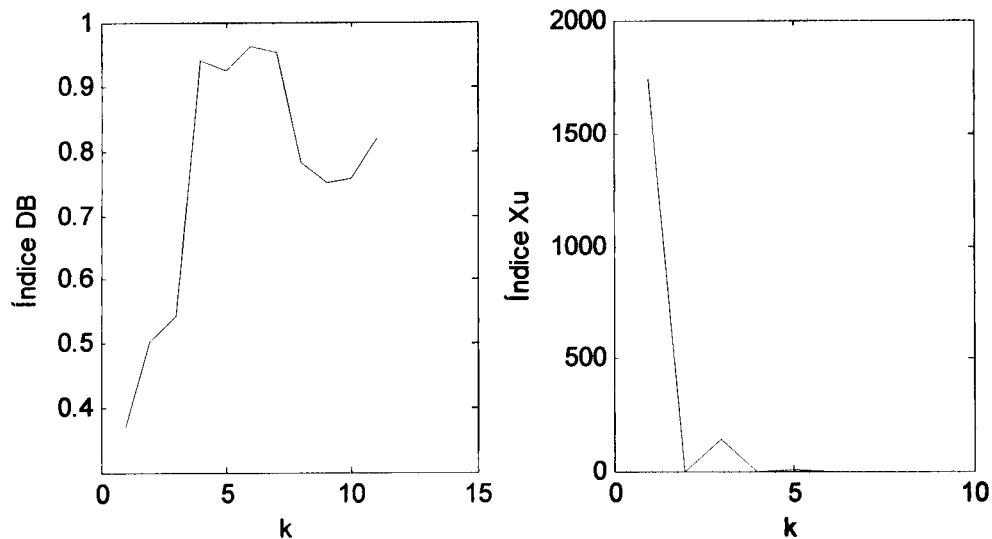


Figura 3.11 Representação dos valores dos índices DB e Xu, para o conjunto V .

Os índices de semelhança voltam a apresentar resultados coerentes (figura 3.11) e de acordo com a partição dos dados esperada, ou seja, ambos consideram a partição em 2 classes como sendo a mais eficaz.

3.2.2. Agregação pelo Método *Complete*

O método *complete* também designado por método do “vizinho mais afastado”, utiliza a medida da distância máxima entre dois *clusters*.

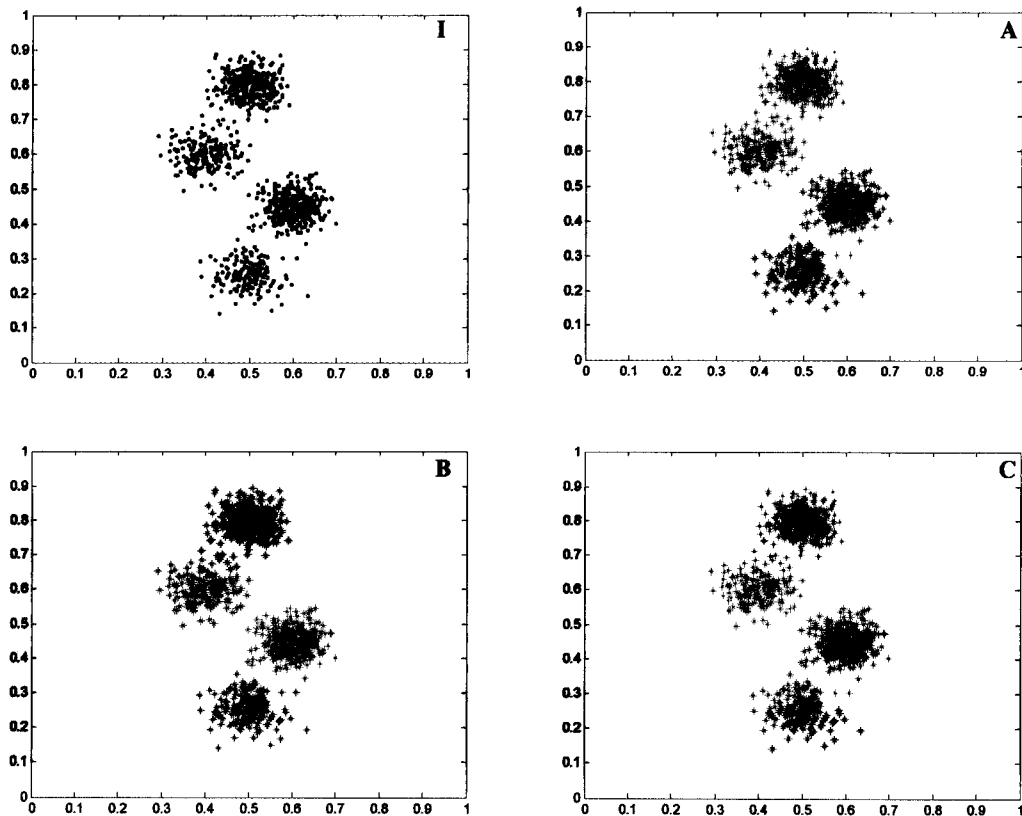


Figura 3.12 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.

Pelo método de agregação *complete*, a partição obtida para 4 classes, figura 3.12 (B), corresponde à partição esperada dos dados. Considerando a distância máxima em lugar da distância mínima, este método tem um desempenho melhor do que o método de agregação *single* neste tipo de dados, já que na classificação com 4 classes, a distância máxima entre conjuntos é diminuída ao separar as quatro classes naturais e aumentada ao formar classes com apenas um elemento.

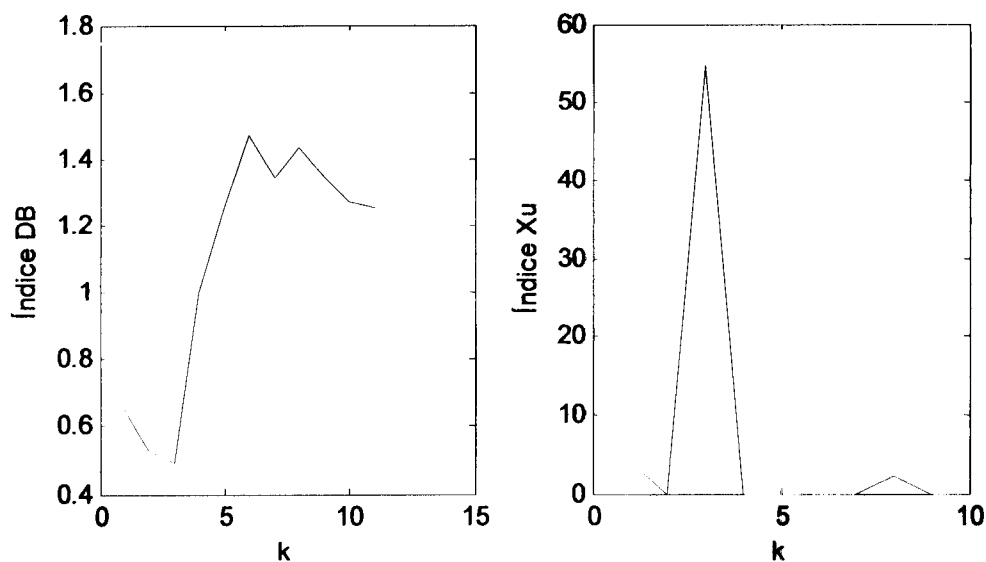


Figura 3.13 Representação dos valores dos índices DB e Xu, para o conjunto I.

O desempenho dos índices DB e Xu é semelhante e corresponde ao pretendido, ambos considerando que 4 é o número de classes naturais do conjunto de dados em estudo (figura 3.13).

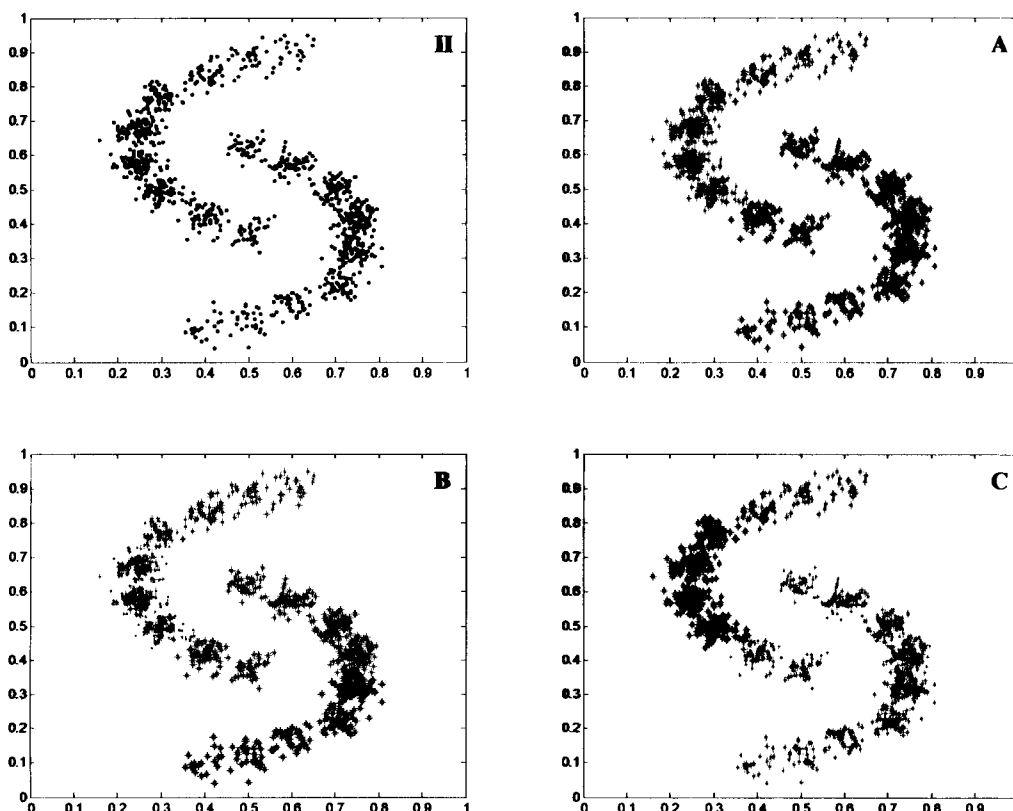


Figura 3.14 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.

Neste conjunto de dados o método de classificação revela-se incapaz de encontrar a partição pretendida dos dados em 2 classes. Ao pretender uma diminuição da distância máxima, este método prejudica os grupos de dados com distribuição alongada subdividindo-os em classes mais pequenas e, por isso, desajustadas.

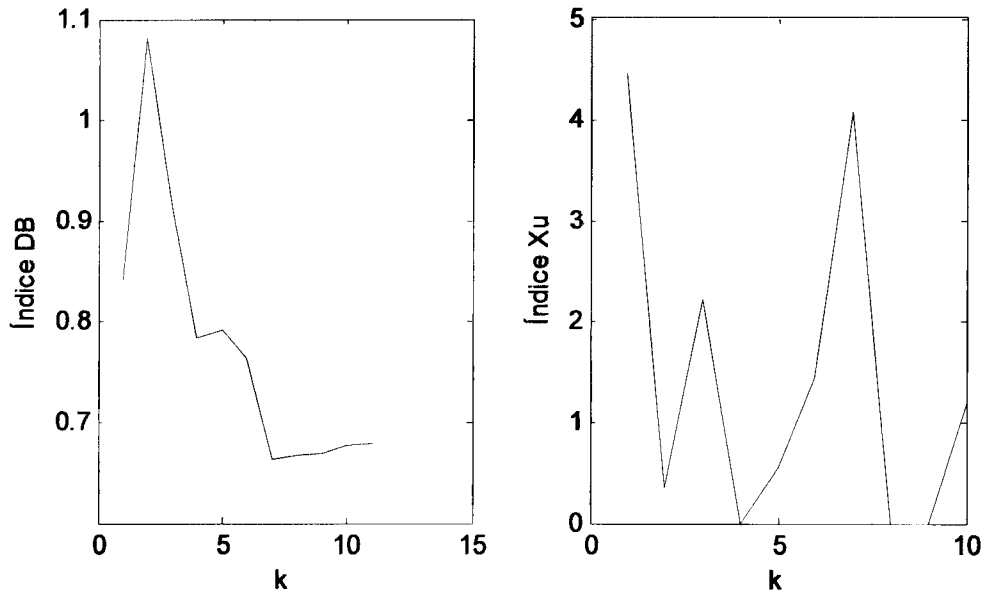


Figura 3.15 Representação dos valores dos índices DB e Xu, para o conjunto II.

Apesar dos resultados a nível de classificação, ambos os índices elegem a partição em 8 classes como a mais adequada (figura 3.15).

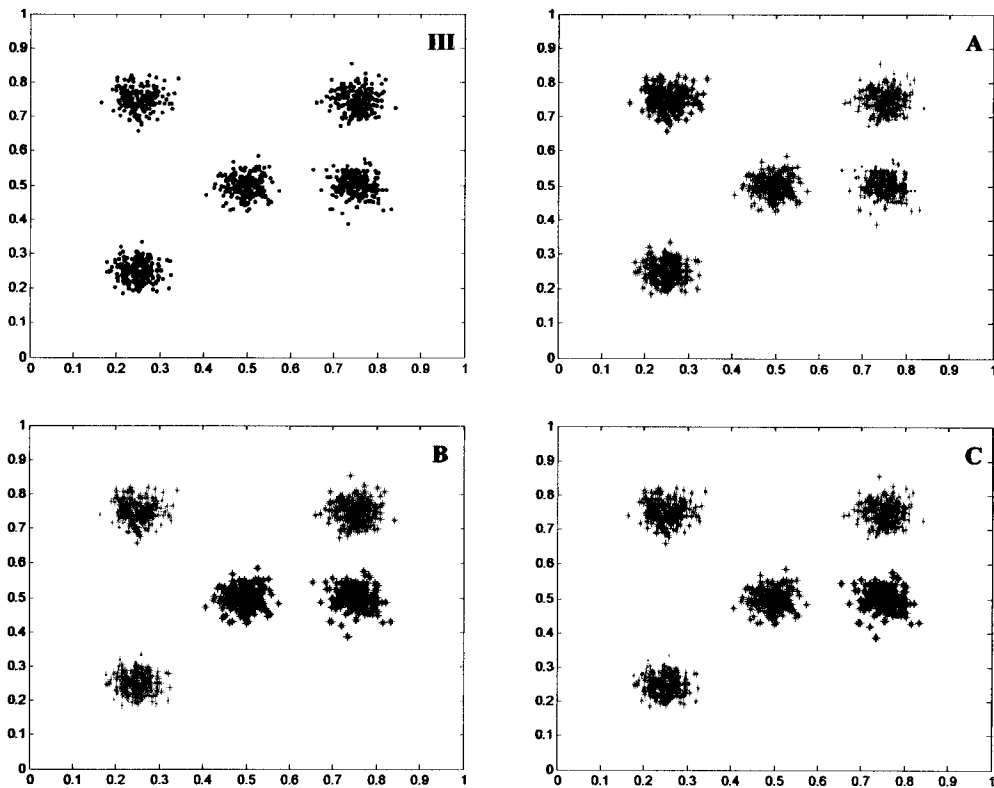


Figura 3.16 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.

O conjunto de dados III é classificado correctamente segundo este método de agregação, para todas as partições (figura 3.16).

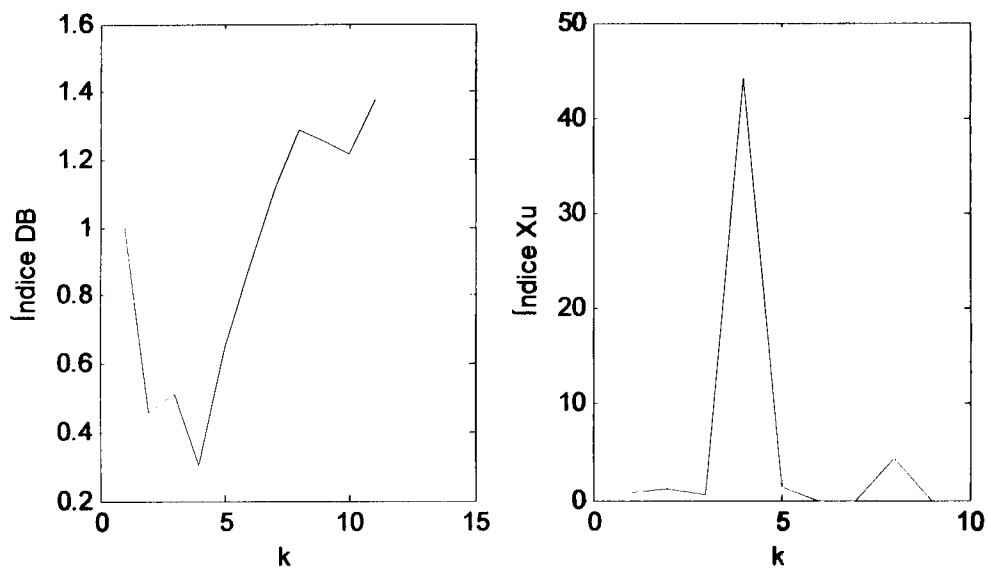


Figura 3.17 Representação dos valores dos índices DB e Xu, para o conjunto III.

O desempenho dos dois índices é semelhante e corresponde ao pretendido, ambos consideram que 5 é o número natural de classes do conjunto de dados em estudo (figura 3.17).

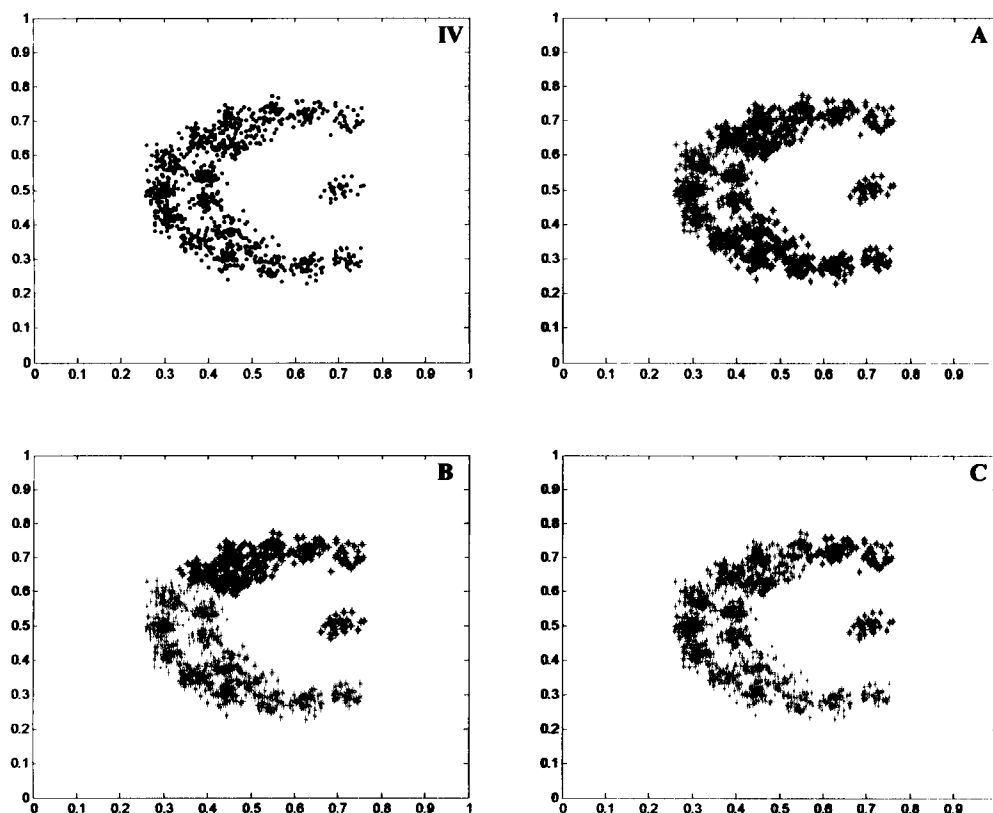


Figura 3.18 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.

Os resultados do método de agregação *complete* neste conjunto de dados revelam mais uma vez a incapacidade em reconhecer conjuntos de dados com distribuição alongada.

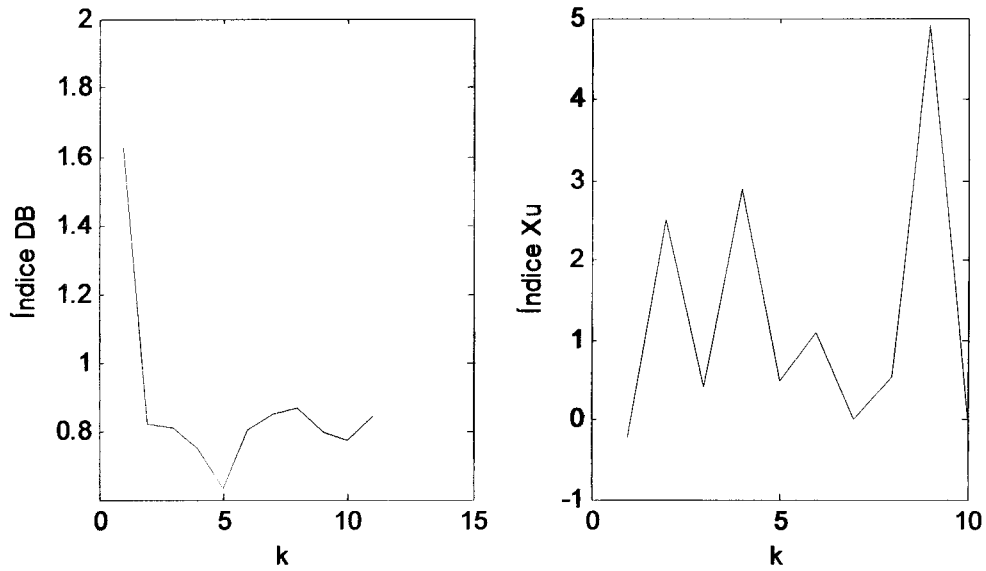


Figura 3.19 Representação dos valores dos índices DB e Xu, para o conjunto IV.

A classificação do conjunto de dados IV produz resultados incoerentes: o índice DB considera que 6 é o número natural de dados enquanto que o índice Xu escolhe a partição em 10 classes (figura 3.19).

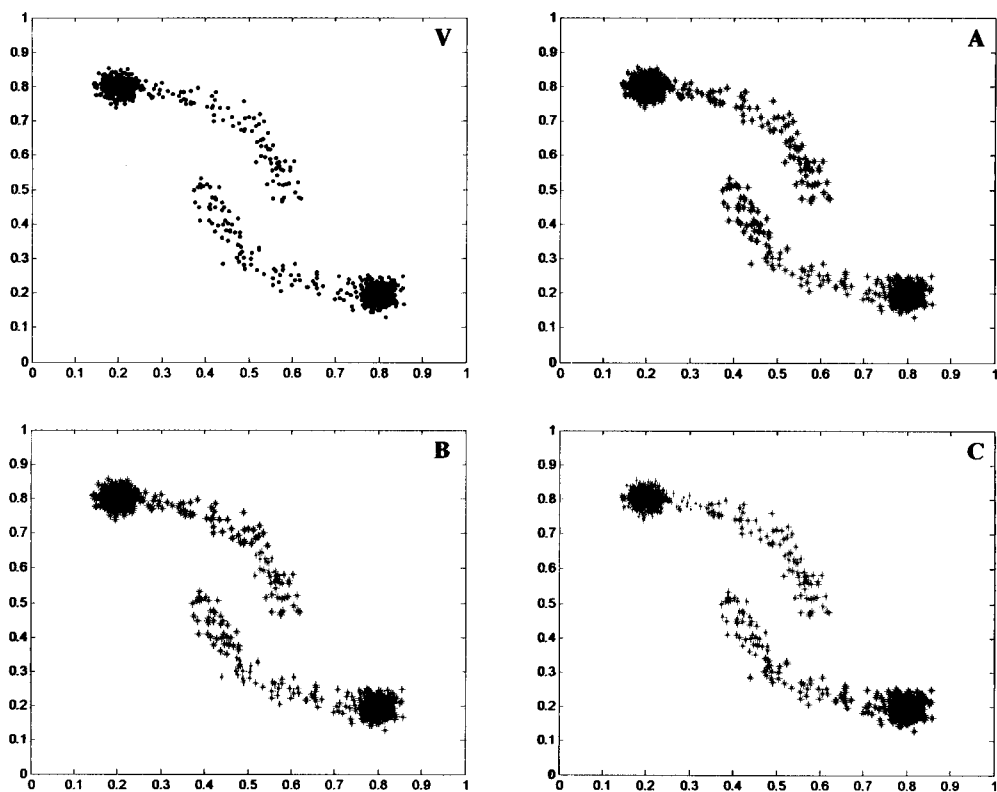


Figura 3.20 A Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.

Os resultados para o conjunto de dados V são de certa forma análogos aos do conjunto de dados IV, para este classificador (figuras 3.20 e 3.18).

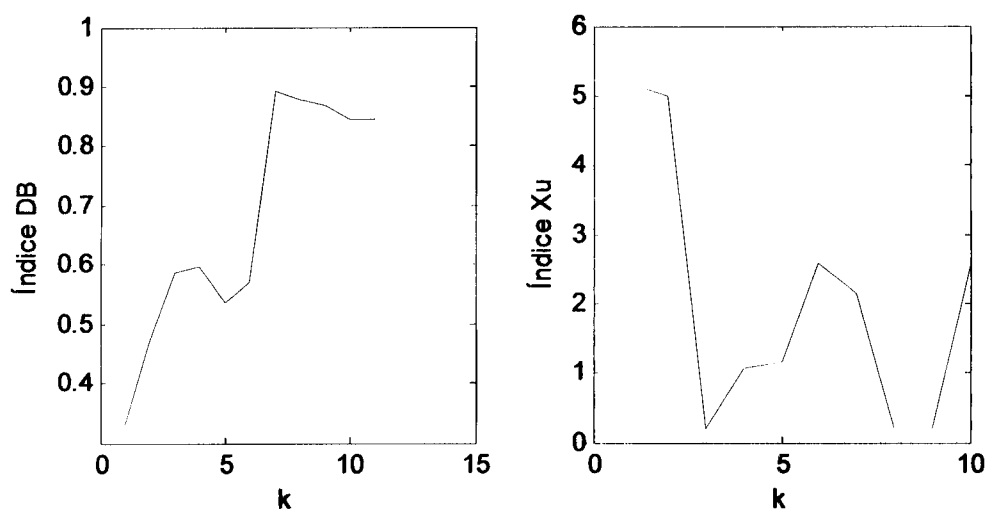


Figura 3.21 Representação dos valores dos índices DB e Xu, para o conjunto V.

Para a classificação efectuada, ambos os índices consideram a partição em duas classes como a mais adequada (figura 3.21). No entanto, conforme se pode ver na figura 3.20 (A), a partição em 2 classes produzida por este método não é a que se esperava à partida.

3.2.3. Agregação pelo Método *Average*

Ao aplicar o método *average* a árvore hierárquica é construída com base na média das distâncias entre cada par de objectos pertencentes a classes diferentes.

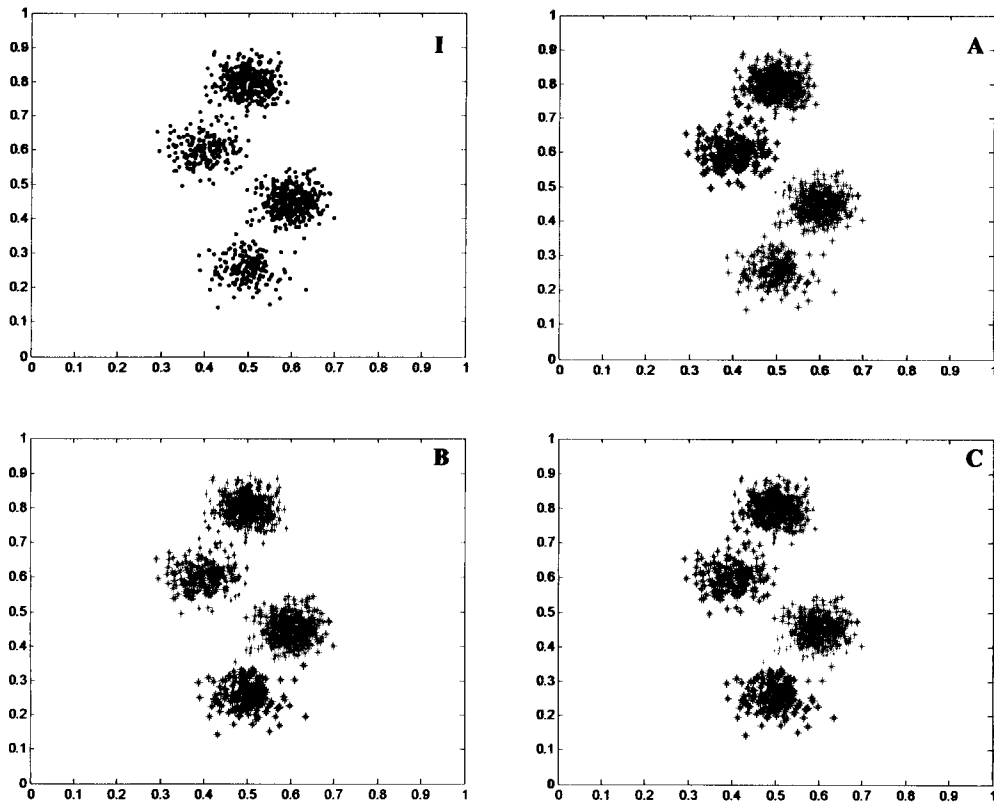


Figura 3.22 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.

Tal como no método de agregação *complete*, o conjunto de dados I da figura 3.22 é classificado em quatro classes (B) de forma a apresentar as quatro classes naturais presentes nos dados.

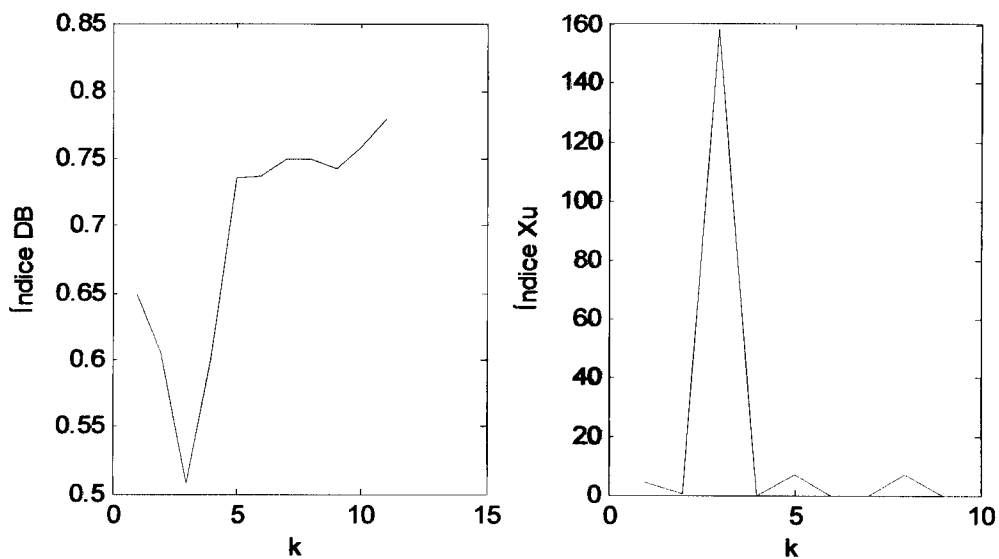


Figura 3.23 Representação dos valores dos índices DB e Xu, para o conjunto I.

Analogamente ao verificado no método *complete* para o conjunto de dados I, o desempenho dos dois índices é semelhante e corresponde ao pretendido, considerando 4 como o número de classes naturais do conjunto de dados em estudo (figura 3.23).

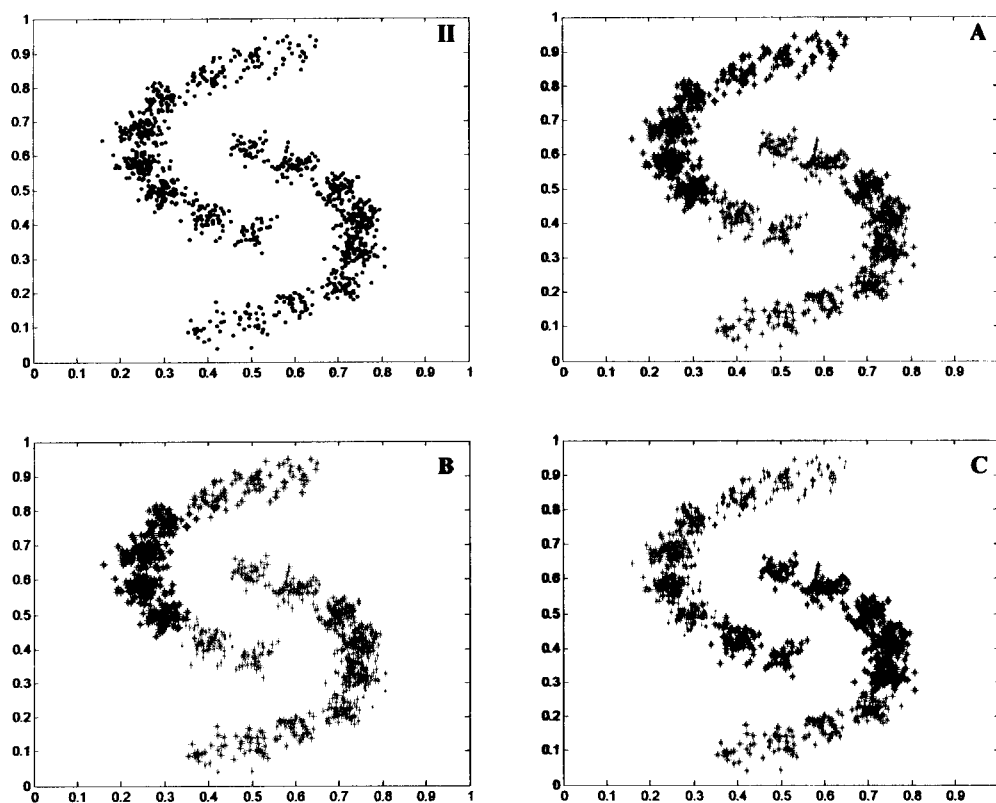


Figura 3.24 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.

Para o conjunto de dados II (figura 3.24) o método de agregação *average* tem também um desempenho semelhante ao do método *complete* (figura 3.14), falhando na captação da correcta distribuição das classes apresentadas.

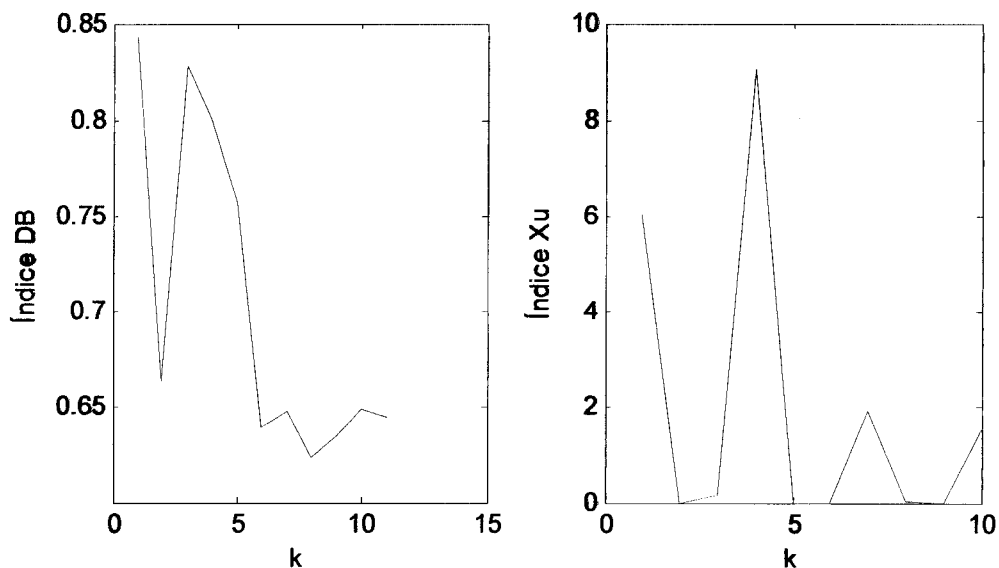


Figura 3.25 Representação dos valores dos índices DB e Xu, para o conjunto II.

Os índices DB e Xu apresentam resultados diferentes: o 1º elege a partição em 9 classes enquanto que o 2º considera a classificação em 5 classes como a mais vantajosa (figura 3.25).

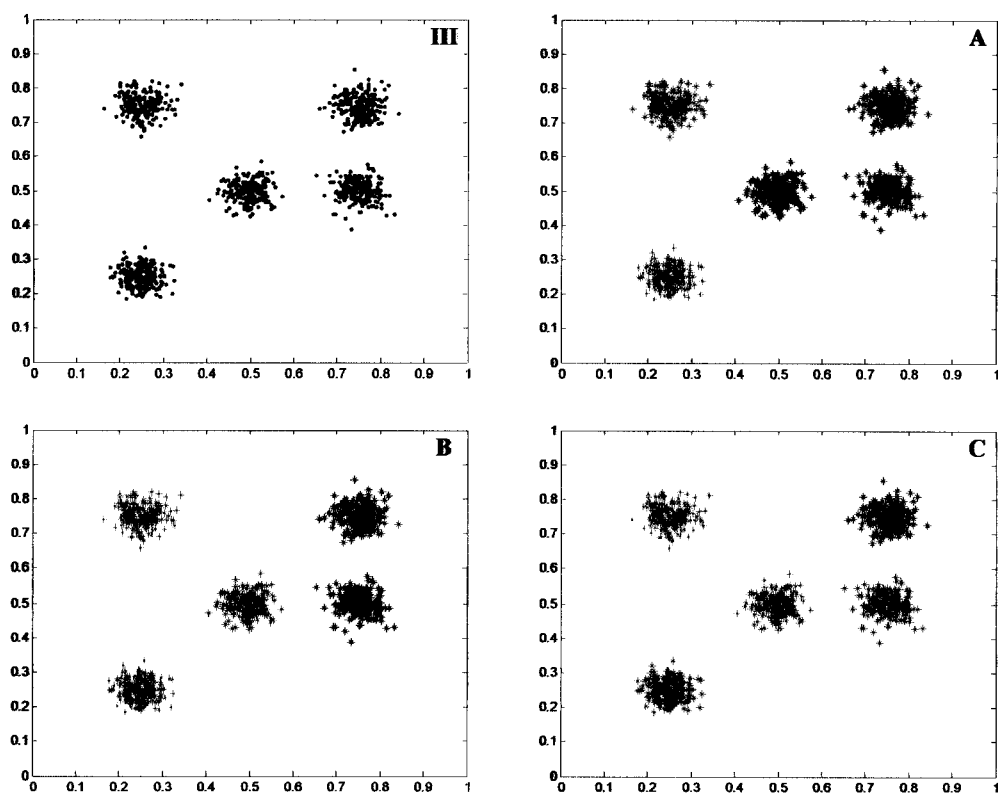


Figura 3.26 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.

Os resultados apresentados nas figuras 3.26 e 3.27 indicam que, para o conjunto de dados III, o método *average* tem um comportamento análogo ao dos métodos anteriormente expostos.

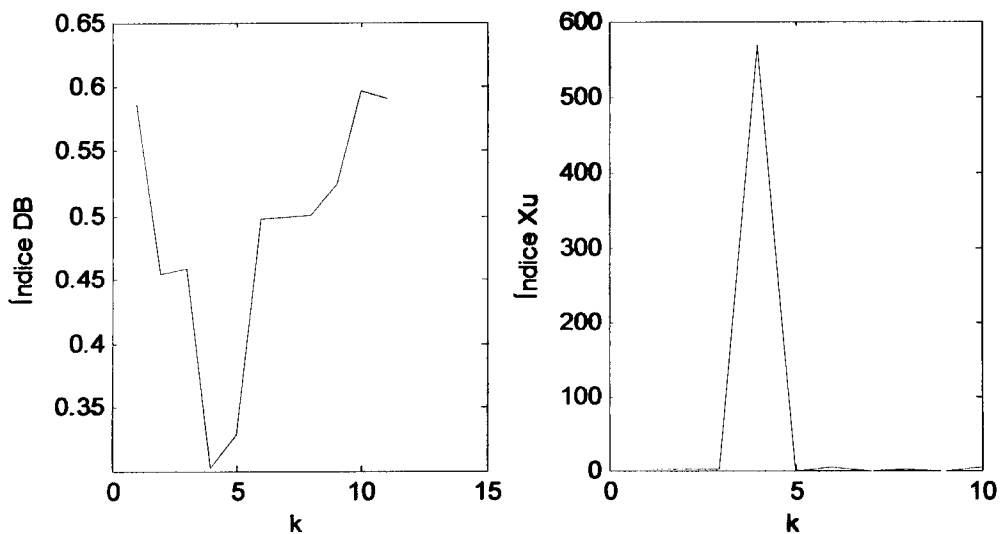


Figura 3.27 Representação dos valores dos índices DB e Xu, para o conjunto III.

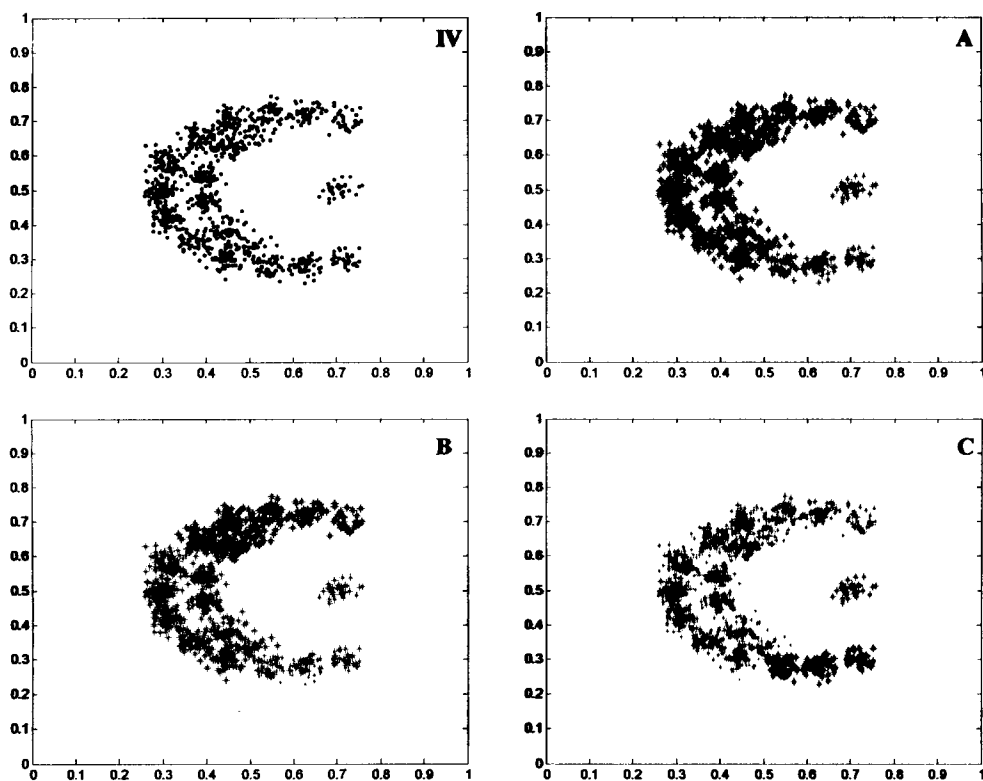


Figura 3.28 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.

Para o conjunto de dados de teste IV, o método *average* demonstra mais uma vez a sua incapacidade em reconhecer conjuntos de dados alongados, produzindo um resultado pouco satisfatório para a partição em 2 classes (figura 3.28 A).

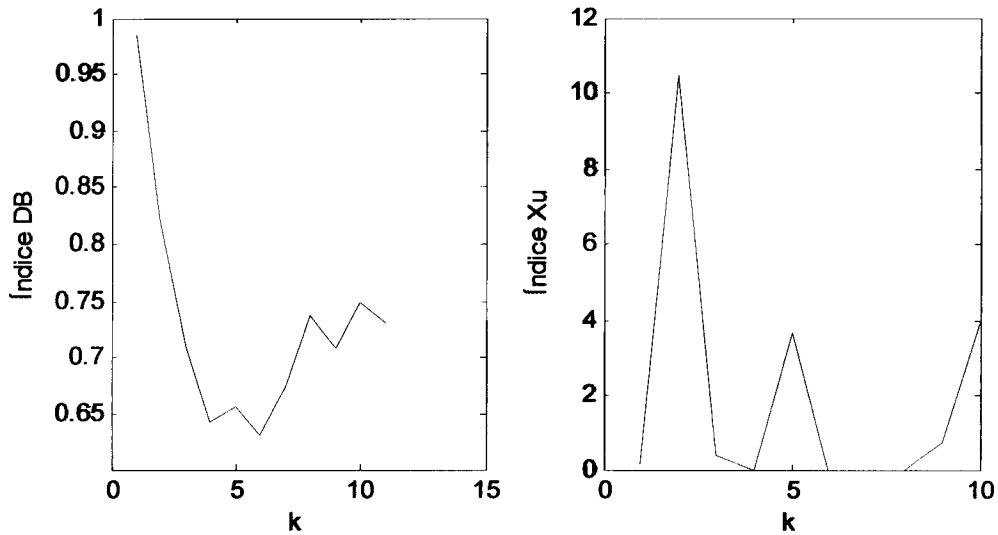


Figura 3.29 Representação dos valores dos índices DB e Xu, para o conjunto IV.

Neste caso o índice DB considera que 7 é o número natural de classes nos dados enquanto que o índice Xu elege a classificação em 3 classes como a mais eficiente (figura 3.29).

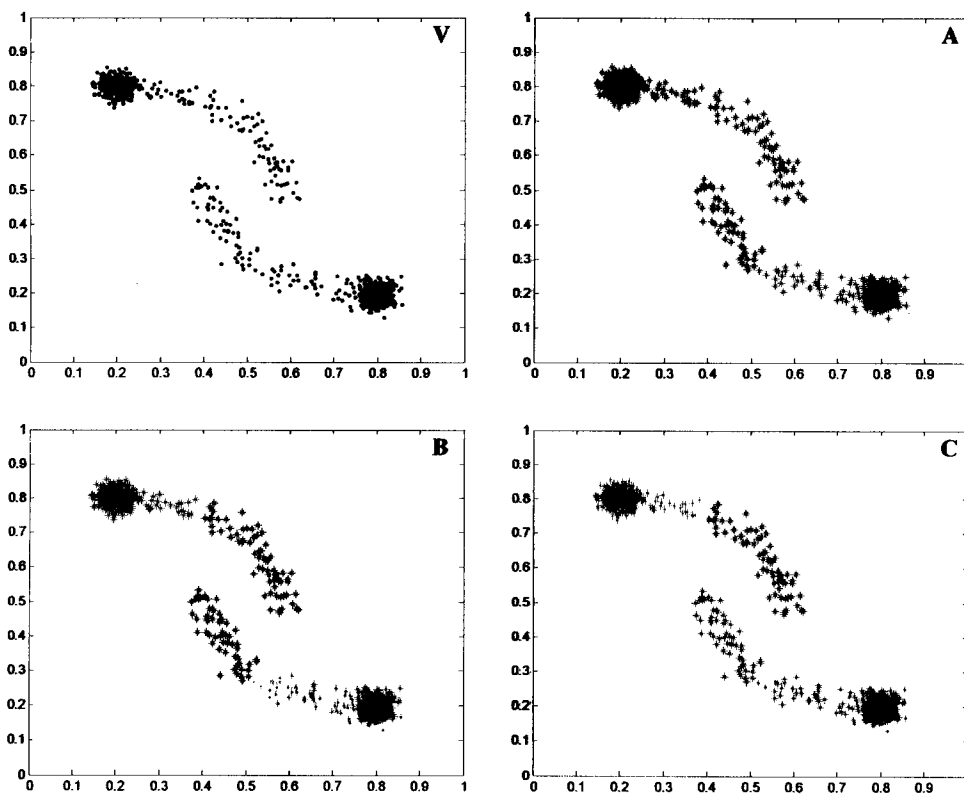


Figura 3.30 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.

Os resultados para o conjunto de teste V, apresentados na figura 3.30, indicam que o método em questão tem um comportamento análogo ao do método anteriormente exposto, falhando na identificação das classes previstas do conjunto de dados da imagem original.

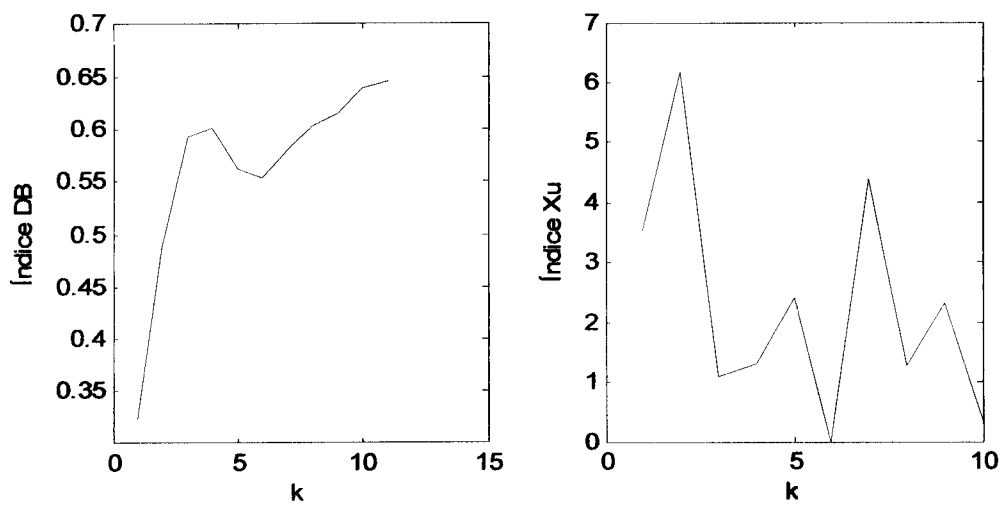


Figura 3.31 Representação dos valores dos índices DB e Xu, para o conjunto V.

A classificação do conjunto V produz resultados incoerentes: o índice DB considera que 2 é o número natural de dados enquanto que o índice Xu escolhe a partição em 3 classes (figura 3.31).

3.2.4. Agregação pelo Método *Weighted*

Este método é análogo ao método *average* sendo a distância utilizada, ponderada de acordo com o número de elementos de cada *cluster*.

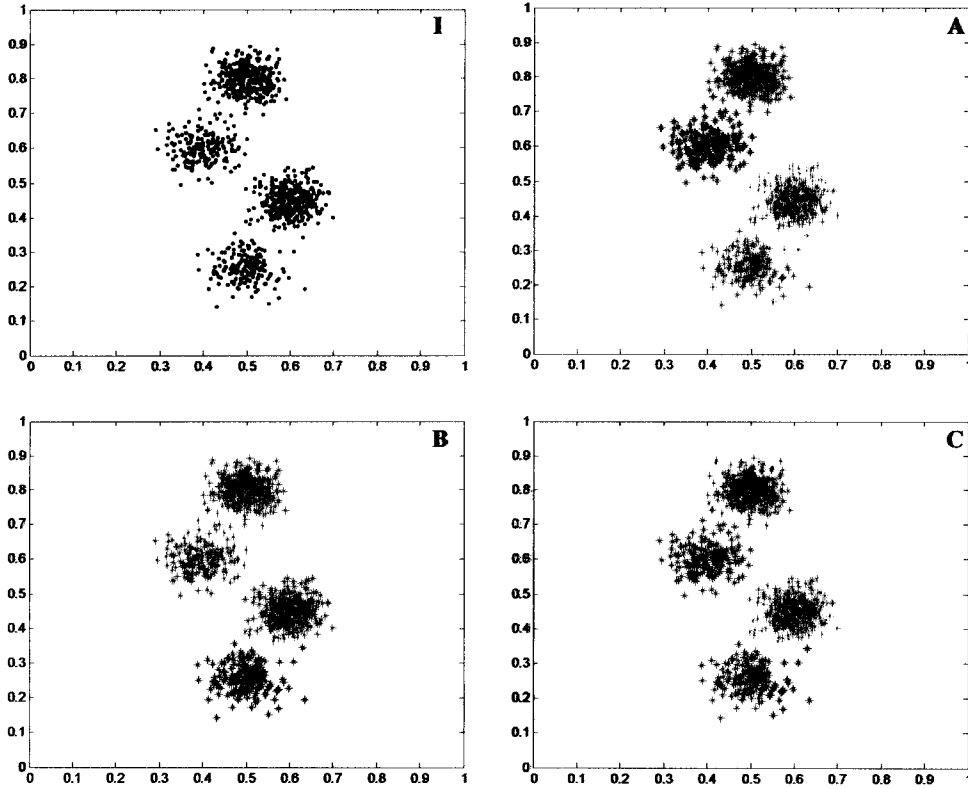


Figura 3.32 Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.

O método de agregação *weighted* produz resultados satisfatórios face ao conjunto I (figura 3.32) já que na classificação em 4 classes reconhece as classes naturais presentes nos dados.

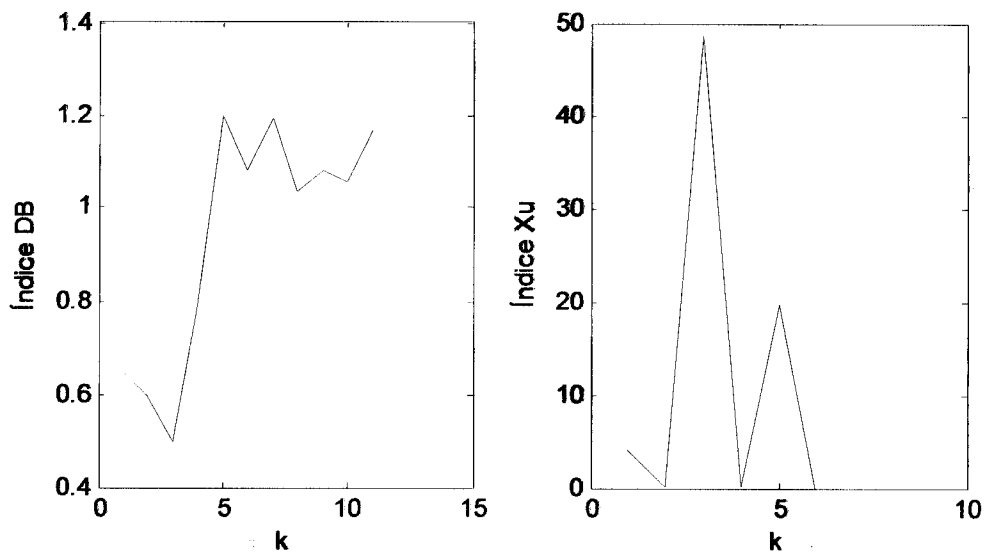


Figura 3.33 Representação dos valores dos índices DB e Xu, para o conjunto I.

O desempenho dos dois índices é semelhante e corresponde ao pretendido, ambos consideram que 4 é o número de classes naturais do conjunto de dados em estudo (figura 3.33).

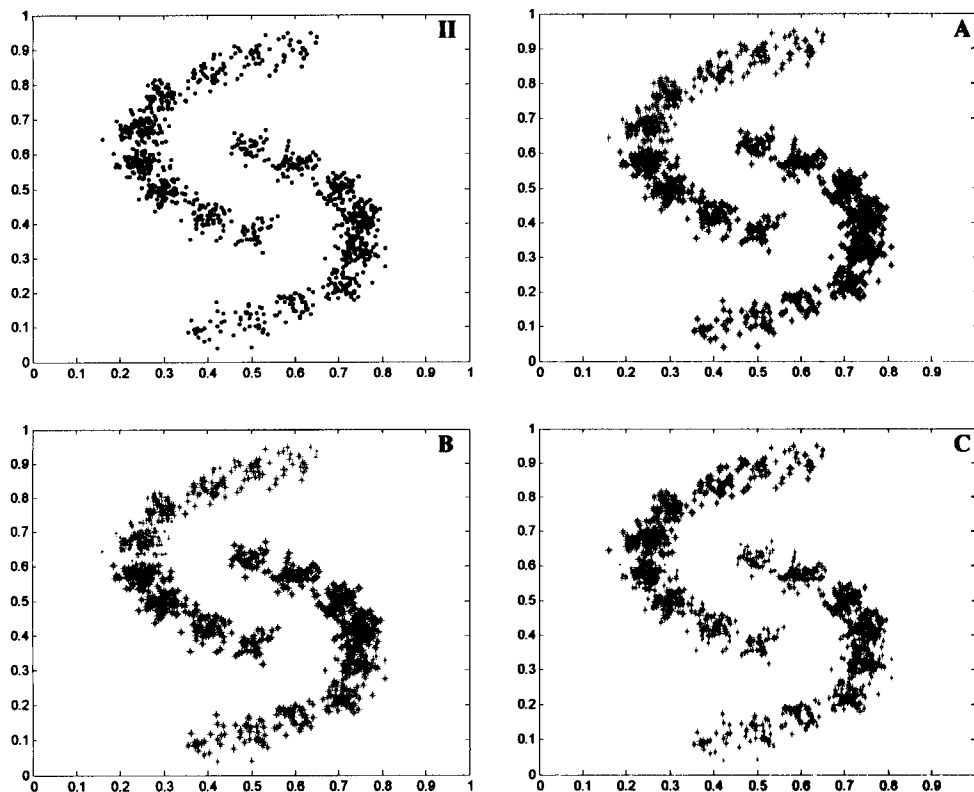


Figura 3.34 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.

No conjunto de dados sintéticos II (figura 3.34) o método de agregação *weighted* tem um desempenho semelhante ao dos métodos *complete* e *average*, falhando na captação da correcta distribuição das classes apresentadas.

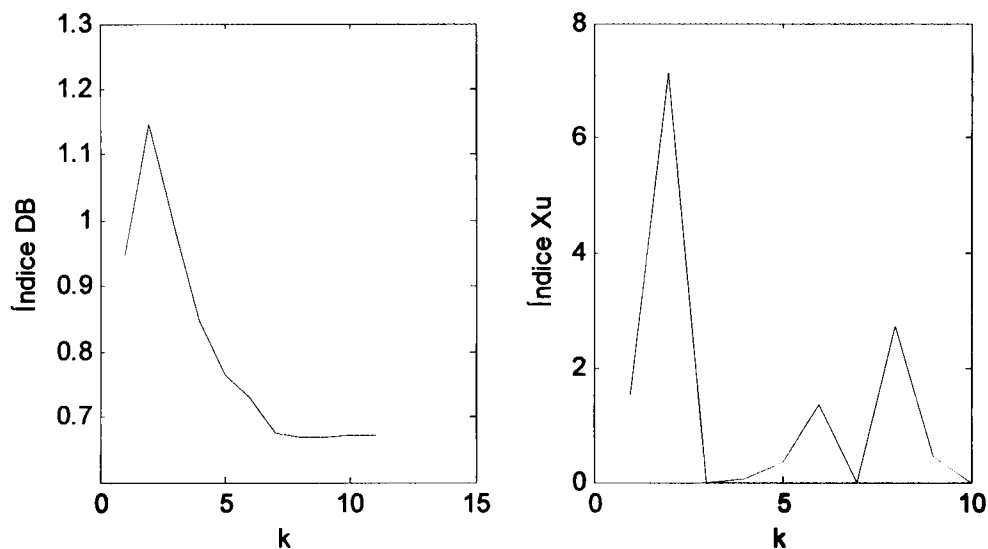


Figura 3.35 Representação dos valores dos índices DB e Xu, para o conjunto II.

Os resultados obtidos pela aplicação dos índices às partições determinadas, revelam pouca consistência, o número natural de classes nos dados é 8 pelo índice DB e 3 pelo índice Xu (figura 3.35).

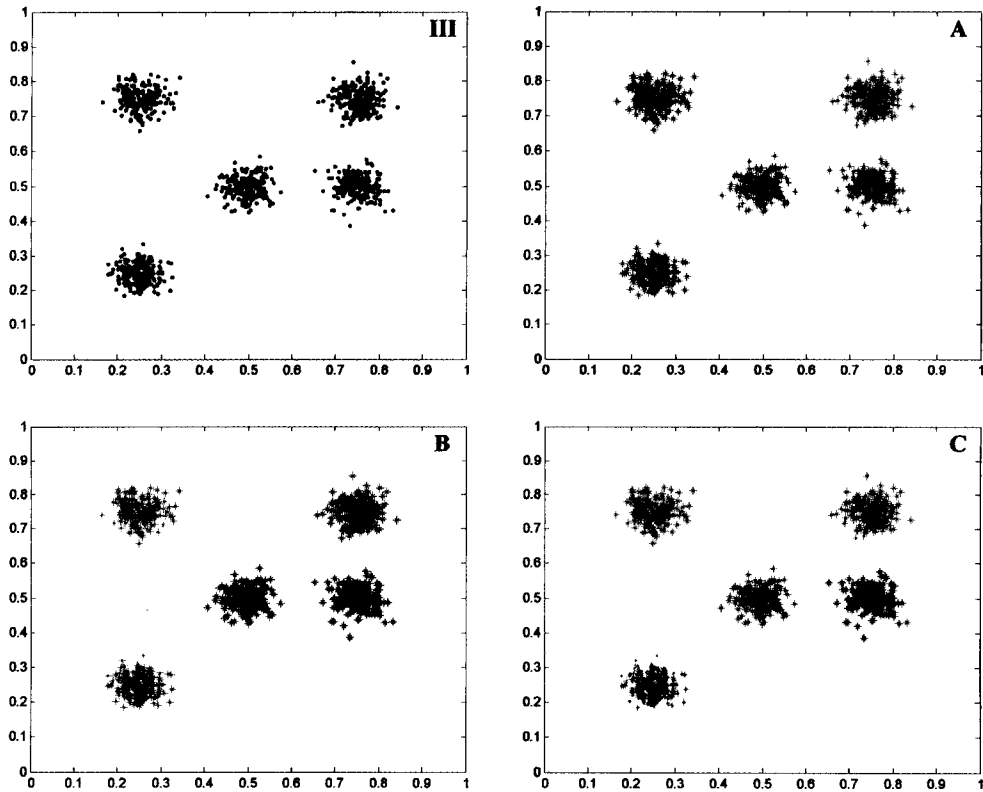


Figura 3.36 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.

Os resultados, para o conjunto III, apresentados nas figuras 3.36 e 3.37, indicam que o método *weighted* tem um comportamento aceitável, analogamente aos métodos anteriormente expostos.

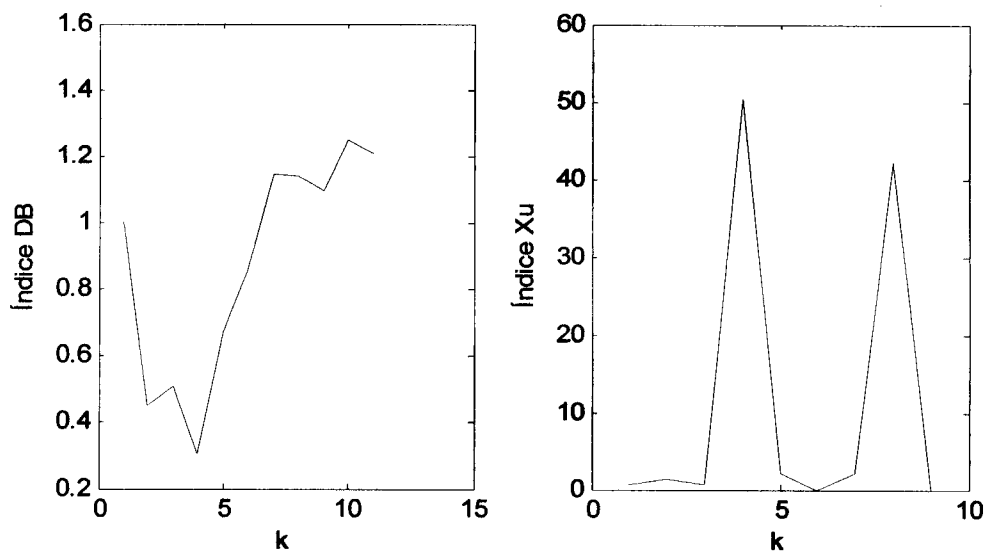


Figura 3.37 Representação dos valores dos índices DB e Xu, para o conjunto III

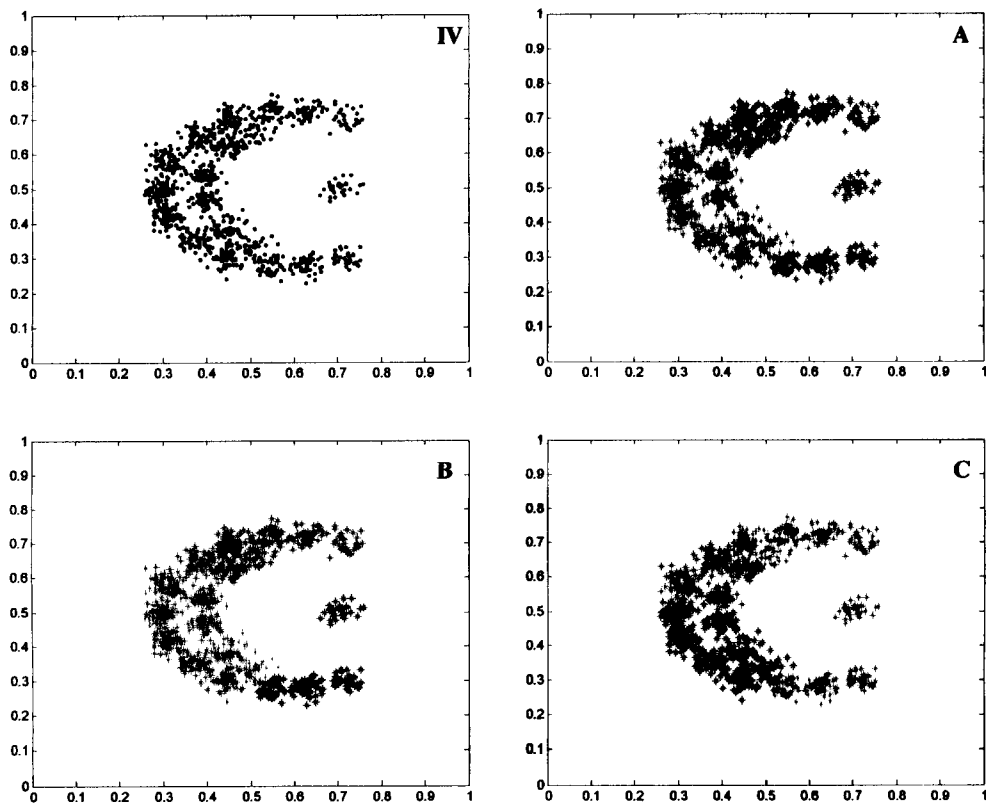


Figura 3.38 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.

A classificação do conjunto de dados IV (figura 3.38) pelo método de agregação *weighted* é análoga à obtida pelo método de agregação *complete*, não produzindo resultados satisfatórios.

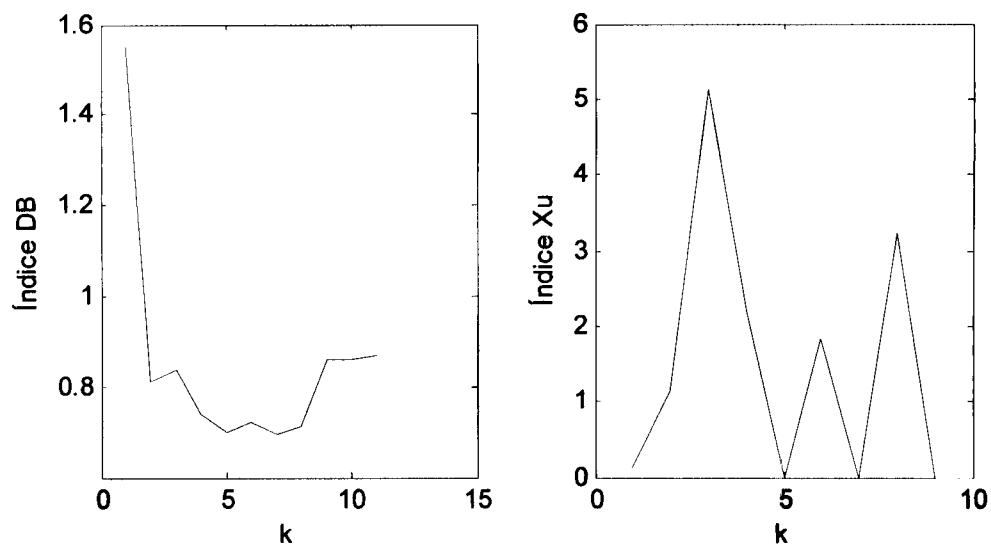


Figura 3.39 Representação dos valores dos índices DB e Xu, para o conjunto IV.

O índice DB considera que o número de classes naturais nos dados será 6 ou 8 enquanto que para o índice Xu a partição mais eficiente será em 4 classes (figura 3.39).

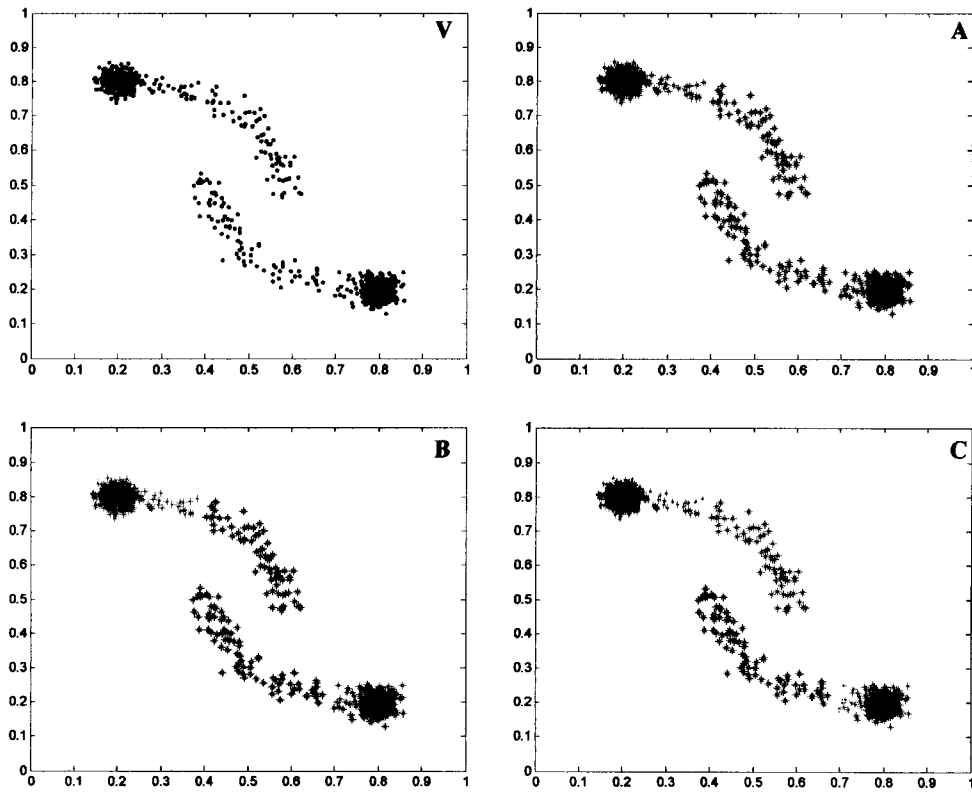


Figura 3.40 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.

Os resultados para o conjunto de dados V, apresentados na figura 3.40, indicam que o método *weighted* tem um comportamento análogo ao dos métodos de agregação *complete* e *average*, falhando na identificação das classes naturais do conjunto de dados da imagem original.

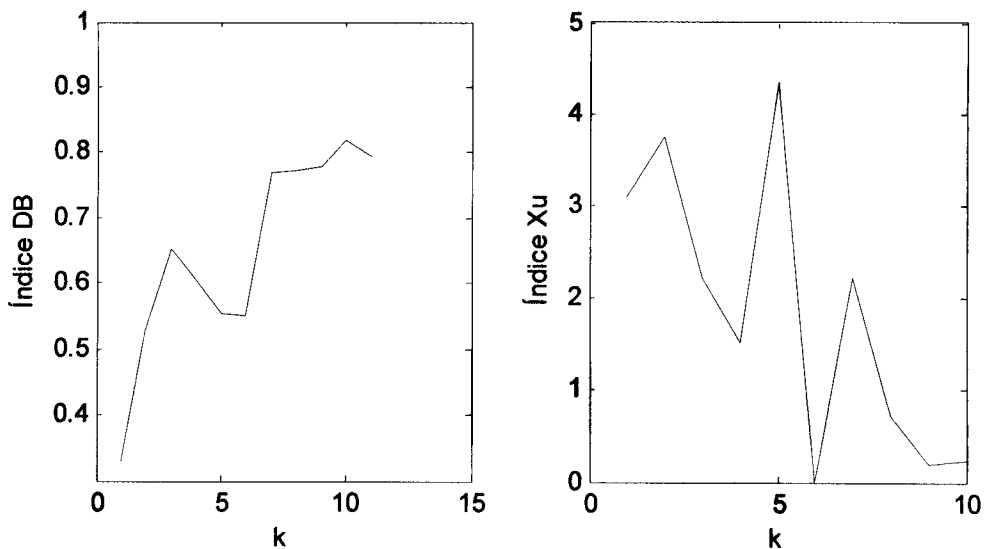


Figura 3.41 Representação dos valores dos índices DB e Xu, para o conjunto V.

Em relação aos índices de semelhança, o número de classes apontado pelo índice DB é 2 enquanto que o índice Xu elege 6 como o número natural de classes existente nos dados (figura 3.41).

3.2.5. Agregação pelo Método *Ward*

O método de Ward utiliza a noção de soma dos quadrados interna de cada *cluster*, ou seja, a soma dos quadrados das distâncias entre todos os objectos no *cluster* e o seu centróide.

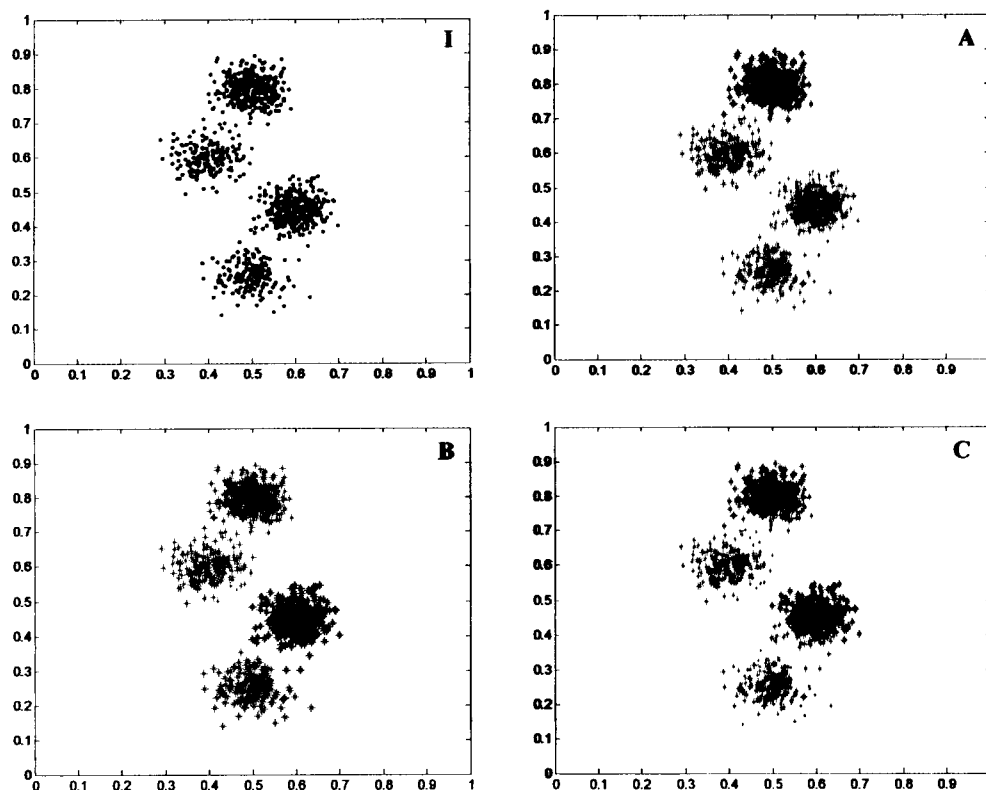


Figura 3.42 A Conjunto de teste (I), classificação em 3 (A), 4 (B) e 5 (C) classes.

Tal como nos métodos de agregação *complete*, *average* e *weighted* o conjunto de dados I é classificado em quatro classes (B), figura 3.42, de forma a apresentar as quatro classes naturais presentes nos dados.

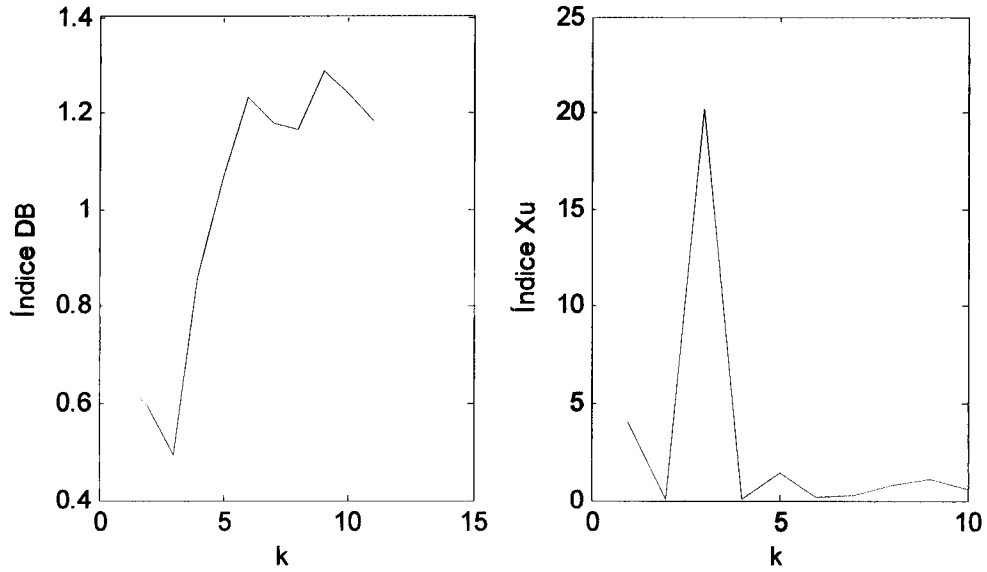


Figura 3.43 Representação dos valores dos índices DB e Xu, para o conjunto I.

Em relação à aplicação dos índices, ambos comprovam a presença de 4 classes nos dados (figura 3.43), o que está de acordo com o que se esperava.

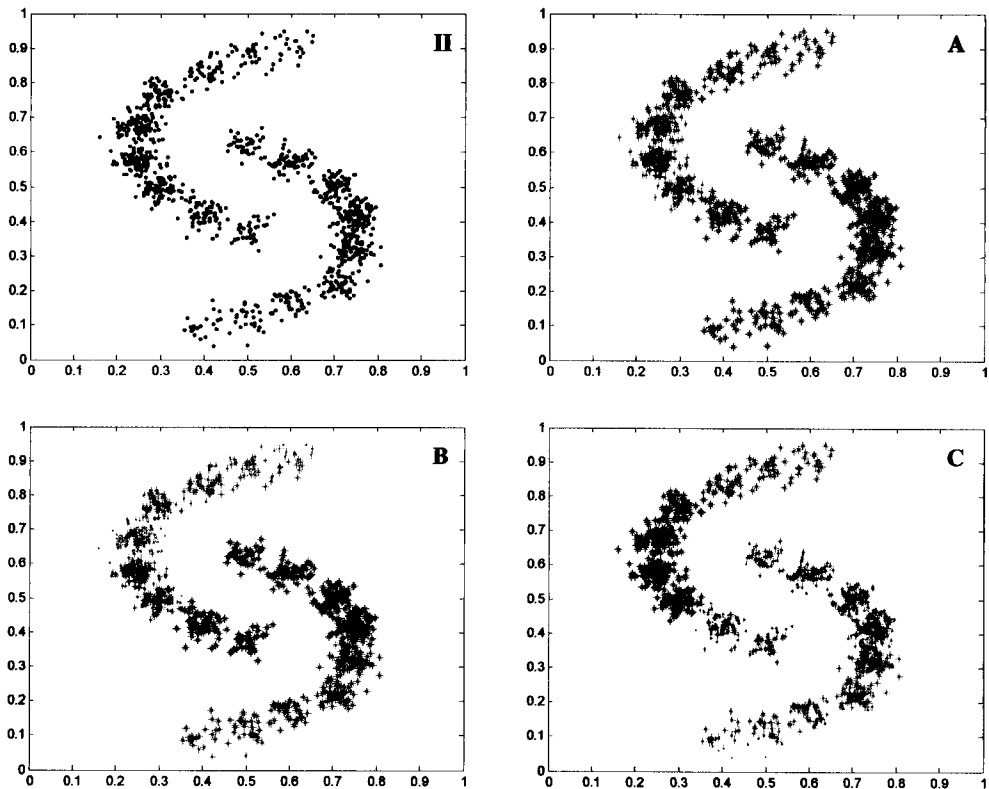


Figura 3.44 Conjunto de teste (II), classificação em 2 (A), 3 (B) e 4 (C) classes.

A classificação do conjunto de dados II (figura 3.44) pelo método de agregação *ward* é análoga à obtida pelo método de agregação *weighted*, não produzindo resultados satisfatórios.

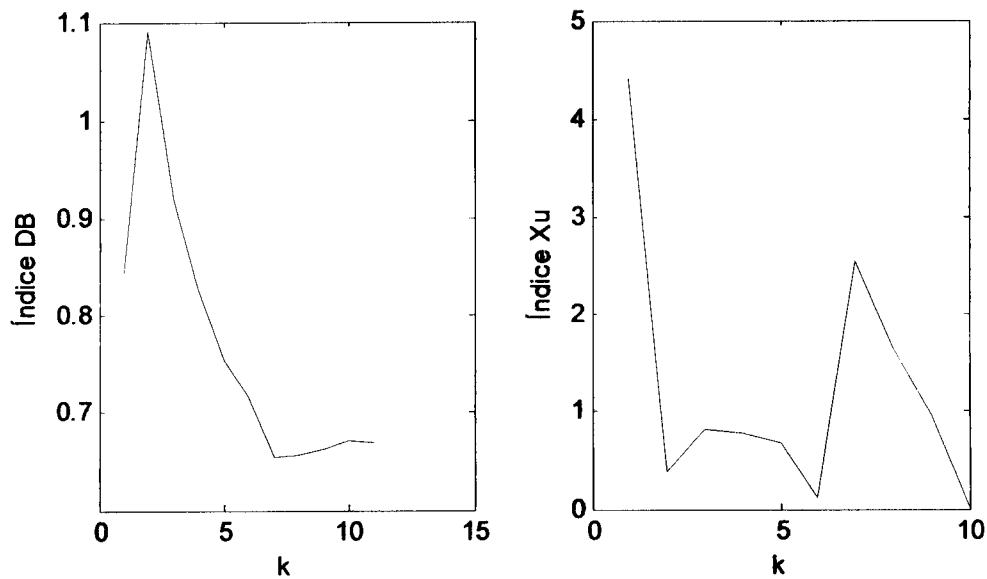


Figura 3.45 Representação dos valores dos índices DB e Xu, para o conjunto II.

Os valores representados pelos índices revelam algumas diferenças (figura 3.45) sendo que o índice DB assinala a presença de 8 classes enquanto que pelo índice Xu o número de classes nos dados deverá ser 2, embora haja também um máximo local para 8 classes.

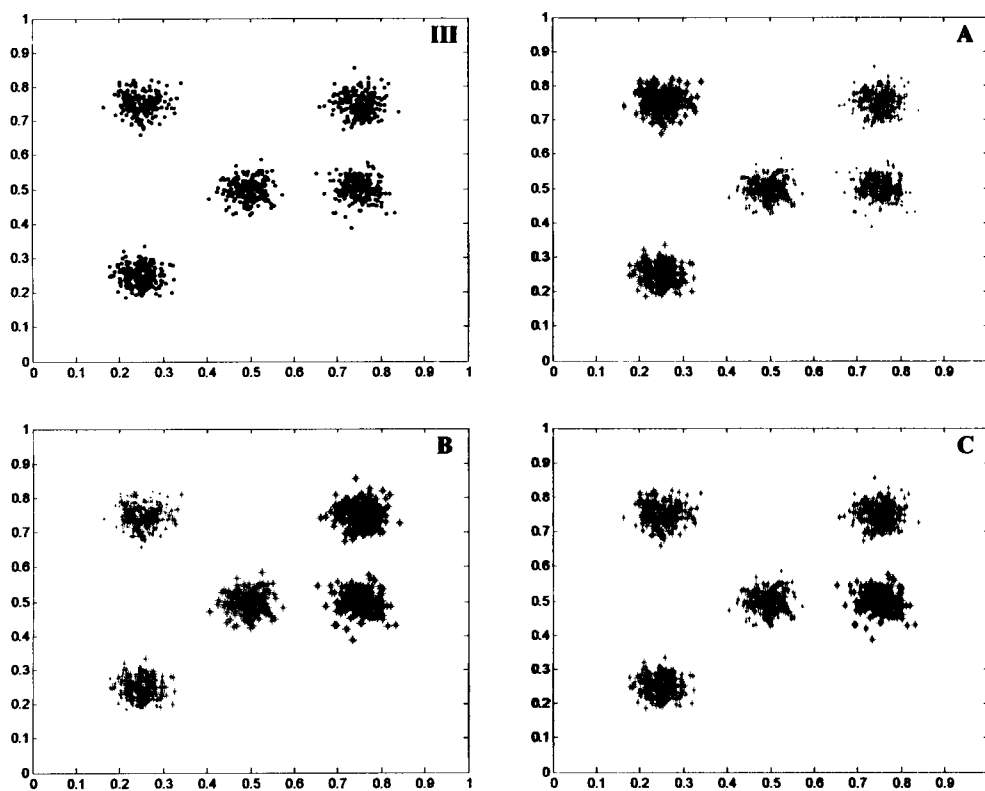


Figura 3.46 Conjunto de teste (III), classificação em 3 (A), 4 (B) e 5 (C) classes.

Os resultados para o conjunto de teste III, apresentados nas figuras 3.46 e 3.47, indicam que o método de Ward tem um bom desempenho, analogamente aos métodos anteriormente expostos, identificando da forma esperada a partição dos dados em 4 classes.

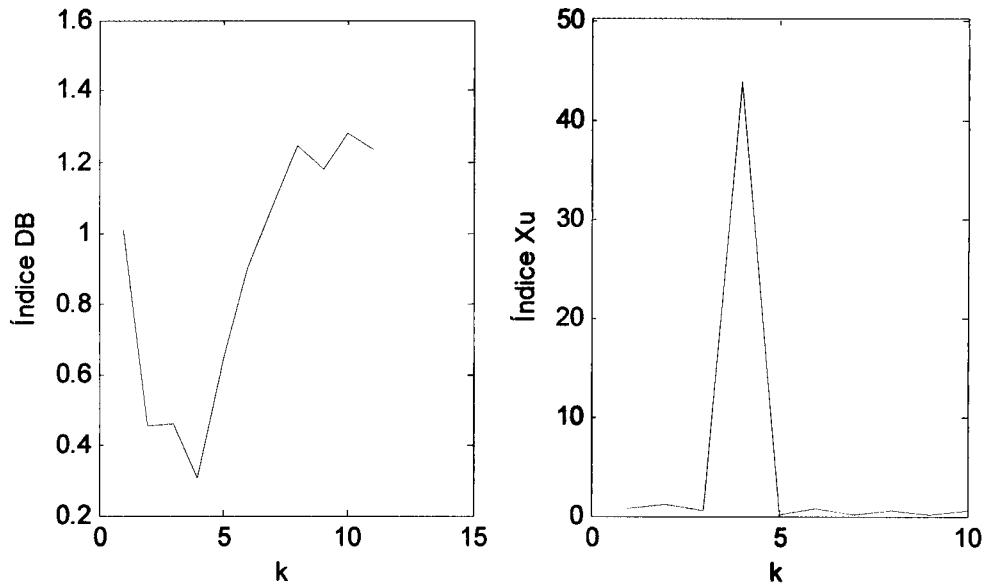


Figura 3.47 Representação dos valores dos índices DB e Xu, para o conjunto III.

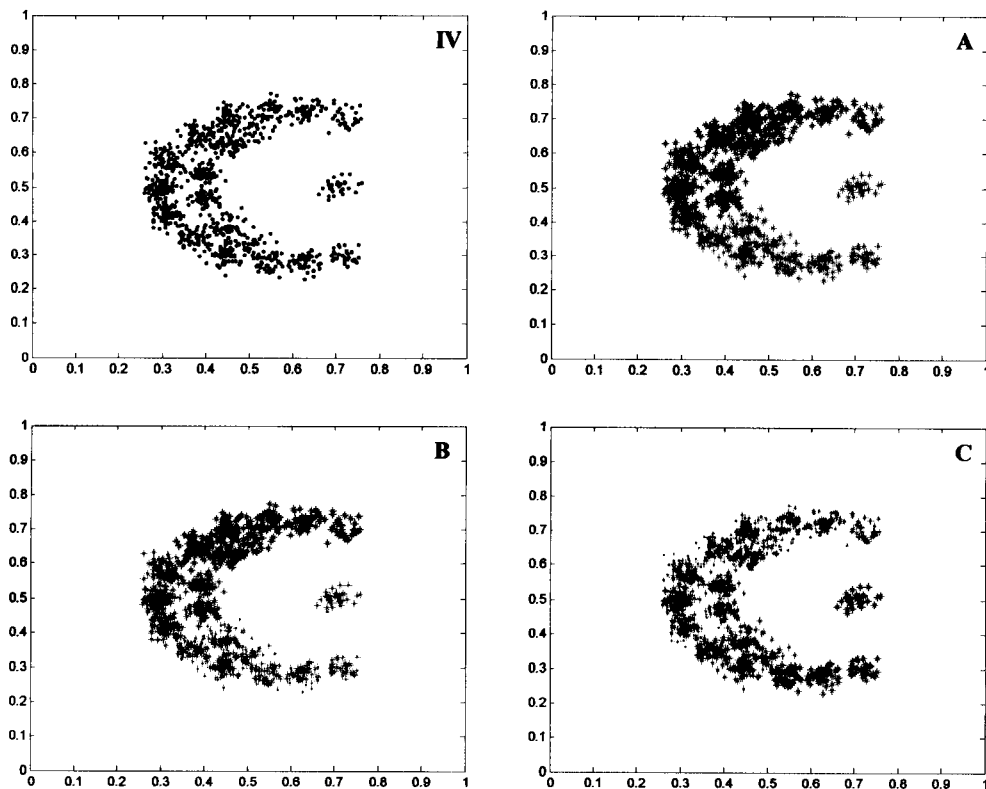


Figura 3.48 Conjunto de teste (IV), classificação em 2 (A), 3 (B) e 4 (C) classes.

A classificação do conjunto de dados IV (figura 3.38) pelo método de agregação *ward* não produziu os resultados desejados, à semelhança do que ocorreu com o método de agregação *average*.

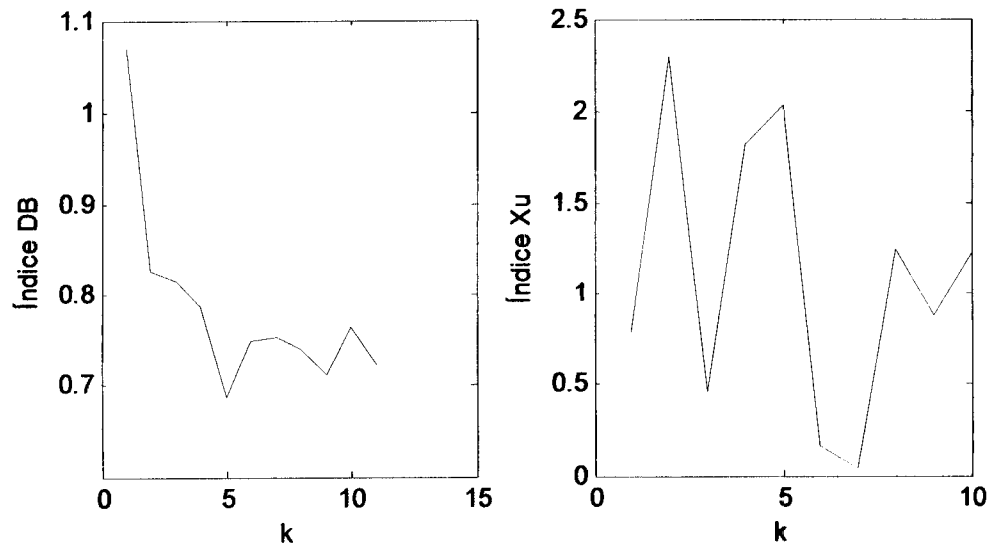


Figura 3.49 Representação dos valores dos índices DB e Xu, para o conjunto IV.

Em relação à eleição da melhor partição, o índice DB considera 6 classes nos dados e o índice Xu aponta a existência de 3 classes naturais (figura 3.49), surgindo eventualmente também a escolha de 6 classes.

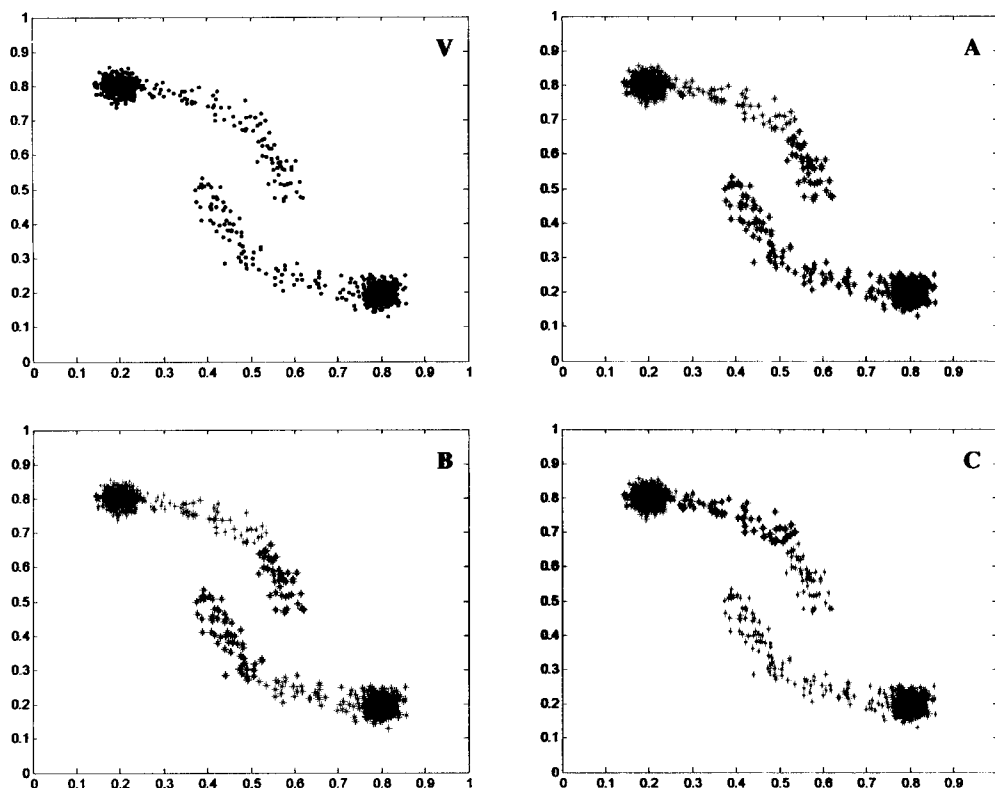


Figura 3.50 Conjunto de teste (V), classificação em 2 (A), 3 (B) e 4 (C) classes.

Os resultados para o conjunto de teste V são apresentados na figura 3.50. Estes resultados indicam que o método de Ward tem um comportamento análogo ao dos métodos de agregação *complete*, *average* e *weighted*, falhando na identificação das classes esperadas do conjunto de dados.

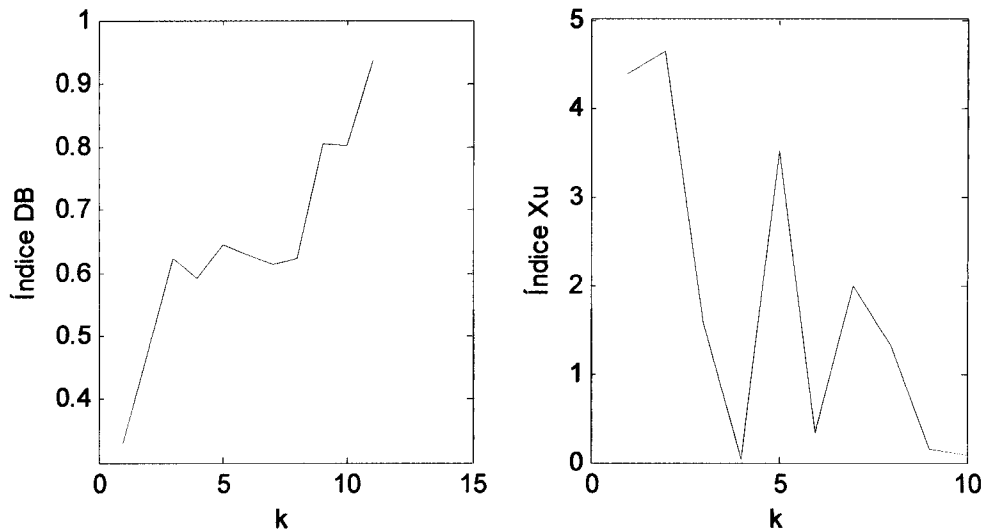


Figura 3.51 Representação dos valores dos índices DB e Xu, para o conjunto V.

Pelo índice DB o número de classes naturais nos dados é 2 enquanto que pelo índice Xu a partição mais eficiente é a que considera 3 classes (figura 3.51) ou, numa 2ª escolha, 6 classes.

3.2.6. Agregação pelos Métodos *Median* e *Centroid*

Em relação aos métodos de agregação de classes *median* e *centroid* foi obtida uma árvore hierárquica não monotónica, o que, segundo o manual da implementação do algoritmo (The MathWorks, 2004), indica que os dois métodos em causa não serão adequados para a classificação do conjunto de dados em análise.

Uma árvore não monotónica é criada quando a medida da distância da união de dois *clusters* a um terceiro *cluster* é menor do que a distância de cada um desses dois *clusters*, ao terceiro *cluster*. Considere-se o seguinte exemplo onde estão representados quatro *clusters* X, Y, Z, W e os respectivos centróides: C_x , C_y , C_z , C_w , sendo W o *cluster* obtido da união dos dois *clusters* X e Y. Ora, da figura constata-se que:

$$d(C_w, C_z) < d(C_x, C_z), d(C_w, C_z) < d(C_y, C_z)$$

e portanto a árvore hierárquica que representa a agregação exposta será uma árvore não monotónica.

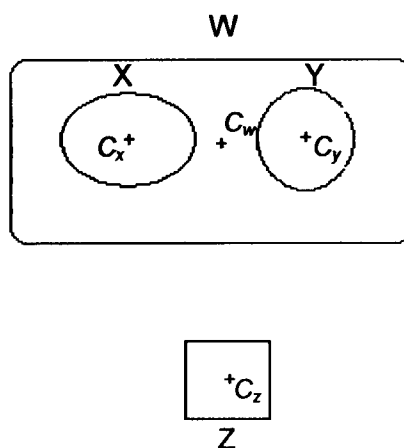


Figura 3.52 Representação de *clusters* e respectivos centróides.

3.3. Conclusões

Numa perspectiva global, a classificação que utiliza o método *single* para a agregação hierárquica, apresenta um desempenho aceitável apenas para os *clusters* de aspecto alongado. Nos conjuntos de dados II, IV e V, este método é o único a produzir resultados aceitáveis, de acordo com a forma como se pretendia efectuar a partição dos dados. No entanto, no caso de conjuntos com *clusters* de forma arredondada, este classificador produz resultados pouco satisfatórios. Os restantes métodos, apesar de apresentarem uma dificuldade notória em reconhecer *clusters* alongados, classificam correctamente grupos de dados com forma arredondada, mostrando um desempenho aceitável nos *clusters* do conjunto de dados I.

Os resultados obtidos para os diferentes métodos vêm reforçar a ideia de não ser possível encontrar um único método capaz de responder correctamente à grande diversidade de formas dos *clusters* possíveis de figurarem em conjuntos de dados reais, como por exemplo dados obtidos a partir de imagens multi-espectrais.

Capítulo 4

Aplicação a imagens multi-espectrais

No caso em que os dados a classificar são uma imagem multi-espectral, cada padrão a classificar corresponde normalmente a um pixel e a dimensão dos dados é em geral o número de bandas da imagem (por exemplo, dimensão 3 para uma imagem RGB).

Ao pretender encontrar a melhor partição da enorme quantidade de padrões (pixels) presentes numa imagem multi-espectral, recorrendo às ferramentas fornecidas pelas *toolboxes* de Processamento de Imagem e de Estatística do MATLAB, o utilizador depara-se com um problema: a impossibilidade em classificar imagens com grande número de dados em tempo útil. A complexidade temporal de todas as técnicas de classificação apresentadas neste trabalho é tanto maior quanto maior o número e a dimensão dos padrões. Por exemplo, numa imagem de 2000×2000 pixels, o número de padrões ascende a 4 milhões.

Neste capítulo é proposto um método de pré-processamento, a aplicar à imagem multi-espectral que se pretende analisar, de forma a ser possível a sua classificação utilizando os métodos anteriormente descritos.

4.1. Método Proposto

O procedimento desenvolvido compreende a diminuição do número de padrões da imagem a classificar, através de um processo de agregação espacial. Em seguida é realizada a partição desse conjunto de dados em vários níveis através de um dos métodos de classificação e é determinada a ou as partições ideais. Por último é efectuada a transposição dos dados classificados para o formato inicial da imagem.

A figura 4.1 apresenta um esquema do processamento efectuado. Na fase 1 são aplicados filtros e é efectuada a segmentação da imagem obtida, a fase 2 consiste na extracção dos vectores para cada padrão correspondente a um objecto segmentado, a fase 3 é a classificação dos dados e a aplicação dos índices de semelhança e a fase 4 é a transposição dos resultados para o formato inicial da imagem.

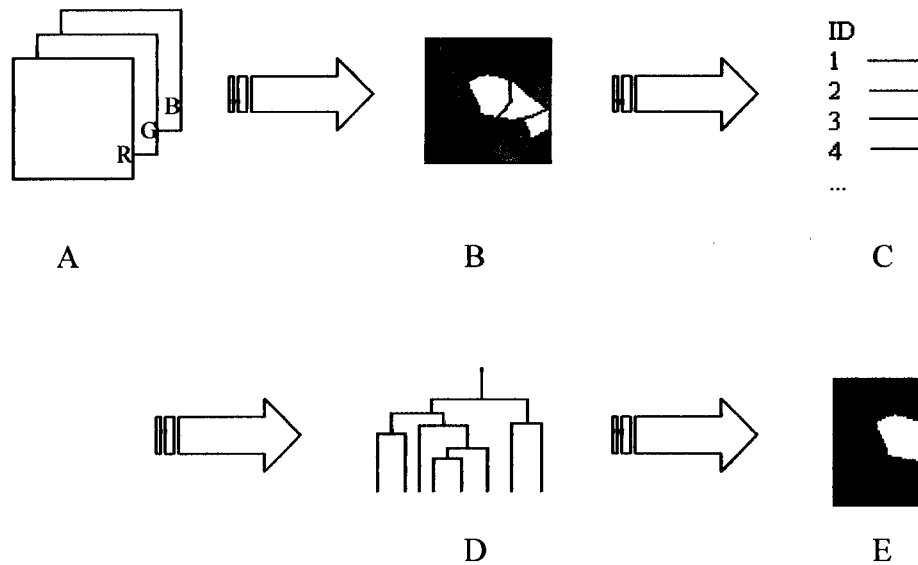


Figura 4.1 Esquema do método proposto para a classificação de imagens multi-espectrais: Imagem Multi-espectral (A), Imagem Segmentada (B), Conjunto de Padrões Reduzido (C), Classificação (D) e Imagem Classificada (E).

4.2. Filtragem

O processo de segmentação de imagem consiste em dividir uma imagem num conjunto de secções (objectos), normalmente com alguma coerência espacial. A reunião de todos os objectos produzidos pelo processo de segmentação deve corresponder à imagem original e a intersecção de quaisquer dois objectos deverá ser o conjunto vazio. Este processo é de certa forma análogo ao processo de classificação de imagem, na medida em que se pretende obter uma partição dos dados, mas difere na relação espacial das regiões, pois duas regiões de intensidade semelhante mas separadas espacialmente são consideradas objectos distintos no processo de segmentação.

A aplicação directa de uma operação de segmentação a uma imagem multi-espectral provoca na maioria dos casos uma sobre-segmentação, ou seja, a produção de demasiados objectos de pequena dimensão, originada por diversos factores entre os quais a existência de ruído na imagem. No sentido de evitar situações deste tipo, o método proposto inicia-se com a aplicação de um filtro passa-baixo à imagem em estudo.

A primeira etapa do método proposto consiste na aplicação de um filtro passa-baixo à imagem de intensidade obtida pela média aritmética das bandas R, G e B da imagem original. Neste caso a imagem será submetida a um filtro espacial linear, ou seja, uma técnica que opera directamente na imagem através de transformações lineares. A aplicação de um filtro passa-baixo a uma imagem (amaciamento ou suavização) provoca uma redução do ruído e remove pequenos detalhes (numa imagem de grandes dimensões), mecanismo que se torna útil quando o objectivo é a extracção de

objectos de tamanho significativo com características espectrais semelhantes (Gonzalez e Woods, 2002).

A aplicação de um filtro espacial pode ser vista como o deslizamento de uma sub-imagem, denominada filtro ou máscara, de pixel em pixel percorrendo todos os pixels da imagem a filtrar, representada por $f(x,y)$ (figura 4.2).

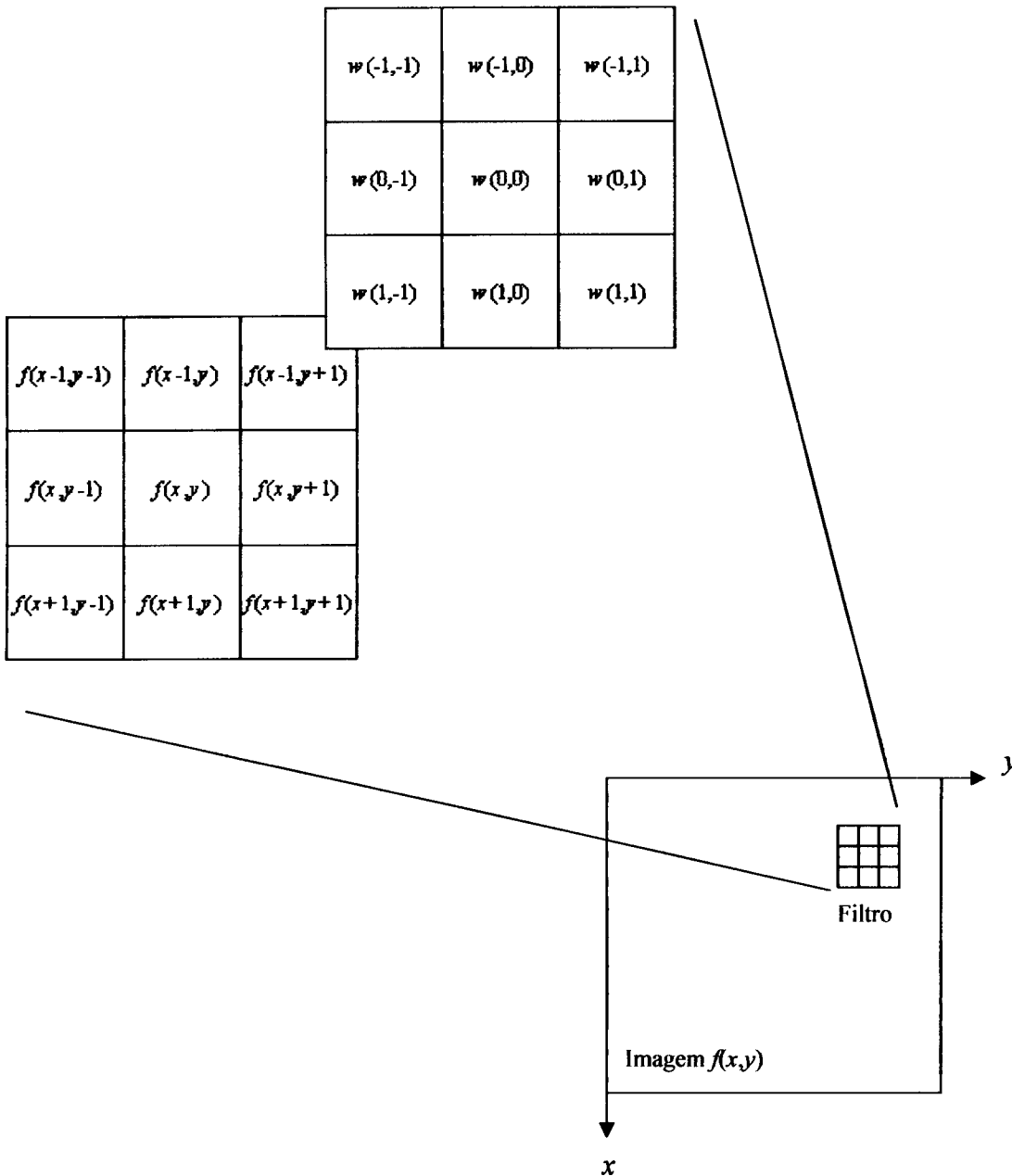


Figura 4.2 O mecanismo de filtragem espacial.

Aos valores que formam um filtro, é atribuída a designação de coeficientes. Para cada ponto (x,y) por onde a máscara passa, o filtro determina um novo valor, R , calculando a soma dos produtos entre os

coeficientes da máscara e a intensidade (que toma valores discretos entre 0 e 255) dos pixels da imagem correspondentes à área abrangida pelo filtro. A equação (10) representa o processo de filtragem utilizando um filtro 3×3 :

$$R = w(-1, -1)f(x-1, y-1) + w(-1, 0)f(x-1, y) + \dots \\ \dots + w(0, 0)f(x, y) + \dots + w(1, 0)f(x+1, y) + w(1, 1)f(x+1, y+1) \quad (10)$$

Para uma máscara de dimensão $m \times n$, assume-se que $m = 2a + 1$ e $n = 2b + 1$, onde $a, b \in \mathbb{N}$; o filtro mais pequeno terá dimensão 3×3 .

Considerando uma imagem $M \times N$ com um filtro de dimensão $m \times n$, o processo de filtragem linear é dado por (11).

$$g(x, y) = \sum_{s=-a}^a \sum_{t=-b}^b w(s, t)f(x+s, y+t), \quad (11)$$

onde $a = (m-1)/2$, $b = (n-1)/2$ e $x, y \in \mathbb{N}_0$ são tais que $a \leq x \leq M-1-a$ e $b \leq y \leq N-1-b$.

Convém notar que, para pontos que se encontrem nos extremos laterais da imagem, a operação de filtragem não pode ser aplicada pelo método descrito. Nesses casos é usual considerar-se uma margem, de acordo com a dimensão do filtro, sendo atribuídos valores a esses pixels através de um critério predefinido.

Para aplicar o filtro à imagem de intensidade (figura 4.7 (B)) utilizou-se a ferramenta `imfilter` disponível na `toolbox` de Processamento de Imagem do MATLAB. Para tal foi produzido o filtro `disk` de dimensão 9×9 pela função `fspecial`. Este é um filtro de média, os seus coeficientes estão distribuídos de forma a darem mais peso ao pixel central e gradualmente menos peso aos restantes pixels à medida que estes se encontram mais distantes do pixel central, circular, de raio 4, contido numa matriz 9×9 . A função `imfilter`, permite a escolha de diversos parâmetros entre os quais opções para o tratamento dos pixels das extremidades da imagem que não são abrangidos pelo filtro. Neste caso foi escolhida a opção `replicate`, a qual associa aos extremos o valor do pixel interno mais próximo. Definidos os parâmetros `input` para `imfilter`, procede-se à filtragem da imagem em teste (exemplo na figura 4.7 (B)) obtendo-se uma imagem menos nítida (figura 4.7 (C)).

Obtida a imagem filtrada, segue-se a aplicação do gradiente, processo bastante utilizado no pré-processamento de imagens para a segmentação (Gonzalez et al, 2004). Este processo passa também pela aplicação de filtros espaciais lineares.

Em processamento de imagem, as primeiras derivadas são implementadas utilizando a magnitude do gradiente. Para uma função $f(x, y)$, o gradiente de f no ponto (x, y) é definido por um vector coluna bidimensional (equação (12)):

$$\nabla \mathbf{f} = \begin{bmatrix} G_x \\ G_y \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} \quad (12)$$

A magnitude (usualmente designada apenas por gradiente) deste vector é dada por:

$$\nabla f = \text{mag}(\nabla \mathbf{f}) = [G_x^2 + G_y^2]^{1/2} = \left[\left(\frac{\partial f}{\partial x} \right)^2 + \left(\frac{\partial f}{\partial y} \right)^2 \right]^{1/2} \quad (13)$$

Considere-se uma região arbitrária de 3×3 de uma imagem, representada por:

z_1	z_2	z_3
z_4	z_5	z_6
z_7	z_8	z_9

Figura 4.3 Região 3×3 de uma imagem.

Ou seja, o ponto central z_5 corresponde ao pixel $f(x, y)$, z_8 representa o pixel $f(x+1, y)$, etc. O conjunto formado pelos pontos $z_1, z_2, z_3, z_4, z_6, z_7, z_8$ e z_9 constitui a vizinhança do pixel central z_5 . Para uma imagem digital, a derivada em cada ponto é considerada em termos de diferenças. No caso em análise, adoptam-se os operadores de *Sobel* (Gonzalez et al, 2004) para definir as derivadas parciais:

$$\frac{\partial f}{\partial x} = (z_7 + 2z_8 + z_9) - (z_1 + 2z_2 + z_3) \quad \text{e} \quad \frac{\partial f}{\partial y} = (z_3 + 2z_6 + z_9) - (z_1 + 2z_4 + z_7) \quad (14)$$

A equação (13) toma então a forma:

$$\nabla f = \left[\left((z_7 + 2z_8 + z_9) - (z_1 + 2z_2 + z_3) \right)^2 + \left((z_3 + 2z_6 + z_9) - (z_1 + 2z_4 + z_7) \right)^2 \right]^{1/2} \quad (15)$$

Ora, para a implementação do gradiente por este método, o MATLAB dispõe de duas máscaras, denominadas operadores de *Sobel*, apresentadas na figura 4.4.

1	2	1
0	0	0
-1	-2	-1

A

1	0	-1
2	0	-2
1	0	-1

B

Figura 4.4 Filtros de *Sobel*: realce de arestas horizontais (A) e de arestas verticais (B).

Para determinar o gradiente da imagem em estudo previamente filtrada (g) basta portanto aplicar a seguinte transformação:

$$\nabla f(x, y) = \left[\left(\sum_{s=-1}^1 \sum_{t=-1}^1 A(s, t) g(x+s, y+t) \right)^2 + \left(\sum_{s=-1}^1 \sum_{t=-1}^1 B(s, t) g(x+s, y+t) \right)^2 \right]^{1/2} \quad (16)$$

Os filtros de *Sobel* são filtros passa-alto que, pela disposição dos seus coeficientes, acentuam essencialmente arestas horizontais (A) e verticais (B). Os pixels da imagem gradiente (exemplo na figura 4.7 (D)) tomam o valor zero em áreas de intensidade constante e valores proporcionais ao grau de alteração de intensidade em áreas cujos pixels possuem valores variáveis.

4.3. Segmentação

Neste momento torna-se novamente necessária a aplicação dum filtro passa-baixo, à semelhança da fase inicial, para, mais uma vez, evitar a produção de um número excessivo de objectos, pelo processo de segmentação. Após esta nova suavização, as condições para a segmentação estão reunidas.

O processo de segmentação consiste na divisão de uma imagem em diferentes regiões, cada uma reunindo determinadas características (Fu e Mui, 1981). Neste trabalho, a técnica de segmentação precede a classificação da imagem multi-espectral, sendo esta abordagem utilizada frequentemente na classificação deste tipo de imagens (Soh e Tsatsoulis, 1999). Existem diversas técnicas de segmentação, as quais se podem agrupar em três grandes grupos: *characteristic feature thresholding* ou partição de histograma, detecção de descontinuidades e extracção de regiões. Cada uma destas abordagens apresenta por si só vantagens e inconvenientes. A técnica utilizada neste estudo designa-se por *watershed* e incorpora muitos dos conceitos das três abordagens mencionadas, produzindo assim resultados mais consistentes e estáveis. A palavra *watershed* exprime a ideia de uma linha que divide sistemas de captação de água ou bacia hidrográfica. Para compreender a forma como este conceito se

pode aplicar a processamento de imagem considere-se uma imagem (g) como uma superfície topográfica onde os valores de intensidade da imagem, $g(x,y)$, são interpretados como alturas no terreno. Os locais de menor elevação são as bacias de captação (*catchments basins*), ou seja, imaginando que a imagem possa ser coberta de água, estes são os locais onde a água se acumulará em primeiro lugar. Cada bacia de captação é formada em torno de um mínimo regional que é uma área conexas da imagem tal que todos os pixels que a compõem possuem a mesma intensidade sendo esta inferior à de todos os pixels vizinhos. (Vincent e Soille, 1991). As linhas que delimitam as bacias de captação (componentes conexas da imagem) são linhas a que se atribui a designação de *watersheds*, figura 4.5.

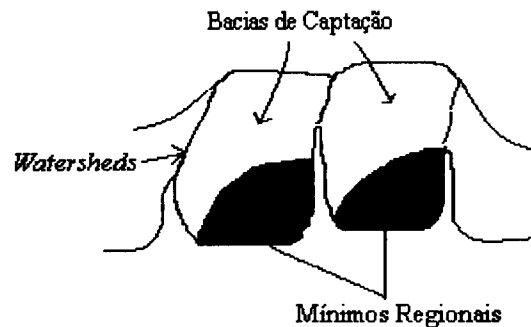


Figura 4.5 Bacias de captação, *watersheds* e mínimos regionais.

À medida que o nível da água vai subindo na imagem, são construídas barreiras ou diques (*dams*) para impedir que as componentes conexas, bacias de captação, sejam agregadas. A técnica de segmentação *watershed* consiste na separação destas bacias de captação definindo as linhas contínuas que as separam.

Sejam M_1, M_2, \dots, M_R os conjuntos dos pontos correspondentes aos mínimos regionais de uma imagem $g(x,y)$ e $C(M_i)$ o conjunto dos pontos (x,y) na bacia de captação associada ao mínimo regional M_i , conjunto esse que forma uma componente conexas. O nível da água, n , vai subindo em passos de uma unidade sendo que o nível mínimo é $n = \min(g(x,y)) + 1$ e o nível máximo será $n = \max(g(x,y)) + 1$. Seja $T[n]$ o conjunto de coordenadas (s,t) para as quais $g(s,t) < n$, ou seja:

$$T[n] = \{(s,t) | g(s,t) < n\}. \quad (17)$$

$T[n]$ é precisamente o conjunto de pontos da imagem que se encontram abaixo do plano $g(x,y) = n$.

Seja $C_n(M_i)$ o conjunto de pontos da bacia de captação associada ao mínimo regional M_i que estão inundadas quando a água está no nível n , ou seja:

$$C_n(M_i) = C(M_i) \cap T[n]. \quad (18)$$

A reunião destes $C_n(M_i)$ designa-se por $C[n]$:

$$C[n] = \bigcup_{i=1}^R C_n(M_i) \quad (19)$$

Considere-se $n_{\max} = \max(g(x,y)) + 1$, então:

$$C[n_{\max}] = \bigcup_{i=1}^R C(M_i) \quad (20)$$

Ora, pelo mecanismo descrito, o número de elementos dos conjuntos $T[n]$ e $C_n(M_i)$ vai aumentando ou permanecendo constante (mas nunca decresce) à medida que o nível da água vai subindo na imagem, sendo que $C[n-1] \subseteq C[n]$. De (17) e (18) decorre que $C[n] \subseteq T[n]$ e portanto $C[n-1] \subseteq T[n]$. Ora esta inclusão permite concluir que cada componente conexa de $C[n-1]$ está contida em exactamente uma componente conexa de $T[n]$.

O algoritmo associado à segmentação *watershed* é inicializado com $C[n_{\min}] = T[n_{\min}]$, onde $n_{\min} = \min(g(x,y)) + 1$, e procede recursivamente assumindo-se, em cada passo n , que $C[n-1]$ está construído. Seja $Q[n]$ o conjunto formado pelas componentes conexas em $T[n]$. Para cada componente conexa $q \in Q[n]$, existem três situações possíveis:

- a) $q \cap C[n-1] = \emptyset$;
- b) $q \cap C[n-1]$ contém exactamente uma componente conexa de $C[n-1]$;
- c) $q \cap C[n-1]$ contém mais do que uma componente conexa de $C[n-1]$.

No momento em que a água está no nível n , se a) ocorrer, então existe uma componente conexa $q \in Q[n]$ que não possui nenhum ponto comum com o conjunto de pontos das bacias de captação formadas até ao nível $n-1$, $C[n-1]$, o que significa o aparecimento de um novo mínimo regional. Neste caso, q é incorporado no conjunto $C[n-1]$ originado $C[n]$. No mesmo nível n , a ocorrência da situação b), indica que existe uma componente $q \in Q[n]$ que faz parte da bacia de captação associada a um determinado mínimo regional já detectado, e portanto q é incorporado no conjunto $C[n-1]$ originado $C[n]$. A situação c) ocorre quando uma linha de separação entre 2 ou mais bacias de captação é encontrada. Para prevenir a união das bacias de captação em causa, provocada pela subida do nível da água, é necessária a construção de um ou mais diques em q .

A construção de um dique, capaz de impedir a junção de duas bacias é feita através da operação morfologia de dilatação. Na figura 4.6 A e B estão representadas porções de dois *catchment basins* nos níveis $n-1$ e n de inundação, respectivamente. A cada *catchment basin* está associado um mínimo regional, sejam M_1 e M_2 os conjuntos de pontos em cada um desses mínimos; $C_{n-1}(M_1)$ e $C_{n-1}(M_2)$ os conjuntos de pontos de cada *catchment basin* no nível $n-1$. Na etapa n , a água passou de uma das componentes conexas para a outra e por isso é necessária a construção de uma barreira entre as duas bacias de captação. Seja $C[n-1]$ a união dos dois conjuntos $C_{n-1}(M_1)$ e $C_{n-1}(M_2)$, a qual é representada por duas componentes conexas na figura 4.6 A. Seja q a componente conexa da figura 4.6 B e considere-se a dilatação dos pontos da figura 4.6 A pelo elemento de estrutura representado na figura 4.6 C, sujeita a: 1) a dilatação deve ser restrita a q , ou seja, o centro do elemento de estrutura durante a dilatação só pode ser colocado em pontos de q ; 2) a dilatação não pode ser efectuada em pontos que

provoquem a união dos conjuntos a dilatar. A figura 4.6 D mostra que a primeira dilatação provocou a expansão uniforme de cada uma das componentes conexas originais pois a condição 1) foi satisfeita sem originar pontos que unissem as duas componentes. Na segunda dilatação os pontos que quebraram a condição imposta em 2) (pontos assinalados com uma cruz na figura 4.6 D) são precisamente os pontos que vão construir o dique pretendido. Esta construção fica completa quando se atribui aos pixels correspondentes ao dique, uma intensidade superior a todas as intensidades da imagem original (em geral atribui-se a intensidade 255 numa imagem 8 bits) para que, em relação à elevação do terreno, estes pontos estejam sempre acima dos pontos das bacias de captação.

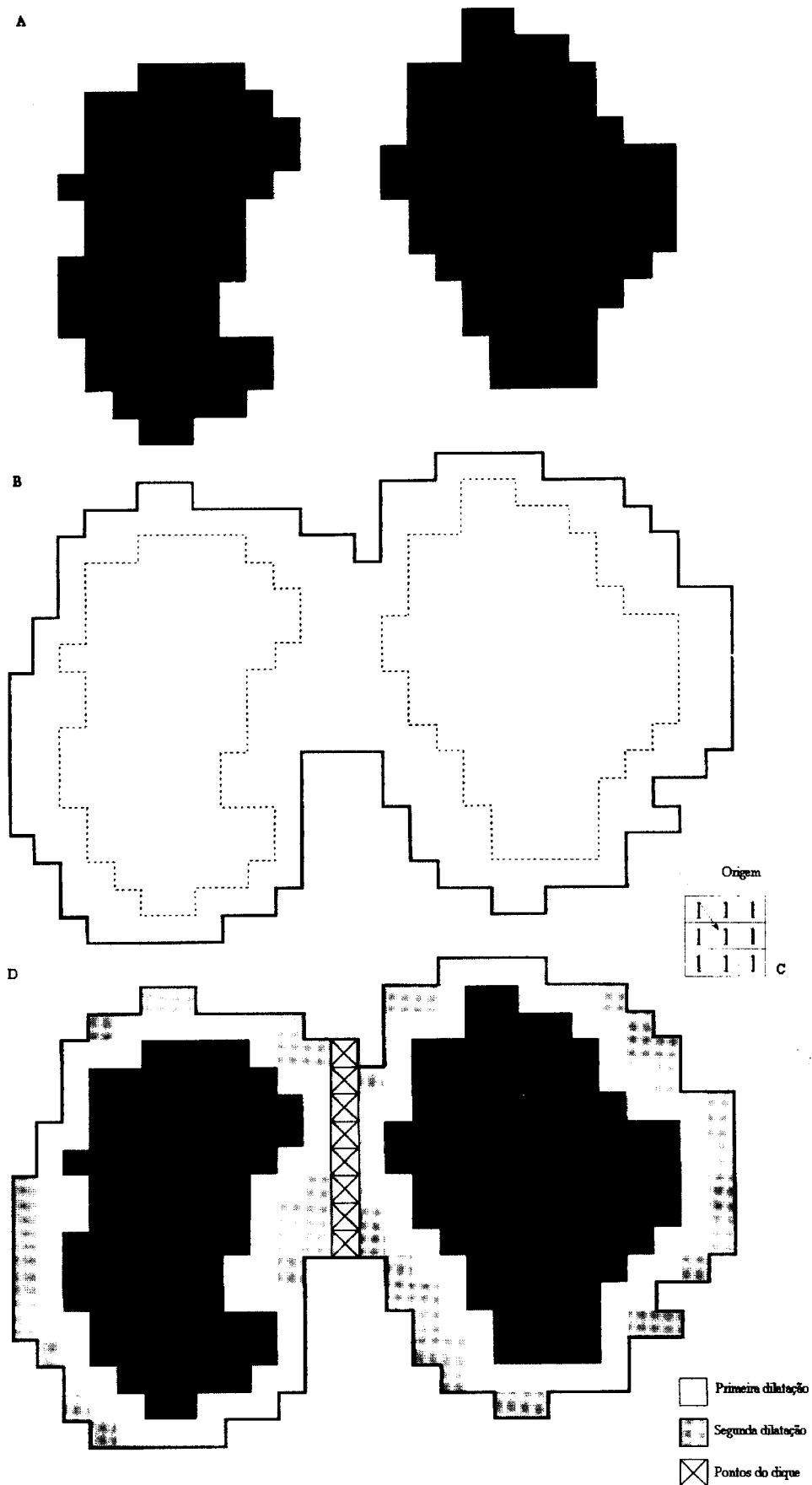


Figura 4.6 A: Duas bacias de captação inundadas na fase $n-1$ da subida do nível da água; B: Fase n da inundação, junção da água das duas bacias; C: Elemento estruturante da dilatação; D: Resultado da dilatação e construção do dique.

A eficiência deste algoritmo pode ser melhorada obrigando n a tomar apenas valores que correspondam a pixels existentes na imagem, valores que podem ser extraídos facilmente utilizando o histograma da imagem. A sua implementação utilizando o MATLAB foi efectuada aplicando a função *watershed* à imagem previamente produzida pelo método descrito. A figura 4.7 (E) ilustra um exemplo de uma imagem segmentada por esta técnica, os objectos (bacias de captação) são as regiões a branco e a preto estão assinaladas as *watersheds*.

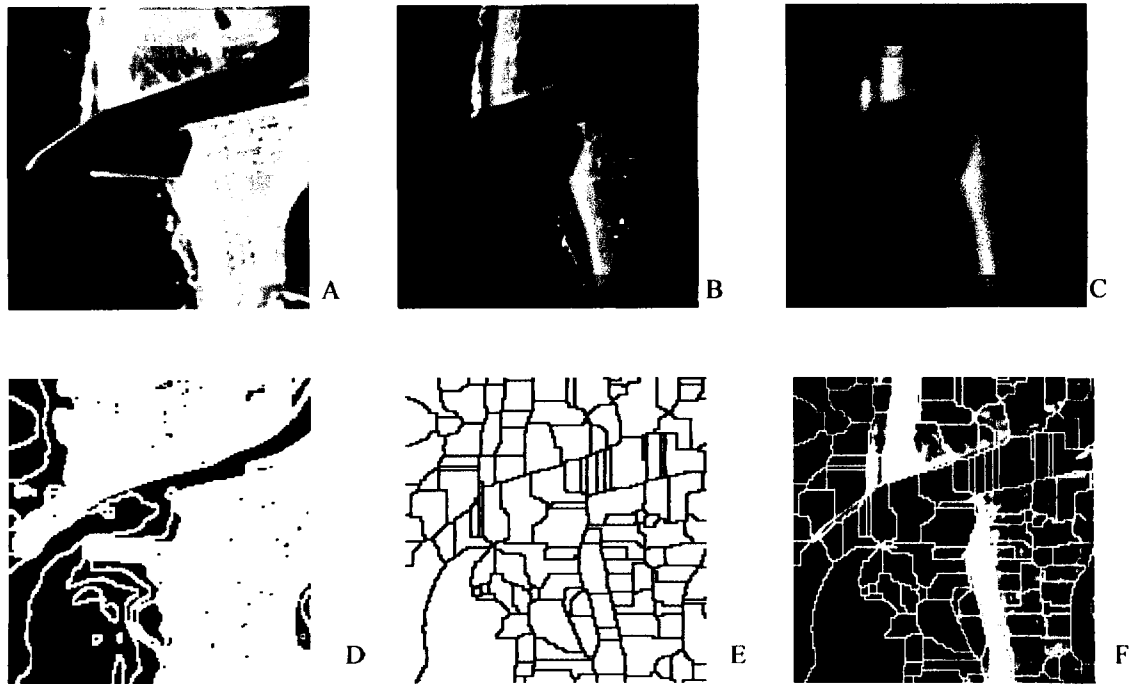


Figura 4.7 Imagem original (A), imagem de intensidade (B), imagem filtrada (C), imagem gradiente (D), imagem segmentada (E) e segmentação sobreposta com imagem original (F).

A imagem obtida pela segmentação *watershed* é constituída por pixels, (x,y) , cujas intensidades, $s(x,y)$, tomam valores entre 0 (atribuído aos pontos das linhas *watershed*) e n_{obj} , onde n_{obj} é o número total de objectos encontrados. Antes do processo de classificação, é necessário atribuir a cada pixel de intensidade nula (*watershed*) a intensidade de um dos objectos vizinhos. Assim, para cada pixel (x,y) , tal que $s(x,y) = 0$, foi atribuída a intensidade que representa a moda do conjunto das intensidades da sua vizinhança, associando aos pixels das vizinhanças horizontais e verticais peso 2 e aos pixels das vizinhanças diagonais peso 1. Convém notar que para a determinação desta moda, os pixels da vizinhança que possuam intensidade nula não vão ser considerados. Um exemplo deste processo de atribuição de objectos aos pixels correspondentes a *watersheds* é apresentado na figura 4.8. Na figura está representada uma secção da imagem classificada onde os pixels z_5 , z_6 , z_8 e z_9 estão sobre a linha *watershed* e por isso têm intensidade 0. Ao pixel z_5 é atribuída a intensidade 12 já que esta representa a moda do conjunto: $\{z_1, z_2, z_2, z_3, z_4, z_4, z_7\} = \{13, 12, 12, 13, 12, 12, 13\}$.

$z_1=13$	$z_2=12$	$z_3=13$
$z_4=12$	$z_5=0$	$z_6=0$
$z_7=13$	$z_8=0$	$z_9=0$

Figura 4.8 Secção 3×3 da imagem resultante da segmentação.

Este procedimento é aplicado a todos os pixels pertencentes às linhas de separação dos objectos e após esta fase o processo de classificação pode ser iniciado.

4.4. Classificação e Transposição dos Resultados

Nesta fase é necessário construir um conjunto de variáveis representativas dos objectos obtidos pelo processo de segmentação, para proceder à sua classificação. Para tal, considerem-se todos os pares de coordenadas dos pixels que pertencem a um determinado objecto da imagem segmentada. De seguida considerem-se as intensidades atribuídas a estes pares de coordenadas na imagem RGB inicial. A ideia da segmentação efectuada é considerar que todos esses pixels formam um conjunto relativamente homogéneo e que por isso pode ser representado como um único elemento no processo de classificação. Assim, a cada elemento representando um objecto, são atribuídas três intensidades, em que cada uma delas corresponde à média das intensidades de todos os pixels atribuídos a esse objecto, nas 3 bandas R, G e B. O resultado deste processo é portanto um conjunto de pontos, cada um representando um objecto (o conjunto tem tantos dados quantos os objectos presentes na imagem segmentada), composto por três componentes, representado numa matriz de dimensão $n_{obj} \times 3$.

Segue-se o processo de classificação onde a matriz $n_{obj} \times 3$ obtida é classificada em 2, 3,... até 40 classes pelos cinco métodos de agregação descritos no Capítulo 2.

A classificação efectuada atribui a cada elemento da matriz, uma determinada classe; ou seja, a cada objecto é atribuída uma classe. Assim, a transposição do resultado da classificação para o formato original da imagem consiste em atribuir a todos os pixels que compõem cada objecto a classe associada a esse mesmo objecto, permitindo desse modo a visualização de uma imagem classificada.

Os resultados obtidos por todo o método descrito de filtragem, segmentação e classificação, podem ser visualizados na imagem classificada de uma forma mais intuitiva se as cores utilizadas para a representação das classes obedecerem a um determinado critério. Para tal, na partição inicial em 40 classes é atribuída uma cor a cada classe. À medida que duas classes são agregadas formando uma nova classe, a cor que é atribuída à nova classe é a cor da classe inicial com mais elementos. Esta atribuição inicial é efectuada através de um método que utiliza de forma alargada a gama de cores

disponíveis, mantendo no entanto a relação de proximidade entre classes através da atribuição de cores visualmente parecidas (Marçal, 2005). Deste modo, da observação de imagens apresentando diferentes partições, é possível perceber que classes foram obtidas da junção (e separação) de duas classes.

Capítulo 5

Teste com Imagem de Satélite

O método de classificação de imagens descrito no Capítulo 4 foi aplicado a uma imagem de satélite de teste. A imagem foi obtida pelo instrumento ASTER (*Advanced Spaceborne Thermal Emission and Reflection Radiometer*) do satélite Terra, lançado pela NASA em 1999 como parte do programa Earth Observing System – EOS (Yamaguchi, 1998), tendo sido seleccionada uma secção de dimensão 800×800 pixels, onde cada pixel corresponde a uma área de 15×15 metros. Esta imagem apresenta 3 bandas nas zonas visível e infravermelho próximo do espectro electromagnético, sendo apresentada uma composição RGB das bandas 1, 2 e 3 na figura 5.1. Apesar desta imagem de teste ser uma pequena secção da imagem original (cerca de 4%), o número de pixels é demasiado grande para se efectuar uma classificação de forma directa usando os classificadores de aglomeração hierárquica do MATLAB.

A imagem em análise (figura 5.1) representa a zona da Ria de Aveiro que compreende uma zona de água profunda (o mar e a ria), zonas de água pouco profundas na ria, zonas de terra com diferentes tipos de vegetação, zonas urbanas ou com vias de comunicação e área de praia arenosa.



Figura 5.1 Imagem de teste (composição RGB das bandas 1, 2 e 3, com histograma modificado).

Do processo de segmentação resultaram 4105 objectos, com dimensão média de 155 pixels. Na figura 5.2 estão representados os pontos que formam o conjunto de dados a classificar, em que cada ponto representa um objecto. Cada um dos eixos da representação gráfica corresponde à intensidade em cada uma das 3 bandas do ASTER, podendo tomar valores entre 0 e 255 (dados a 8 bits).

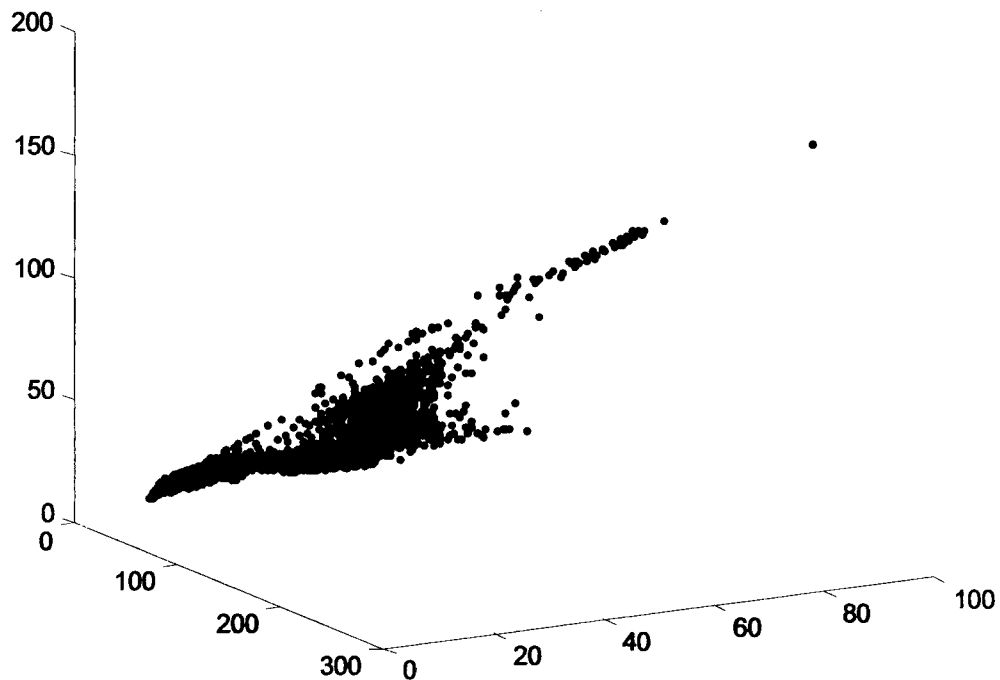


Figura 5.2 Representação a 3 dimensões dos pontos que representam cada objecto obtido da segmentação.

O método de classificação, descrito no Capítulo 2, foi aplicado a este conjunto de dados, sendo os resultados transpostos para a imagem de acordo com o método descrito no Capítulo 4. Nas secções seguintes são apresentados os resultados obtidos para os diferentes métodos de agregação hierárquica assim como os gráficos resultantes da aplicação dos índices DB e Xu para cada partição obtida.

5.1. Método de Agregação *Single*

O primeiro método de classificação testado de acordo com os procedimentos descritos nos capítulos 2 e 4 utilizou a agregação hierárquica com o parâmetro *single*. O resultado da classificação desta imagem em 40 classes é apresentado na figura 5.3, onde se usou uma tabela de cor para melhorar a visualização dos dados.

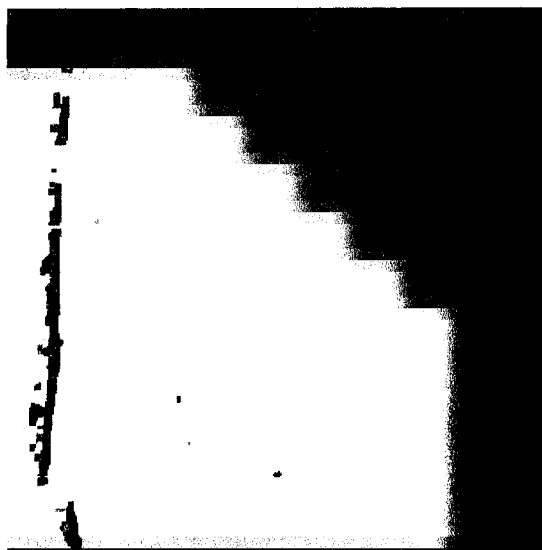


Figura 5.3 Imagem classificada em 40 classes pelo método de agregação *single*.

O método de agregação *single* produziu resultados pouco satisfatórios tendo formado uma classe com a maioria dos pixels da imagem agregando zonas de água, zonas verdes e urbanas. A maioria das restantes classes foi atribuída a zonas de areia. Face a estes resultados, a representação dos valores dos índices DB e Xu para este método não foi considerada, uma vez que a classificação inicial em 40 classes é claramente insatisfatória.

5.2. Método de Agregação *Complete*

O segundo método de agregação testado foi o *complete*. Os resultados da respectiva classificação em 40 classes são apresentados na figura 5.4.



Figura 5.4 Imagem classificada em 40 classes pelo método de agregação *complete*.

Uma análise visual da imagem RGB original e desta, permite verificar que foram identificadas de forma bastante satisfatória as zonas de água (distingue zonas de maior e menor profundidade), as zonas de areia e verdes; as zonas urbanas formam áreas bastante heterogéneas, característica captada de forma razoável por este método.

As figuras 5.5 e 5.6 apresentam os índices de DB e Xu para o método *complete* para o número de classes entre 2 e 40 para o índice DB e entre 2 e 39 classes para o índice Xu. O valor apresentado no eixo das abcissas, k , corresponde à partição em $k+1$ classes.

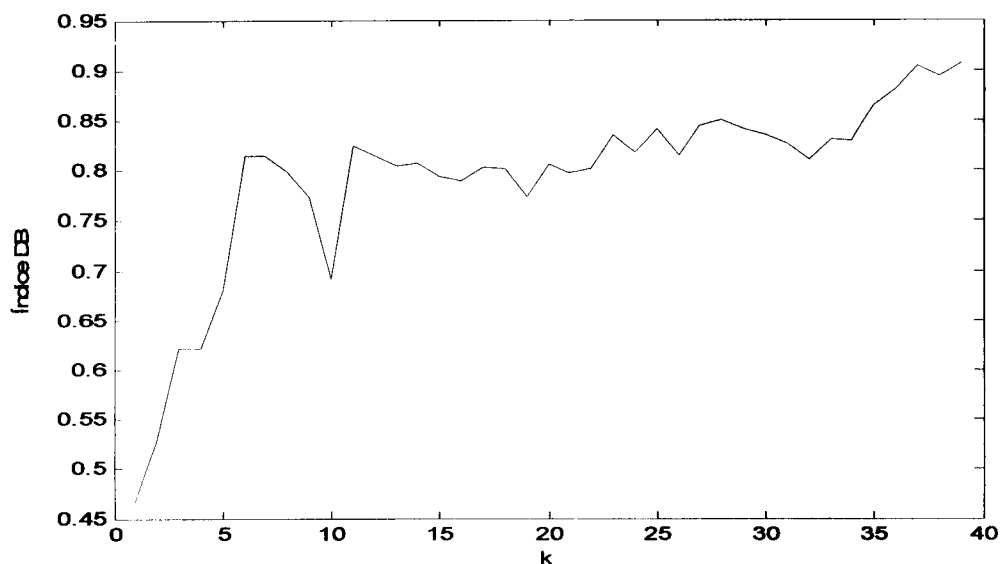


Figura 5.5 Representação dos valores do índice DB para a classificação pelo método *complete*.

Pelos valores do índice DB para as diferentes partições, representados na figura 5.5, conclui-se que as melhores partições serão as que dividem os dados em apenas 2 ou em 11 classes.

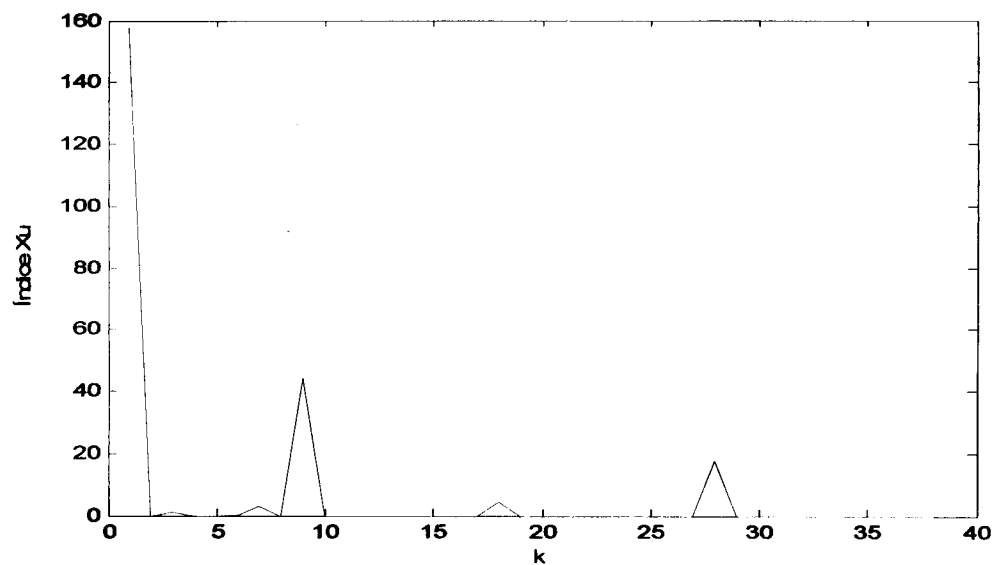


Figura 5.6 Representação dos valores do índice Xu para a classificação pelo método *complete*.

O índice Xu elege para melhores partições aquelas que dividem os dados em 2, 10 e 29 classes. As imagens resultantes da classificação em 2, 10, 11 e 29 classes através do método de agregação *complete* são apresentadas na figura 5.7.

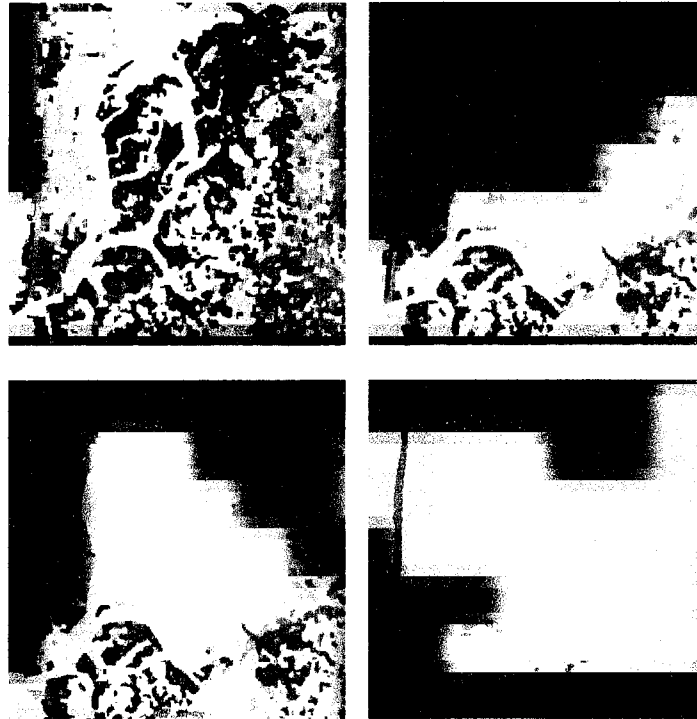


Figura 5.7 Imagem classificada em 29, 11, 10 e 2 classes pelo método de agregação *complete*.

5.3. Método de Agregação *Average*

Outro método de agregação testado, foi o *average*. A classificação da imagem da figura 5.1 em 40 classes, utilizando o método de agregação *average*, produziu o resultado que se pode visualizar na figura 5.8.



Figura 5.8 Imagem classificada em 40 classes pelo método de agregação *average*.

A análise visual deste resultado e da imagem RGB permite concluir que a partição em 40 classes identifica de forma razoável as zonas de água mais profunda, as zonas da ria com menos profundidade, as zonas arenosas, as zonas verdes e as regiões urbanas mais densas.

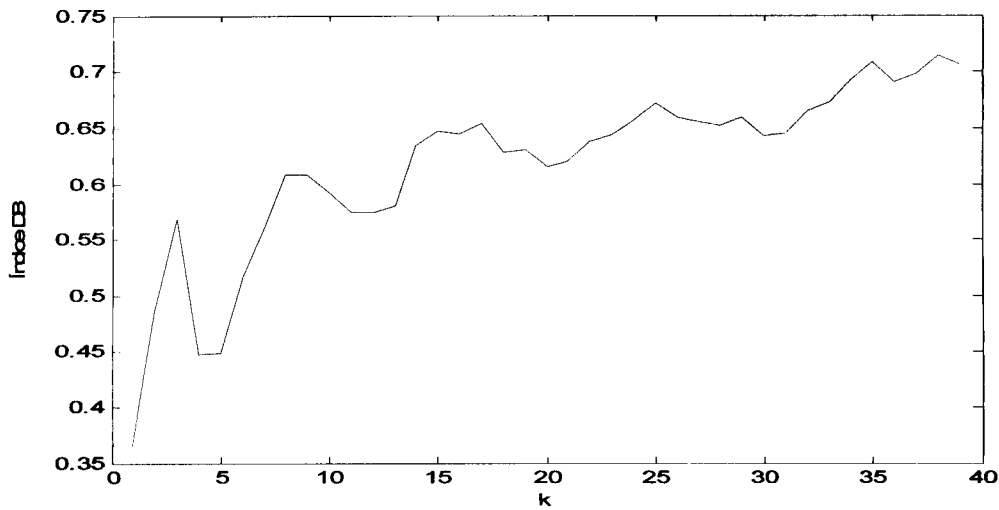


Figura 5.9 Representação dos valores do índice DB para a classificação pelo método *average*.

Os gráficos dos índices DB e Xu para este método são apresentados nas figuras 5.9 e 5.10.

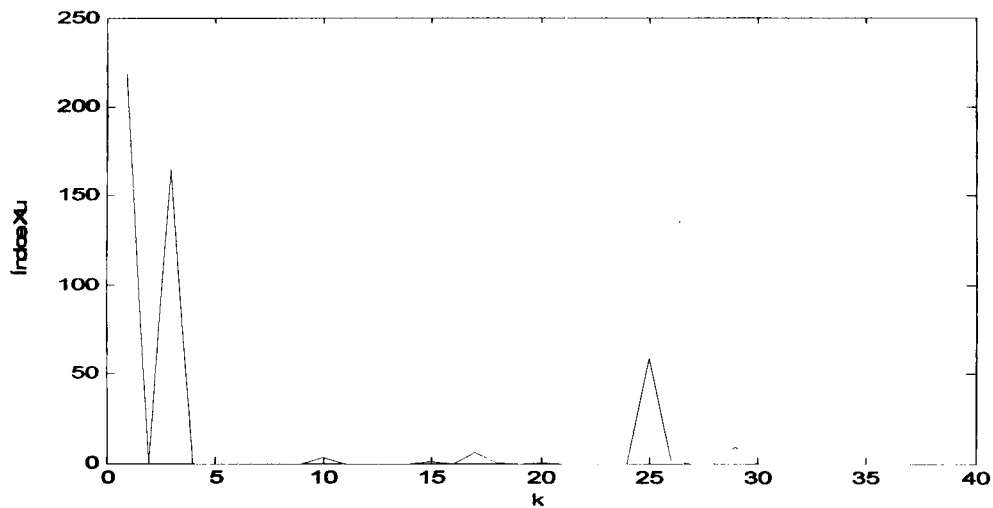


Figura 5.10 Representação dos valores do índice Xu para a classificação pelo método *average*.

O índice DB aplicado a todas as partições produzidas, desde 2 ($k=1$) até 40 ($k=39$) classes, assinala as partições em 2, 5 e 6 classes como as mais eficientes (figura 5.9). Por outro lado, o índice Xu aponta as partições dos dados em 2, 4 ou 26 classes (figura 5.10) como sendo as que representam melhor o conjunto de dados em análise. A partição em 2 classes é a única apoiada em simultâneo pelo índice DB. Estas 4 soluções estão apresentadas na figura 5.11.

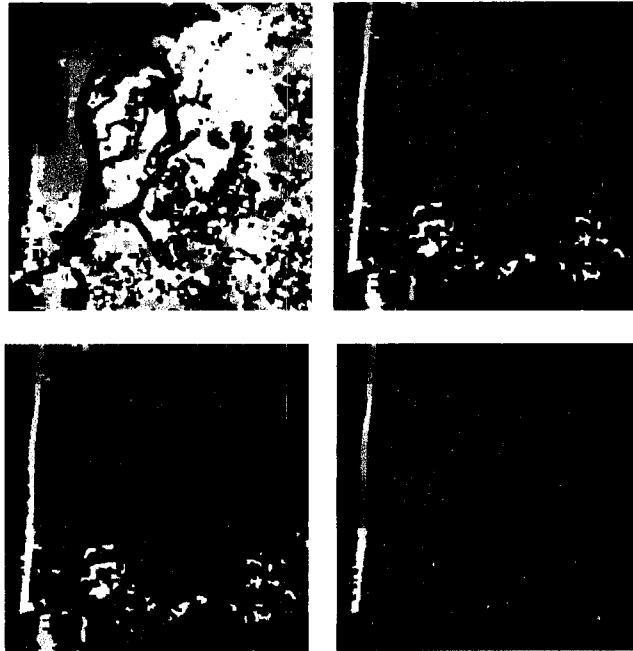


Figura 5.11 Imagem classificada em 26, 5, 4 e 2 classes pelo método de agregação *average*.

5.4. Método de Agregação *Weighted*

O resultado da classificação em 40 classes utilizando o método de agregação *weighted* é bastante satisfatório já que é clara, pela interpretação visual da imagem RGB (figura 5.1) e da imagem classificada (figura 5.12), a distinção efectuada entre as várias zonas que cobrem esta área.

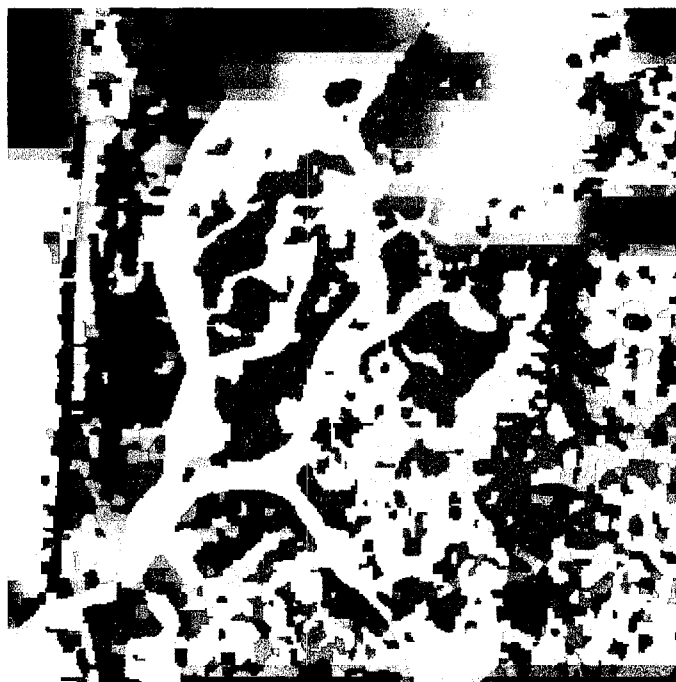


Figura 5.12 Imagem classificada em 40 classes pelo método de agregação *weighted*.

Os gráficos dos índices de DB e Xu são apresentados nas figuras 5.13 e 5.14. Segundo o índice DB (figura 5.13) a partição em 5 classes é a mais apropriada aos dados.

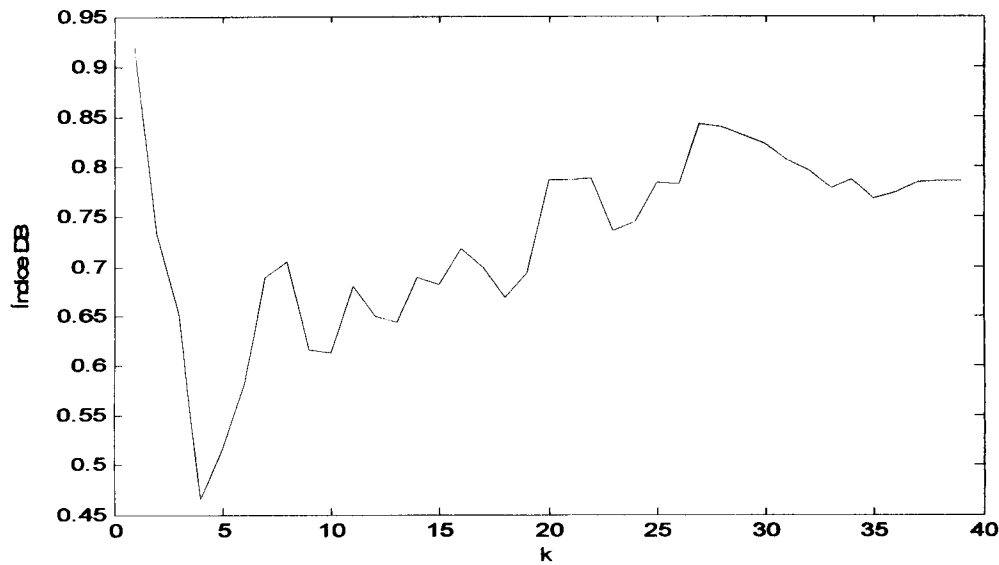


Figura 5.13 Representação dos valores do índice DB para a classificação pelo método *weighted*.

A estrutura que melhor representa os dados em análise de acordo com os resultados do índice Xu (figura 5.14) é a partição em 4 classes sendo a de 31 classes uma opção a não excluir.

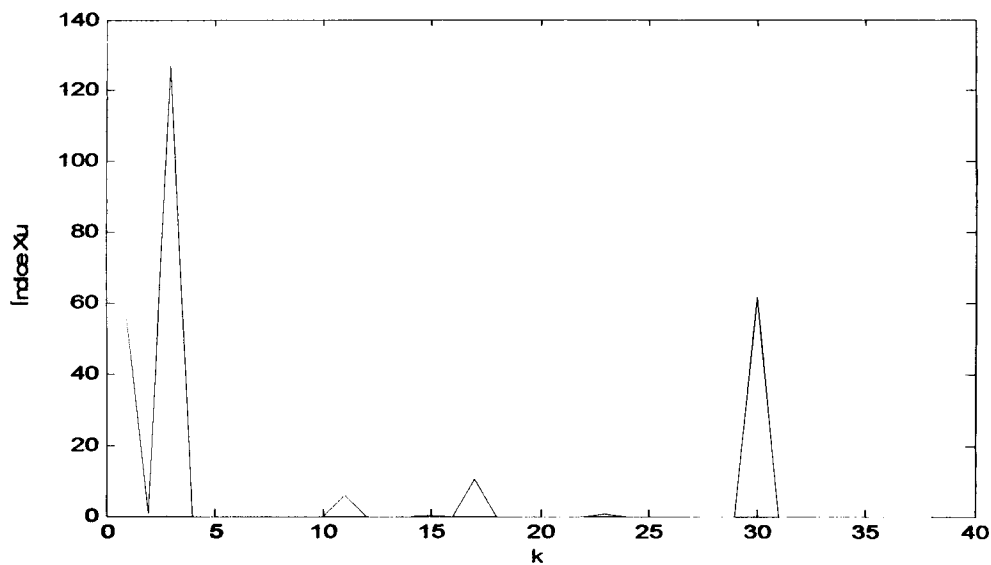


Figura 5.14 Representação dos valores do índice Xu para a classificação pelo método de agregação *weighted*.

Na figura 5.15 apresentam-se as imagens classificadas em 4, 5 e 31 classes.



Figura 5.15 Imagem classificada em 31, 5 e 4 classes pelo método de agregação *weighted*.

5.5. Método de Agregação *Ward*

O último método de agregação hierárquica testado foi o método *ward*. O resultado da aplicação deste método para 40 classes é apresentado na figura 5.16.



Figura 5.16 Imagem classificada em 40 classes pelo método de agregação *ward*.

Esta classificação distingue-se das anteriores no sentido em que as classes formadas são, em geral, de menores dimensões do que as classes apresentadas pelos restantes métodos. No canto superior esquerdo da imagem é possível visualizar, pela primeira vez, uma tentativa de separação da zona de rebentação (mar) da zona de mar alto. O detalhe é nesta imagem uma característica dominante, o que torna a partição em 40 classes difícil de avaliar visualmente. No entanto a avaliação do desempenho deste classificador poderá ser menos dificultada através de uma análise visual comparativa das

partições em menor número de classes com as áreas dominantes (água, vegetação, areia e urbanização) que formam a imagem RGB da figura 5.1.

Os índices DB e Xu são apresentados nas figuras 5.17 e 5.18.

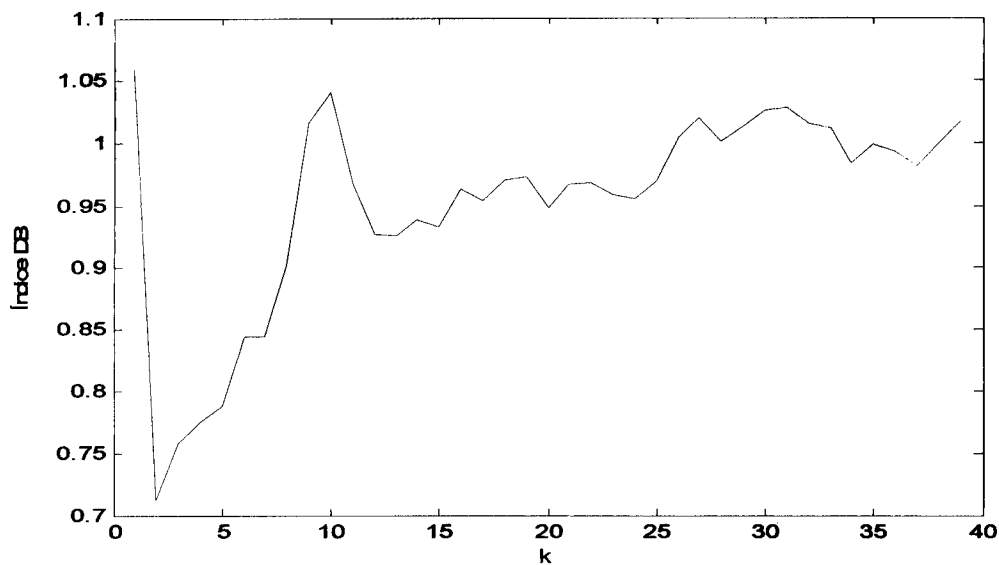


Figura 5.17 Representação dos valores do índice DB para a classificação pelo método de agregação *ward*.

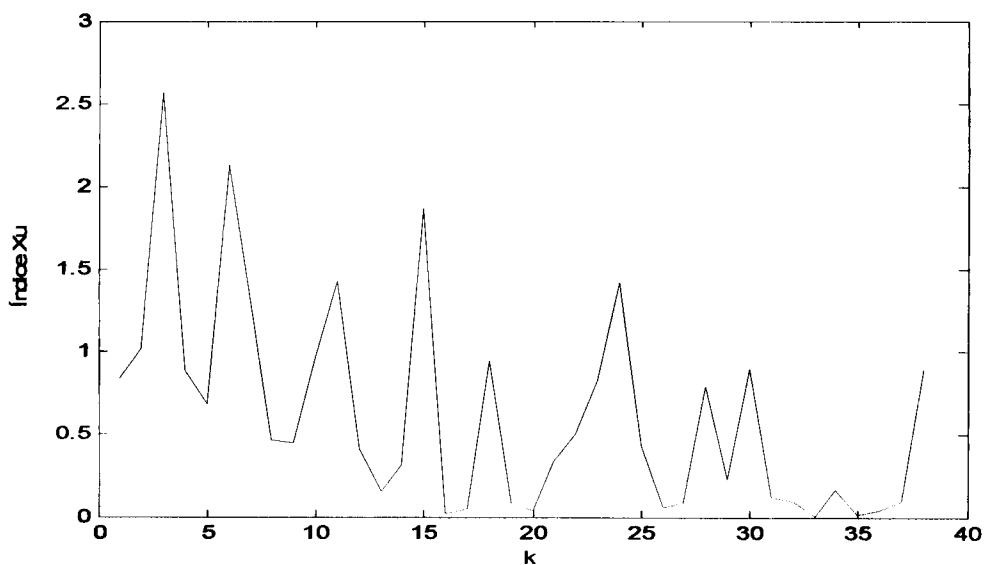


Figura 5.18 Representação dos valores do índice Xu para a classificação pelo método de agregação *ward*.

Uma breve análise ao gráfico representado na figura 5.17 indica que pelo índice DB a melhor partição para representar o conjunto de dados em estudo é a que os divide em 3 classes. O índice Xu elege as partições em 4, 7 e 16 classes como as que melhor representam o conjunto de dados referente à imagem da figura 5.1. As imagens classificadas nestes 4 níveis são apresentadas na figura 5.19.

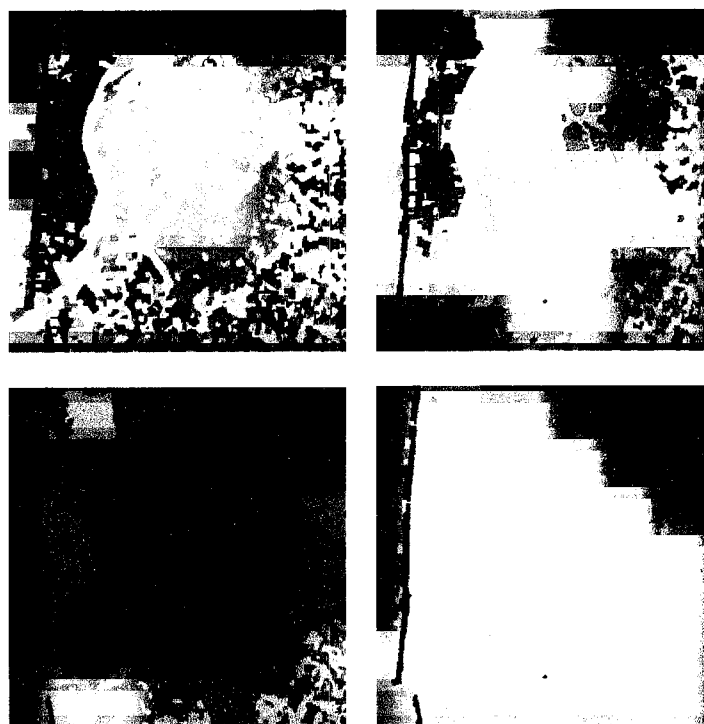


Figura 5.19 Imagem classificada em 16, 7, 4 e 3 classes pelo método de agregação *ward*.

Uma vez que este foi o método de agregação que produziu resultados mais interessantes, tendo como base a interpretação visual da imagem original, os resultados desta classificação são apresentados com mais detalhe em 15 níveis, na figura 5.20. As tabelas de cor utilizadas para estas imagens foram produzidas de acordo com a descrição efectuada na secção 4.4, sendo a cor da classe dominante mantida no processo de fusão entre duas classes.

É possível verificar, na figura 5.20, que a classificação em 2 classes distingue as zonas de terra e de água e, à medida que o número de classes vai aumentando, cada uma destas áreas é dividida com um grau de detalhe crescente. Na imagem resultante da classificação em 6 classes é possível verificar que a zona de água foi separada em duas: a zona mais profunda, que mantém a cor original da classe de água e a zona da ria de menor profundidade. Da mesma forma a zona de terra é dividida em quatro tipos de solo, entre os quais se distinguem: zonas de areia, zonas de urbanização e zonas verdes.

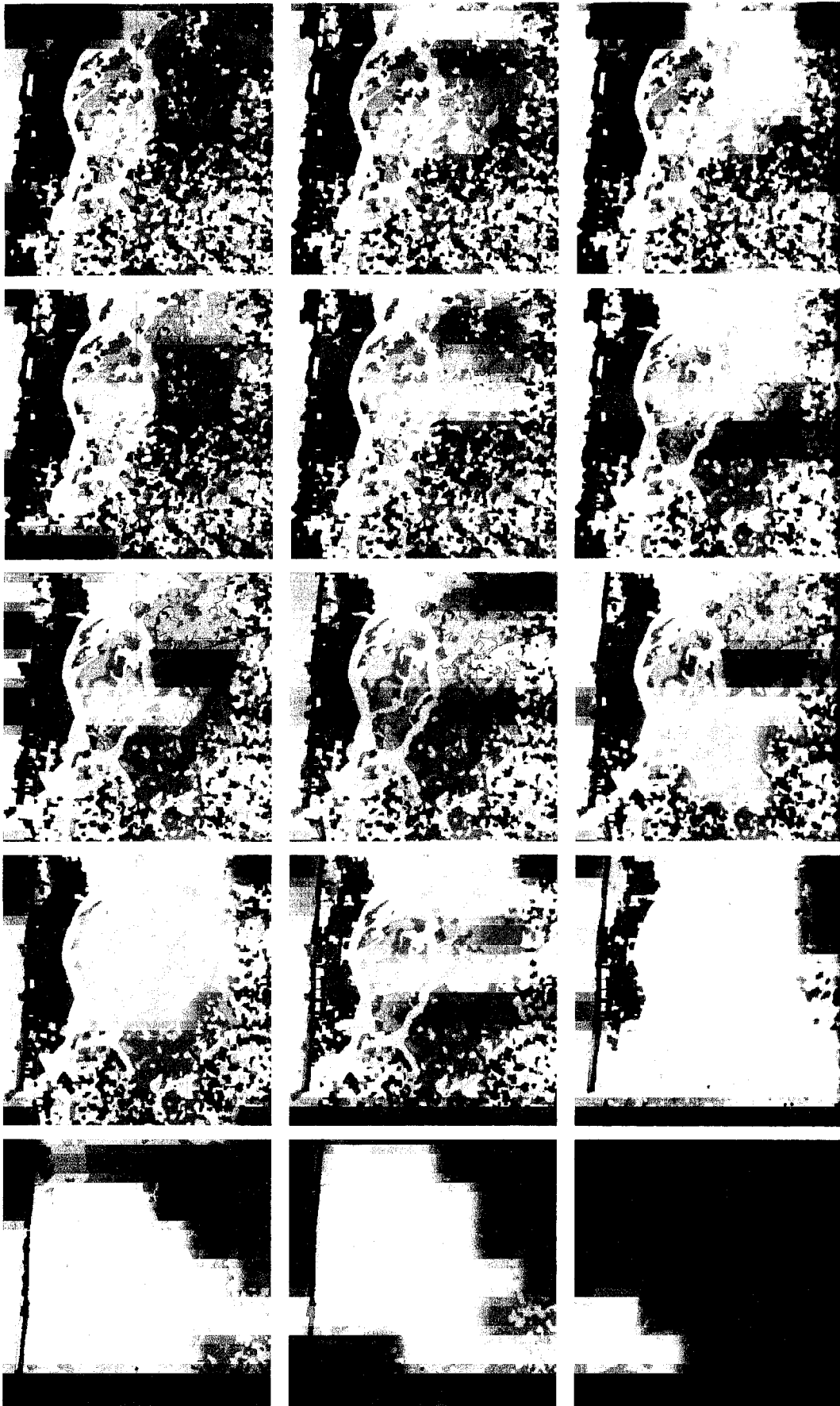


Figura 5.20 Imagem classificada em 40, 37, 34, 31, 28, 25, 22, 19, 16, 13, 10, 8, 6, 4, e 2 classes pelo método de agregação *ward*.

Capítulo 6

Conclusão

Numa primeira fase foram criados conjuntos de dados sintéticos com o objectivo de testar o desempenho dos métodos de agregação hierárquica disponíveis no MATLAB. Para cada método aplicado a cada conjunto de dados de teste, foram determinadas as partições que melhor representam a estrutura intrínseca esperada destes conjuntos.

Os métodos de classificação abordados, distinguíveis pela técnica de agregação de classes utilizada, produziram resultados bastante divergentes no que diz respeito ao conjunto de dados sintéticos. O que teve o melhor desempenho neste conjunto de dados foi o método de agregação pela distância mínima, *single*, já que foi este que produziu resultados mais próximos dos esperados para a maioria dos conjuntos de dados sintéticos. No entanto esta forma de agregação possui uma lacuna extremamente importante quando utilizada sobre conjuntos com classes de forma arredondada, produzindo resultados pouco satisfatórios por não captar a estrutura natural das classes. A classificação utilizando os restantes métodos de agregação, *complete*, *average*, *weighted* e *ward*, teve um comportamento bastante aceitável na identificação de classes de forma arredondada, falhando na identificação de conjuntos com aspecto alongado, divergindo do método de agregação *single*.

A aplicação dos índices DB e Xu para a selecção da melhor partição para cada conjunto de dados sintéticos, produziu resultados coerentes. Nas classificações consideradas bem sucedidas por aproximarem de forma satisfatória a estrutura natural das classes em alguma partição, os índices conseguiram na sua maioria identificar de forma aceitável essa partição.

Numa abordagem com dados reais, foi escolhida uma imagem de satélite multi-espectral da zona da Ria de Aveiro. Devido ao elevado volume de dados da imagem escolhida, foi efectuado um processo de redução de dados utilizando processos de filtragem e segmentação. Ao conjunto de dados reduzido são então aplicados os métodos de classificação por agregação hierárquica.

Na análise da imagem multi-espectral, os resultados da classificação pelo método de agregação *single* foram mais uma vez bastante diferentes dos resultados obtidos pela aplicação dos restantes métodos. Por este método, obteve-se uma separação notória desde a partição em 40 classes, entre a zona de areia e as restantes. Provavelmente devido ao facto da

areia possuir características espectrais bastante distintas das restantes zonas, este método identificou na zona de areia a presença de cerca de 39 classes, agregando quase todas as zonas de água, verdes e urbanização numa mesma classe. Em relação ao desempenho dos métodos *complete*, *average* e *weighted*, nota-se uma distribuição satisfatória das 40 classes iniciais pelas diversas zonas de ocupação do solo. No entanto, o sucedido para o método *single*, manifesta-se também para estes métodos quando o número de classes se vai reduzindo: na partição em duas classes, divide a área total da imagem em zonas de areia numa classe e restantes zonas noutra classe. Por sua vez, pelo método de agregação *ward*, obtiveram-se resultados mais interessantes. Os índices DB e Xu, ao contrário do sucedido para os conjuntos de dados sintéticos, foram quase sempre discordantes na eleição da melhor partição, com excepção dos métodos *complete* e *average*, para os quais a partição em 2 classes é a escolhida pelos dois índices. A distinção apresentada pelo método de agregação *ward*, na partição em duas classes, entre zonas de terra e zonas de água, permite avançar a hipótese de que este será provavelmente o método com melhor desempenho por os seus resultados serem aqueles que melhor se aproximam da classificação esperada da imagem de satélite ASTER.

O método proposto e desenvolvido com recurso às ferramentas disponíveis no MATLAB, permite a classificação de conjuntos com elevado número de elementos, como é o caso de imagens multi-espectrais de satélite, através da redução do número de padrões a classificar por segmentação de imagem. Desta forma é então possível a aplicação de classificadores de aglomeração hierárquica a imagens, produzindo resultados satisfatórios. Em linha com o trabalho realizado, surge a necessidade de uma análise aprofundada do desempenho de cada classificador, de forma a seleccionar o que melhor se adequa aos conjuntos de dados reais a classificar.

Referências

- Dubes, R. C., *How many clusters are best – an experiment*, Pattern Recognition, Vol 20, No 6, pp. 645-663, 1987.
- Duda, R. O., Hart, P. E., Stork, D. G., *Pattern Classification*, Second Edition, John Wiley & Sons, 2001.
- Fu, K. S., Mui, J. K., *A Survey on Image Segmentation*, Pattern Recognition, Vol 13, pp. 3-16, 1981.
- Gonzalez, R. C., Woods, R. E., *Digital Image Processing*, Second Edition, Prentice-Hall, 2002.
- Gonzalez, R. C., Woods, R. E., Eddins, S. L., *Digital Image Processing using MATLAB*, Pearson Prentice-Hall, Inc., 2004.
- Marçal, A. R. S., *Automatic Color Indexing of Hierarchically Structured Classified Images*, IEEE IGARSS '05 Proceedings, Vol 7, IEEE, pp. 4976-4979, 2005.
- Soh, L., Tsatsoulis, C., *Segmentation of Satellite Imagery of Natural Scenes Using Data Mining*, IEEE Transactions on Geoscience and Remote Sensing, Vol 37, No 2, pp. 1086-1099, Março 1999.
- Sung, Wing-Kin, 2003, *Combinatorial methods in bioinformatics 2003/2004 Semester 1, Lecture 8: Phylogenetic Tree Reconstruction: Distance Based*, http://www.comp.nus.edu.sg/~bioinfo/phylogenetic%20tree%20reconstruction_2_8.pdf, visitado em Junho 2006.
- The MathWorks: MATLAB The Language of Technical Computing – Using MATLAB: version 7. The MathWorks, Inc., 2004.
- Vincent, L., Soille, P., *Watersheds in Digital Spaces: An Efficient Algorithm Based on Immersion Simulations*, IEEE Transactions of Pattern Analysis and Machine Intelligence, Vol. 13, No. 6, pp. 583-598, June 1991.
- Xu, S., Kamath, M. V., Capson, D. W., *Selection of partitions from a hierarchy*, Pattern Recognition Letters, Vol. 14, No. 1, pp. 7-15, 1993.
- Yamaguchi, Y., Kahle, A. B., Tsu, H., Kawakami, T. Pniel, M., *Overview of Advanced Spaceborne Thermal Emission and Reflection Radiometer (ASTER)*, IEEE Transactions on Geoscience and Remote Sensing, Vol 36, pp. 1062-1071, 1998.

9

TESE Nº 168
CD-ROM

TESE Nº 168
CD-ROM

TESE M. ENG. NAT.

MARIA LUISA CASTRO