# Chapter 4

# A protocol for classifying ecologically relevant marine landscapes, a statistical approach

# 4 A protocol for classifying ecologically relevant marine landscapes, a statistical approach

**Abstract**

Mapping ecologically relevant zones in the marine environment has become increasingly important. However, biological data are scarce and alternatives are being sought in optimal classifications of abiotic variables. The concept of 'marine landscapes' is based on a hierarchical classification of geological, hydrographic and other physical data. However, this approach is subject to many assumptions and subjective decisions.

Here, an objective protocol is being proposed where abiotic variables are subjected to a statistical approach, using principal components analysis (PCA) and a cluster analysis. The optimal number of clusters is being defined using the Calinski-Harabasz criterion. The methodology has been applied on datasets of the Belgian part of the North Sea (BPNS), a shallow sandy shelf environment with a sandbank-swale topography.

The BPNS was classified into 8 marine landscapes that represent well the natural variability of the seafloor. The internal cluster consistency was validated with a split-run procedure, with more than 99% correspondence between the validation and the original dataset. The ecological relevance of 6 out of the 8 clusters was demonstrated, using indicator species analysis.

The proposed protocol, as exemplified for the BPNS, can easily be applied to other areas and provides a strong knowledge basis for environmental protection and management of the marine environment. A SWOT-analysis, showing the strengths, weaknesses, opportunities and threats of the protocol was performed.

*Keywords:* marine landscapes, abiotic variables, macrobenthic species, Principal Components Analysis, cluster analysis, Belgian part of the North Sea

## 4.1  <u>Introduction</u>

Biodiversity is of utmost importance to maintain the long-term stability of ecosystems, certainly with changing environmental conditions, such as global warming (Keytsman and Jones 2007). This applies to both terrestrial and marine habitats, species or communities, many of which are threatened by the ever-growing pressure on their environment.

Several techniques to map the environment are in place and if we consider predictive modelling and classification techniques for habitat mapping, all of them are based on the assumption that the biological value of an area is related to its abiotic characteristics. Generally, species and communities are linked to their substrate type, topographic position and energy regime. The information on the abiotic environment is generally more widely available than biological information itself; as such, it is aimed at distinguishing ecological landscapes on the basis of specific combinations of these abiotic variables.

Terrestrial examples of classifications of abiotic variables can be found in Fairbanks et al. (2000), Jobin et al. (2003), Rosa-Freitas et al. (2007) and Svoray et al. (2007). Similar methodologies can be applied to the marine environment.

In the framework of marine protection or management, available biotic data (e.g. absence/presence of benthic organisms) are often patchy and highly variable in nature. Moreover, offshore areas are generally devoid of samples. In response, the mapping of "marine landscapes" was developed, as a surrogate of biologically driven habitat mapping. If reliable, this methodology would facilitate the development of management measures for offshore areas in the absence of biological data (obtaining biological data in offshore areas is extremely expensive and time consuming). This hierarchical abiotic classification was first proposed for Canadian waters by Roff and Taylor (2000) and Roff et al. (2003). In this concept, biological data are used only passively, as a validation tool afterwards.

The integration of abiotic datasets (e.g. seabed substrata, depth, slope) that lead to a classification of seabed features, can be performed in a Geographic Information System (GIS) and is now applied widely (e.g. Golding et al. 2005; Connor et al. 2006; and Al-Hamdani and Reker 2007). The advantages are that abiotic parameters and processes are relatively easy to observe and monitor. Moreover, they can often be correlated with biological species or communities (Zacharias and Roff 2000). Another advantage is that the GIS process is quite simple, compared to statistical techniques (e.g. clustering as proposed in this study).

Unfortunately, this approach still lacks objectivity: in several stages of the methodology, subjective decisions have to be made: (1) 'Ecologically relevant abiotic variables' have to be selected as input for the GIS analyses. However, since biological data are generally sparse, this selection is not straightforward; (2) the analysis needs abiotic variables, classified into 'relevant classes in terms of biology'. It is very hard to define both the relevant class breaks and the number of classes; and (3) the 'Queries' step that combines the predefined classes of the abiotic variables into new combinations, being the final marine landscapes. As such, there is a strong need for a more objective and repeatable methodology.

This paper proposes a protocol to increase the objectivity of the marine landscapes approach, based on a statistical analysis for the grouping of full coverage abiotic data. The performance of a combination of PCA and cluster analysis will be demonstrated.

The proposed protocol aims for an unsupervised classification of purely abiotic variables. The ecological validation is done afterwards, independently from the PCA

and the cluster analysis to test the ecological relevance of the marine landscapes. The Belgian part of the North Sea (BPNS) is an ideal case study area, because of its extensive availability of both abiotic and biotic variables. However, in most cases, abiotic datasets will be available and only a few or no biological datasets.


## 4.2   Material and Methods


### 4.2.1   Study area

The BPNS (3600 km²) is situated on the North-West European Continental Shelf. The shelf is relatively shallow and dips gently from 0 to 50 m. The seabed surface is characterized by a highly variable topography, with a series of sandbanks and swales. The sandbanks can be subdivided into four major groups: the Coastal Banks and the Zeeland Banks are quasi-parallel to the coastline, whereas the Flemish Banks and the Hinder Banks have a clear offset in relation to the coast (Lanckneus et al. 2001). The seabed is sandy; the sand fraction (0.063 - 2 mm) is merely found on the sandbanks, whereas coarser sands, gravel (> 2 mm) and higher silt-clay fractions (< 0.063 mm) are found also in the swales (Lanckneus et al. 2001). The sandbanks and the swales are both covered with ripples and dunes. The height of the dunes commonly ranges between 2 and 4 m, though dune heights of up to 11 m are found in the most offshore areas.

Five macrobenthic communities (four subtidal and one intertidal community) are discerned within the mobile substrates of the BPNS (Degraer et al. 2003; and Van Hoey et al. 2004).

On the BPNS, various abiotic datasets are widely available (Van Lancker et al. 2007). In addition, a large dataset of 741 macrobenthic samples (Marine Biology Section, Ugent – Belgium, 2008) can be used for an ecological validation. As such, the BPNS is an ideal test area to develop a new classification method and to validate its ecological relevance.

16 abiotic variables are available for the BPNS (Table 4.1). All of them have a resolution of 250 m, except maximum Chlorophyl a (Chl a) concentration and maximum Total Suspended Matter (TSM) with a resolution of 1000 m. All data grids of the abiotic variables were resampled to 54307 pixels with a resolution of 250 m. Although other abiotic variables (e.g. salinity, temperature, stratification) could be important as well for explaining the presence of benthic species, they were not available for this study. Still, the current dataset represents well the abiotic variability. In the Discussion Section, this topic is discussed in more detail.


### 4.2.2   Research strategy

The protocol starts with a PCA for data reduction (step 1). The resulting components are then subjected to a hierarchical cluster analysis (step 2) and the cluster centres from step 2 are used as starting positions for a *K*-means partitioning (step 3). In step 4, the optimal number of clusters is calculated; in step 5, a validation of the internal cluster consistency is performed; and in step 6, a species indicator analysis (INDVAL) is done (Dufrêne and Legendre 1997), defining for each cluster a number of

significant indicator species and as such offering the possibility for an ecological validation of the classification.

Software used is SPSS version 12 for PCA, ClustanGraphics version 8.03 for the hierarchical and *K*-means clustering, R version 2.5.1 for the calculation of the Calinski and Harabasz (1974) indices (called C-H in this paper) and PC-ORD 4.41 (McCune and Mefford 1999) for the INDVAL analysis.


### 4.2.3    Step 1: PCA analysis

For data reduction and to avoid multicollinearity (i.e. high degree of linear correlation) of the abiotic variables, a PCA was performed (theoretical background e.g. in Jongman et al. 1987; Legendre and Legendre 1998). PCA computes a reduced set of new, linearly independent variables, called principal components (PCs) that account for most of the variance of the original variables. The PCs are a linear combination of the original variables. The PCA was based on a correlation matrix, implying that the Kaiser-Guttman criterion could be applied (Legendre and Legendre 1998). This means that PCs with eigenvalues larger than 1 were preserved as meaningful components for the analysis. To maximize the independence of each PC, a Varimax rotation of the PCs was computed. The PCs were the input variables for the cluster analysis.

Similar applications of PCA for data reduction of abiotic variables are found in Cardillo et al. (1999), Fairbanks (2000), Moreda-Piñeiro et al. (2006), and Frontalini and Coccioni (2008).


### 4.2.4    Step 2: Hierarchical cluster analysis based on Ward's method

To group the pixels with abiotic data on a statistical basis, a hierarchical clustering, based on Ward's (1963) or Orlóci 's (1967) minimum variance method was applied on the PCs (theoretical background e.g. in Jongman et al. 1987; Legendre and Legendre 1998). This method is an agglomerative clustering algorithm that minimizes an objective function which is the same "squared error" criterion that is used in multivariate analysis of variance and results into clusters with a minimal variance between each cluster. At each clustering step, this method finds the pair of objects or clusters whose fusion increases as little as possible the sum, over all objects of the squared Euclidean distances between objects and cluster centroids (Legendre and Legendre 1998). The Euclidean distance is an appropriate model for the relationships among abiotic variables (Legendre and Legendre 1998). Applications of Ward's method for the clustering of abiotic variables can be found in Cao et al. (1997) and Frontalini and Coccioni (2008).

**Table 4.1: Abiotic variables as input for the PCA and cluster analysis.**

| Abiotic variable | Unit | Reference or procedure |
|---|---|---|
| _Sedimentology_ | | Reference: sedimentological database ('sedisurf@') hosted at Ghent University, Renard Centre of Marine Geology. |
| • Median grain-size of sand fraction (63-2000 μm) or $d_s50$ | μm | Reference: Verfaillie et al. (2006) |
| • Silt-clay percentage (0-63 μm) | % | Reference: Van Lancker et al. (2007) |
| • Sand percentage (63-2000 μm) | % | Reference: Van Lancker et al. (2007) |
| • Gravel percentage (> 2000 μm) | % | Reference: Van Lancker et al. (2007) |
| _Topography_ <br> • Digital terrain model (DTM) of bathymetry | m | Reference: Flemish Authorities, Agency for Maritime and Coastal Services, Flemish Hydrography <br> All other topographic variables are derived from the DTM |
| • Slope = a first derivative of the DTM | ° | Procedure: Evans (1980); Wilson et al. (2007) |
| Aspect = a first derivative of the DTM Indices of northness and eastness provide continuous measures (−1 to +1) describing orientation of the slopes. | | Procedure: Wilson et al. (2007); Hirzel et al. (2002a) |
| • Eastness = sin (aspect) | / (no unit) | |
| • Northness = cos (aspect) | / | |
| • Rugosity = ratio of the surface area to the planar area across the neighbourhood of the central pixel | / | Procedure: Jenness (2002); Lundblad et al. (2006); Wilson et al. (2007) |
| Bathymetric Position Index (BPI) = measure of where a location, with a defined elevation, is relative to the overall landscape | | Procedure: Lundblad et al. (2006); Wilson et al. (2007) |
| • BPI (large scale) | / | |
| • BPI (small scale) | / | |
| _Hydrodynamics_ <br> • Maximum bottom shear stress = frictional force exerted by the flow per unit area of the seabed | N/m² | Reference: Management Unit of the North Sea Mathematical Models and the Scheldt estuary (MUMM) |
| • Maximum current velocity | m/s | |
| _Satellite derived variables_ <br> • Maximum near-surface Chlophyl a (Chl a) concentration over a 2-year period (2003-2004) | mg/m³ | Reference: MERIS data processed by MUMM in the framework of the BELCOLOUR-2 project (ESA ENVISAT AOID3443) |
| • Maximum near-surface Total Suspended Matter (TSM): measure for turbidity over a 2-year period (2003-2004) | mg/l | |
| • Distance to coast | m | Computed in GIS |

### 4.2.5 Step 3: K-means partitioning

Although the result of a hierarchical cluster analysis on its own is prone to multiple errors, a hierarchical clustering, based on Ward's method, can generate excellent starting positions (i.e. cluster centroids used as cluster seeds) for a *K*-means partitioning (Milligan 1980; Legendre and Legendre 1998; and Wishart 1987). Partitioning clustering methods produce clusters in a predefined number of groups (K). *K*-means is the most widely used numerical method for partitioning data (examples from the marine environment are found in Legendre et al. (2002); Legendre (2003); Preston and Kirlin (2003); Hewitt et al. (2004); and Zharikov et al. (2005)). Pixels from clusters are allocated to a cluster in which the distance to its centre is minimal. The procedure stops if all pixels have been allocated. A *K*-means procedure exists of 3 steps: the initiation of the starting cluster centres, the allocation of pixels to the initial clusters and the re-allocation of pixels to another cluster. The starting positions and the allocation of the pixel to the initial clusters were taken from the hierarchical clustering, based on Ward's method. Those pixels are clustered that show the smallest increase in the Euclidean Sum of Squares.

### 4.2.6 Step 4: Number of clusters

The most difficult and most subjective decision in the cluster analysis is the number of clusters. Several indices to calculate the optimal number of clusters exist. From a simulation study comparing 30 indices, Milligan and Cooper (1985), proposed the C-H criterion as giving the best results. C-H is the *F*-statistic of multivariate analysis of variance and canonical analysis. *F* is the ratio of the mean square for the given partition, divided by the mean square for the residuals. The number of clusters corresponding with the highest C-H value is the optimal solution in the least-squares sense. C-H was also used as stopping criterion for cluster analysis in the marine environment in Legendre et al. (2002); Hewitt et al. (2004); and Orpin and Kostylev (2006).

### 4.2.7 Step 5: Validation of internal cluster consistency

A cluster analysis automatically allocates each individual to a cluster. To evaluate the internal consistency of the cluster composition, a validation with a split-run procedure was performed. For this procedure, the cluster analysis was first done for the whole dataset. After that, the optimal number of clusters was computed with the C-H criterion. Next, the dataset was randomly split into 2 equal validation parts to which the cluster analysis was applied with the same number of clusters. Finally, the cluster compositions from both validation parts were compared with the original cluster composition by calculating the number of differently classified pixels.

### 4.2.8 Step 6: Indicator species analysis of the clusters

To evaluate whether the obtained clusters have an ecological relevance, a species indicator analysis or INDVAL (Dufrêne and Legendre 1997) was performed. This method identifies indicator species for each of the clusters: if indicator species are

identified then the cluster should have an ecological relevance, whereas if no indicator species can be identified, (most probably) the cluster has no ecological significance. The index is maximum when all individuals of a species are found in a single group of sites and when the species occurs in all sites of that group. The INDVAL index is defined as follows:

$$INDVAL_{ij} = A_{ij} \times B_{ij} \times 100$$

with  $A_{ij}$ = Nindividuals$_{ij}$/Nindividuals$_i$ or the mean abundance of species $i$ in the sites of group $j$, compared to all groups in the study. $A_{ij}$ is a measure of specificity and is maximum when species $i$ is only present in cluster $j$.
$B_{ij}$ = Nsites$_{ij}$/Nsites$_j$ or the relative frequency of occurrence of species $i$ in the sites of group $j$. $B_{ij}$ is a measure of fidelity and is maximum when species $i$ is present in all sites of cluster $j$.

The index is maximal when all individuals of a species are found in a single group of sites and when the species occurs in all sites of that group. The statistical significance for the species indicator values is evaluated using a Monte Carlo permutation procedure. 1000 random permutations were used for this study.
Examples of applications of INDVAL to test the ecological relevance of predefined clusters can be found in Mouillot et al. (2002); Heino and Mykrä (2006); and Perrin et al. (2006).

## 4.3   Results

### 4.3.1   Step 1: PCA analysis

Retaining only those PCs with eigenvalues larger than 1; PCA resulted in 6 PCs, explaining 78.0% of the total variance. The rotated component matrix (Table 4.2) shows the factor loads, being the correlations between the rotated PCs and the original variables.
In decreasing order, PC 1 has high loads (r < -0.5 or r > 0.5) for the variables distance to coast, DTM, maximum TSM, $d_s50$, maximum Chl a, silt-clay % and gravel %; PC 2 for maximum bottom shear stress and maximum current velocity; PC 3 for slope and rugosity; PC 4 for BPI large scale and BPI small scale; PC 5 for eastness and northness; and PC 6 for sand % and gravel %. Gravel % is the only variable that has a high load for 2 PCs, meaning that this relationship is not exclusive.

### 4.3.2   Step 2: Hierarchical cluster analysis based on Ward's method

The 54307 cases with 6 PC variables were clustered to achieve a hierarchical partition tree. This tree is not at all appropriate as end result of the clustering, but the partitions are very useful as starting positions for the $K$-means partitioning.

### 4.3.3    Step 3: K-means partitioning

The cluster centres of the partition tree based on Ward's method, were used as input for the *K*-means partitioning. Subsequently, new cluster centres based on the *K*-means algorithm were computed forming a cascade from 2 to 20 clusters. Those centres were used to compute the C-H criterion.
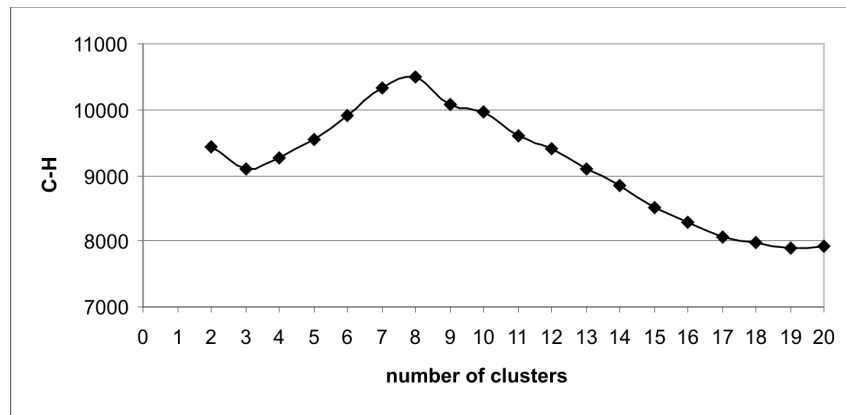
**Table 4.2: Component matrix showing correlations between the Varimax rotated PCs and the original variables.**
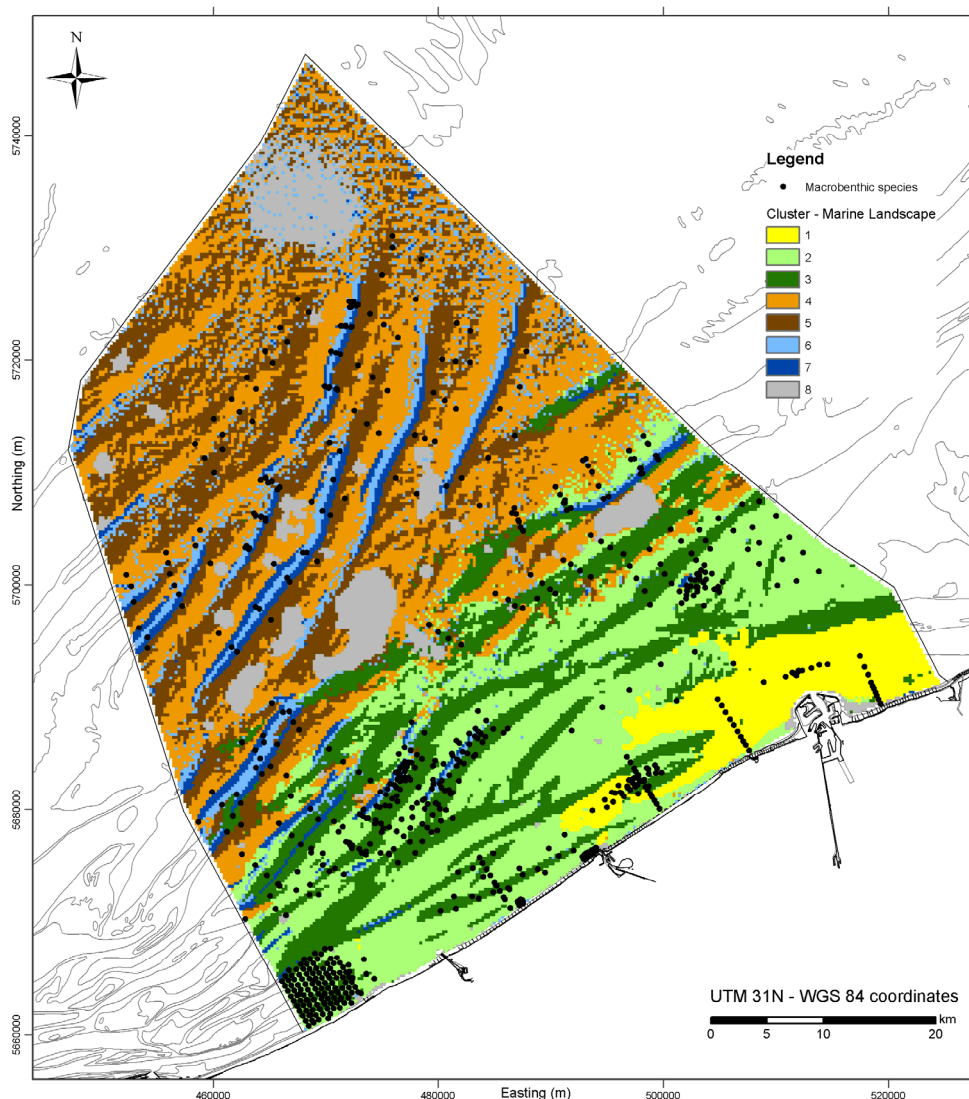**High factor loads (r < -0.5 or r > 0.5) are indicated in bold. Information of the variables can be found in Table 4.1.**

|  | Principal component | | | | | |
|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 |
| $d_s50$ | **-0.894** | -0.094 | 0.088 | 0.052 | 0.064 | 0.133 |
| silt-clay % | **0.668** | 0.467 | -0.135 | -0.069 | -0.079 | -0.285 |
| sand % | -0.230 | -0.487 | 0.149 | 0.098 | 0.064 | **0.748** |
| gravel % | **-0.514** | 0.071 | -0.069 | -0.034 | -0.016 | **-0.665** |
| DTM | **0.932** | -0.075 | 0.028 | 0.212 | -0.048 | 0.039 |
| slope | -0.054 | 0.014 | **0.958** | 0.031 | 0.019 | 0.045 |
| eastness | -0.021 | 0.041 | -0.005 | 0.004 | **0.828** | 0.040 |
| northness | 0.105 | -0.027 | -0.050 | 0.000 | **-0.798** | 0.046 |
| rugosity | -0.184 | 0.037 | **0.909** | 0.186 | 0.037 | -0.042 |
| BPI large scale | 0.074 | 0.048 | 0.170 | **0.862** | -0.001 | 0.048 |
| BPI small scale | -0.116 | 0.003 | 0.023 | **0.851** | 0.005 | -0.048 |
| Max. bottom shear stress | -0.029 | **0.918** | 0.040 | 0.089 | 0.142 | 0.032 |
| max current velocity | -0.184 | **0.912** | 0.055 | -0.005 | -0.029 | 0.021 |
| max chl a | **0.794** | -0.133 | -0.035 | -0.079 | -0.028 | 0.081 |
| max TSM | **0.921** | -0.097 | -0.151 | -0.052 | -0.034 | -0.009 |
| distance to coast | **-0.944** | 0.091 | 0.074 | 0.068 | 0.043 | -0.032 |

### 4.3.4    Step 4: Number of clusters and resulting clusters

Applying the C-H criterion (Figure 4.1), an optimum of 8 clusters was found. The result of the 8 cluster solution is presented in Figure 4.2 and Table 4.3. The clusters or marine landscapes represent well the natural environment and clear relationships with the original abiotic variables are visible.

**Figure 4.1: The number of clusters versus the C-H criterion.
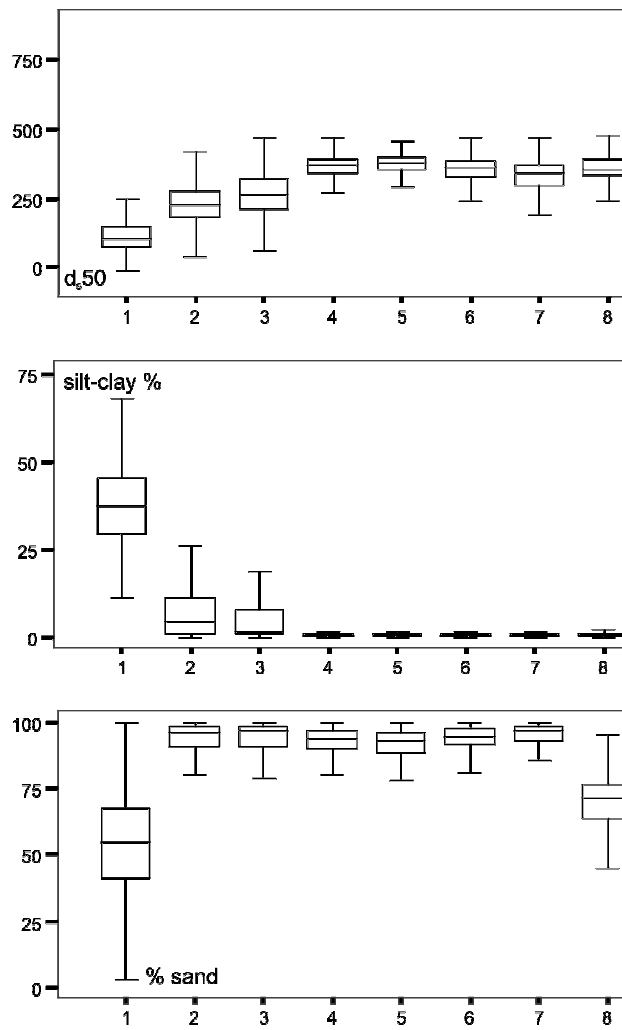C-H reaches an optimum for 8 clusters.**



**Figure 4.2: Belgian part of the North Sea with 8 clusters.
The location of macrobenthic community samples are plotted for validation.
Important patterns of the original abiotic variables are clearly visible on the
map: e.g. high silt-clay % in cluster 1, alternation of sandbanks and flats-
depressions in clusters 2, 3, 4, 5, 6 and 7; patches of gravel and shell fragments in
cluster 8.**

Boxplots (Figure 4.3) show the contribution of the original variables against the clusters. A clear example is the boxplot representing slope. This variable is approximately the same for all of the clusters, except for cluster 7 with higher values.

**Table 4.3: The 8 clusters and their characteristics based on the boxplots (Figure 4.3).**

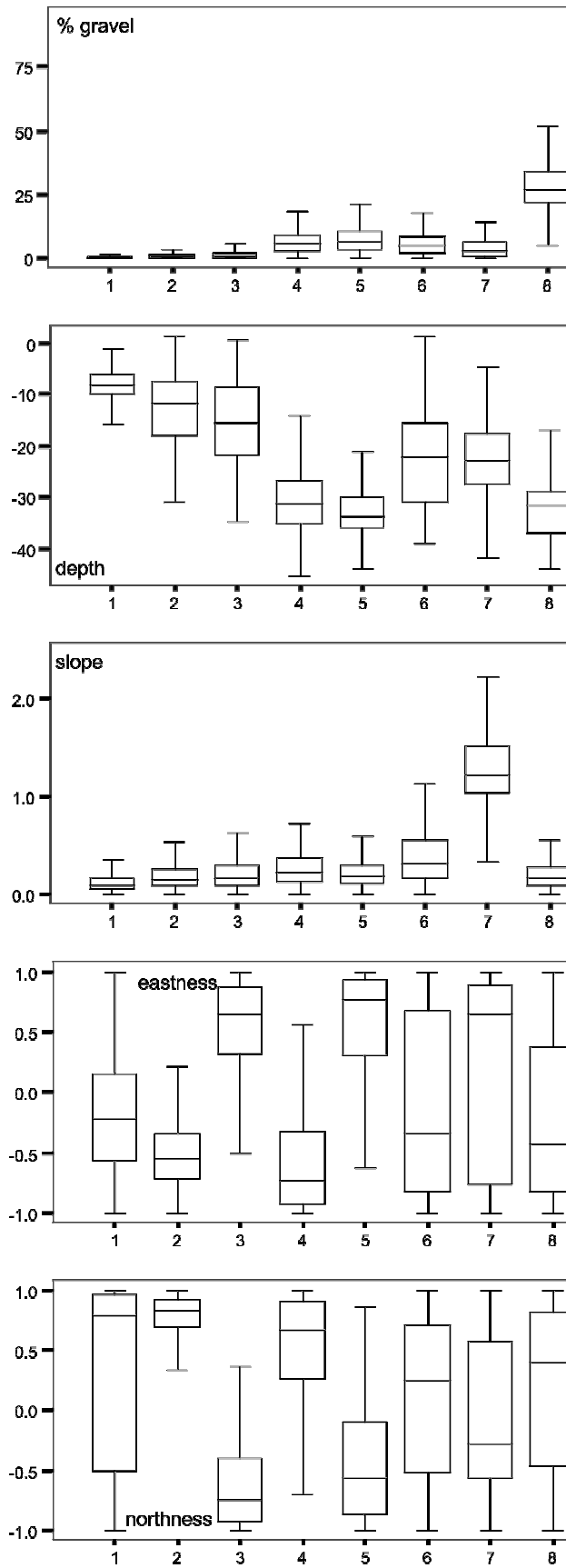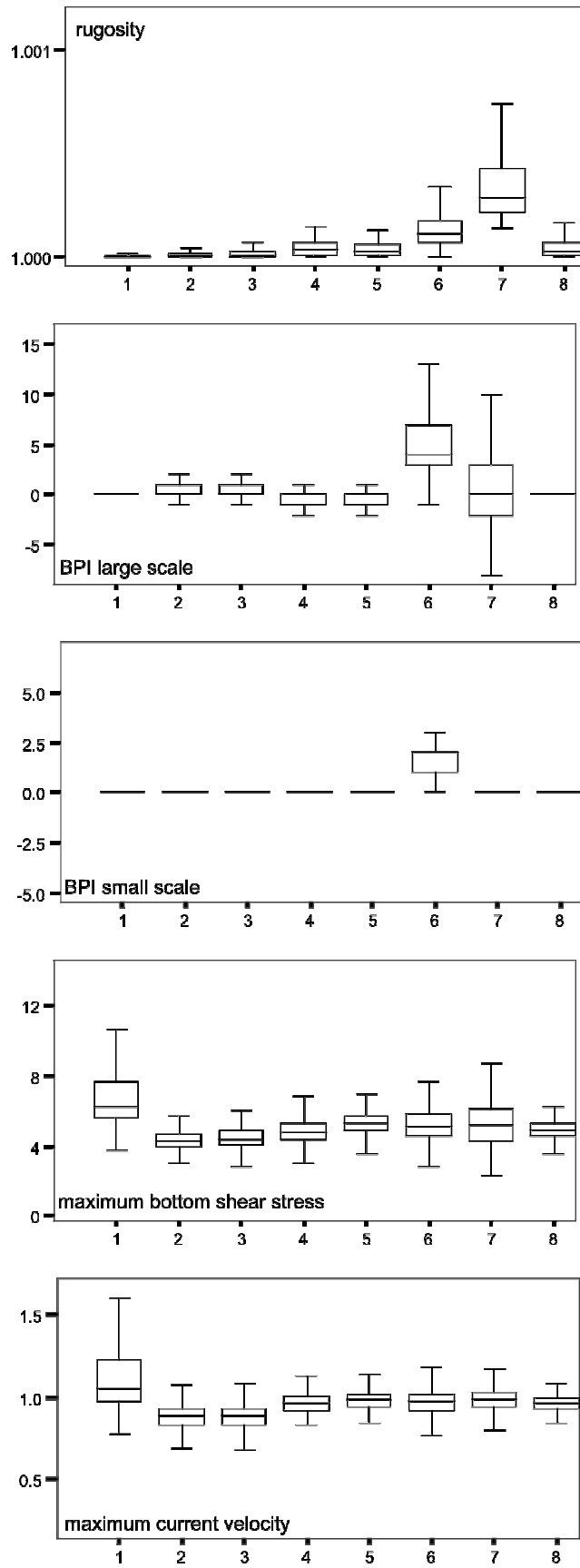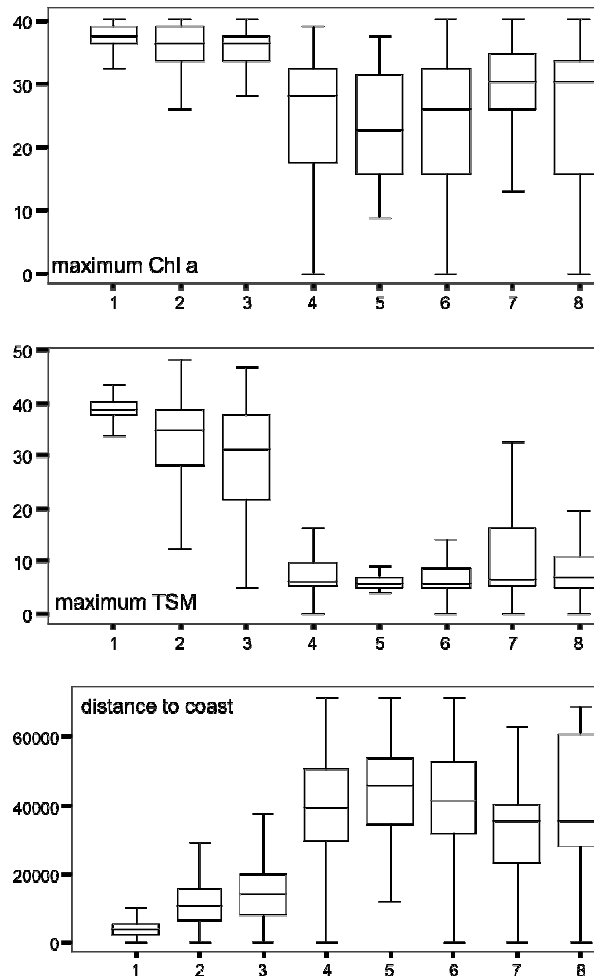| Cluster | Characteristics |
| --- | --- |
| 1 | Shallow, high silt-clay percentage, high current velocity, high bottom shear stress, turbid, high Chl a concentration |
| 2 | Shallow NW orientated flats and depressions, fine sand, slightly turbid, high Chl a concentration |
| 3 | Shallow SE orientated sandbanks, fine to medium sand, slightly turbid, high Chl a concentration |
| 4 | Deep NW orientated flats and depressions, medium sand |
| 5 | Deep SE orientated flats and depressions, medium sand |
| 6 | Crests of sandbanks, medium sand |
| 7 | Slopes of sandbanks, medium sand |
| 8 | High percentage of gravel – shell fragments |



**Figure 4.3a**

**Figure 4.3b**

**Figure 4.3c**

**Figure 4.3d: Boxplots of clusters (X-axis) against abiotic variables (Y-axis). An overview of the abiotic variables and their units is given in Table 4.1. The middle line in the box is the median, the lower and upper box boundaries mark the first and third quartile. The whiskers are the vertical lines ending in horizontal lines at the largest and smallest observed values that are not statistical outliers (values more than 1.5 interquartile range).**

### 4.3.5 Step 5: Validation of the internal cluster consistency

The split-run analysis showed very high correlations between the clusters obtained for the subsets and the clusters obtained for the whole dataset. Subset 1 contains 27153 cases, of which 159 have been classified differently as for the complete dataset. Subset 2 contains 27154 cases of which 184 have been classified differently. This is respectively 99.4 % and 99.3 % correspondence with the complete dataset for subset 1 and subset 2. The misclassified cases of both subsets were randomly distributed.
As shown by the split-run procedure, the internal cluster consistency is very good.

### 4.3.6 Step 6: Indicator species analysis of the clusters

Of the 123 species present in the 741 samples, randomization identified 25 species having a significant indicator value (5% level of significance) for 6 of the 8 defined clusters (Table 4.4). No indicator species could be found for cluster 3 and 6. This

102

means that both clusters do not show significant ecological differences with the other clusters.

Species with indicator values higher than 20 are *Cirratulidae* spp. and *Macoma balthica* for cluster 1; *Lanice conchilega* and *Spisula subtruncata* for cluster 2; *Echinocyamus pusillus* for cluster 4; *Tellina pygmaea*, *Gastrosaccus spinifer* and *Bathyporeia spp.* for cluster 5; and *Ophiura* spp. for cluster 7.

**Table 4.4: Significant indicator species analysis of the defined clusters.**

| species | cluster | INDVAL (%) | Randomised INDVAL (%) Mean | SD | p* | A (%) | B (%) |
|---|---|---|---|---|---|---|---|
| *Cirratulidae spp.* | 1 | **51.4** | 10.2 | 4.49 | 0.001 | **73** | **70** |
| *Macoma balthica* | 1 | **26.8** | 4.4 | 2.37 | 0.001 | **67** | 40 |
| *Glycera alba* | 1 | 14.5 | 5 | 2.55 | 0.011 | 48 | 30 |
| *Nephtys hombergii* | 2 | 19.7 | 8 | 2.6 | 0.005 | 41 | 48 |
| *Ensis spp.* | 2 | 19.4 | 7.3 | 3.35 | 0.015 | **65** | 30 |
| *Lanice conchilega* | 2 | **21.3** | 8.8 | 3.86 | 0.016 | **62** | 34 |
| *Phyllodoce mucosa / Phyllodoce maculata* | 2 | 17.7 | 7.7 | 3.6 | 0.026 | **57** | 31 |
| *Eumida spp.* | 2 | 12.3 | 5.5 | 2.87 | 0.035 | **57** | 22 |
| *Donax vittatus* | 2 | 9.3 | 4.4 | 2.43 | 0.042 | 45 | 21 |
| *Spisula subtruncata* | 2 | **23.3** | 14.1 | 4.98 | 0.046 | **75** | 31 |
| *Glycera capitata = Glycera lapidum* | 4 | 17.8 | 4.2 | 2.31 | 0.002 | 41 | 44 |
| *Echinocyamus pusillus* | 4 | **22.8** | 3.8 | 2.45 | 0.003 | **59** | 39 |
| *Branchiostoma lanceolatum* | 4 | 8.9 | 2.1 | 1.56 | 0.005 | **73** | 12 |
| *Pisione remota* | 4 | 8 | 2.2 | 1.82 | 0.017 | **76** | 11 |
| *Aonides oxycephala* | 4 | 10.1 | 3.4 | 2.61 | 0.03 | **64** | 16 |
| *Hesionura elongata* | 4 | 8.3 | 4 | 2.61 | 0.041 | 40 | 21 |
| *Thia scutellata* | 4 | 7.1 | 3.3 | 2 | 0.043 | 31 | 23 |
| *Tellina pygmaea* | 5 | **31.7** | 3.2 | 1.96 | 0.001 | **60** | **53** |
| *Gastrosaccus spinifer* | 5 | **23.3** | 7.1 | 2.77 | 0.002 | 35 | **66** |
| *Bathyporeia spp.* | 5 | 23 | 11.4 | 4.57 | 0.022 | 30 | **76** |
| *Pisidia longicornis* | 5 | 6.7 | 2.1 | 1.88 | 0.028 | **51** | 13 |
| *Ophiura spp.* | 7 | 25 | 11.6 | 5.13 | 0.019 | **62** | 40 |
| *Nephtys cirrosa* | 7 | 18.3 | 13.3 | 2.41 | 0.038 | 20 | **90** |
| *Aonides paucibranchiata* | 8 | 8.1 | 2.8 | 1.93 | 0.029 | 45 | 18 |
| *Bivalvia spp.* | 8 | 6.8 | 2.6 | 1.82 | 0.037 | **75** | 9 |

**p* Statistically significant at the 0.05 level; SD = standard deviation; A = specificity; B = fidelity; INDVAL values higher than 20% are marked in bold; A and B values higher than 50% are marked in bold.**

## 4.4  <u>Discussion</u>

This paper proposes an objective protocol to define ecologically relevant zones, solely on the basis of abiotic datasets. These zones are called 'marine landscapes', as they show a strong correlation with the abiotic variables and, in particular, the topography.

### 4.4.1  An objective method to define marine landscapes

The classical Marine Landscape methodology, as proposed by Roff and Taylor (2000); and Roff et al. (2003); and applied by Golding et al. (2004); Schelfaut (2005); Connor et al. (2006); and Al-Hamdani and Reker (2007) is highly subjective because of three reasons. First, the selection of ecologically relevant abiotic variables is biased. For the present protocol, no selection is necessary as input for a PCA; because PCs are constructed as linear combinations of the available, original abiotic variables (e.g. Cardillo 1999; and Fairbanks 2000). Secondly, there is a difficulty of classifying the selected abiotic variables into relevant classes. In this paper, a solution is proposed that abandons the classification and uses the continuous abiotic variables as input for the further analysis (e.g. Wilson 2007). Thirdly, the 'Queries' step is highly subjective because new combinations (the clusters or 'marine landscapes') are chosen arbitrarily from the predefined classes of the abiotic variables. This can be overcome by combining all possible classes, but this would lead rapidly to too many classes (e.g. 6 variables with 5 classes already means 30 landscapes). As such, this paper uses the C-H criterion to define a relevant number of clusters to automatically cluster the continuous abiotic variables (e.g. Legendre et al. 2002; Hewitt et al. 2004; and Orpin and Kostylev 2006).

With the objective approach proposed in this paper, there are still some decisions to be made during the analysis. First, for the cluster analysis, the number of groups has to be decided. Out of own physical knowledge of the BPNS, the solution of 8 marine landscapes seems to represent well the natural environment and none of the clusters seems to be useless. Their relation with the overall environment is clear, which was also exemplified by boxplots indicating the contribution of each abiotic variable to the clusters (Figure 4.3). The C-index (Hubert and Levin 1976), being a very good stopping criterion comparable to C-H (Milligan and Cooper 1985), has been tried as stopping criterion on this dataset, but it does not work for large datasets as used for this study. Secondly, for the $K$-means procedure, the Euclidean Sum-of-Squares clustering criterion was used as a distance index. As Punj and Stewart (1983) demonstrated, the choice of the (dis)similarity or distance index is of minor importance, compared to the clustering algorithm.

### 4.4.2  Abiotic datasets

Degraer et al. (2008) already discussed the many abiotic variables that might explain the distribution of macrobenthic communities on the BPNS. For the present study, not only typical variables, such as bathymetry and sedimentological information (e.g. Wu and Shin 1997; Van Hoey et al. 2004; and Willems et al. 2008) were used, but also hydrodynamical data (e.g. Caeiro et al. 2005), turbidity (e.g. Akoumianaki and Nicolaidou 2007), topographically derived features such as BPI (e.g. Lundblad et al. 2006; Wilson et al. 2007), eastness and northness (e.g. Hirzel et al. 2002a; Wilson et al. 2007) and rugosity (e.g. Jenness 2002; Lundblad et al. 2006; Wilson et al. 2007). Still, other abiotic variables could be used, such as curvature (e.g. Wilson et al. 2007), primary productivity (e.g. Smith et al. 2006), organic matter (e.g. Verneaux et al. 2004), salinity (e.g. Al-Hamdani and Reker 2007), temperature (e.g. Connor et al. 2006) and stratification (e.g. Connor et al. 2006). In addition, Guisan and Thuiller (2005); Baptist et al. (2006) and Wilson et al. (2007) stress the importance of spatial scales for predicting the distribution of fauna.

The more abiotic variables become available as input for habitat mapping, the more potential habitats can be classified and potentially new habitats could be identified. However, the relevance of additional classes may not always be clear. It remains important that the variables can be measured or obtained easily and that a sound evaluation of the end products is guaranteed. Another difficulty is the spatial and temporal bias of both biotic and abiotic datasets. Most of the ground-truth data are taken close to the coast and harbours and are strongly biased towards topographic locations. On the BPNS, most samples were taken on the sandbanks, because of their economic potential (e.g. aggregate extraction). Here, the samples are often closely spaced, while other locations are mostly under-sampled. In the most offshore areas, samples are commonly scarce. Apart from the spatial complexity, the samples are also subject to a temporal bias. Gregr and Bodtker (2007) stress the importance of the temporal dynamics (i.e. seasonal variations) for abiotic variables.

In a short time-span, extreme events, such as storms can cause completely different situations of e.g. current regime or suspended matter, causing differences in species composition. Therefore, for the present study, it was decided to work with maximal values of abiotic variables, as those datasets are best suited to represent extreme events (maximum bottom shear stress, maximum current velocity, maximum Chlorophyl a and maximum total suspended matter; cfr. Table 4.1).

On the BPNS, sedimentological samples have been taken from 1976 until now, whilst biological samples are all from a more recent date. Most of the datasets do not cover the same period. Some abiotic datasets are the result of a compilation over many years (e.g. map of $d_s50$ and silt-clay %), whereas others represent a very limited time span (e.g. maximum bottom shear stress; based on data from a spring-neap tidal cycle, 14.8 days). In an ideal situation, all abiotic and biotic datasets would cover the same spatial and temporal scale.

Misleading conclusions can be drawn because of the inappropriate use of some datasets. Sedimentological samples are very suitable to define the sand and to a lesser extent the silt-clay fraction. The gravel fraction (> 2 mm) might be underestimated when grab samples only have been obtained. Gravel can be detected with acoustical classification techniques, but only minor parts of the BPNS have been covered until known (Van Lancker et al. 2007). However, gravel is a part of very interesting habitats with generally high biodiversities (e.g. relation of gravel occurrence with *Ostrea edulis* and *Clupea harengus* (Houziaux et al. 2007a; and Houziaux et al. 2007b); with scallops (Kostylev et al. 2003); and with algae and *Crepidula fornicata* (Brown et al. 2002)).

Therefore, marine landscape mapping 'suggests' only possible ecologically interesting areas, and its predictive power remains dependent on the nature, quality and stability of the abiotic variables.


### 4.4.3  Ecological relevance

The BPNS is an ideal test case for the proposed methodology as both abiotic and biological datasets are widely available. Since the marine landscapes in the present study are rather limited in surface area, they might be considered as habitats and the results might be similar than those that would be obtained with habitat mapping. However, the difference between them is that he Marine Landscape approach is top-down and the habitat mapping approach is bottom-up. This means that for the top-down approach biotic data are used at the end of the process for the validation (or, in

the case of no samples, not at all for some marine landscapes). For the bottom-up approach, abiotic and biotic data are used from the beginning of the process to create a habitat model, centering around the relationships between both (e.g. Willems et al. 2008). Still, for the top-down approach, abiotic data have to be selected that are at least assumed to have an ecological relevance. This knowledge may be derived from literature or expert judgement, but also from a visual inspection at the beginning of the process, comparing possible abiotic input layers with the number of biotic samples. In this paper, no prior selection of abiotic variables has to be done, as all of them are used as PCs.

The ecological validation for this study was based on an indicator species analysis, defining significant indicator species for the predefined clusters. The results showed that for each cluster, except for cluster 3 and 6, significant indicator species could be found.

As such, the clusters are a good proxy for biological predictions. Still, it must be clear that it is not the absolute aim of the marine landscape mapping to predict the biology as such, therefore other and better predictive modelling techniques exist (e.g. Guisan and Zimmermann 2000). Marine landscapes give an indication about the biology, derived solely from abiotic datasets, and offer a valuable in alternative in areas where biological data are scarce or absent.

There seems to be a discrepancy between the number of landscapes (8) and the number of clusters with significant indicator species (6): if a landscape is ecologically meaningful, then this landscape should be populated by specific biota or, in other words, every landscape should be uniquely linked to the biology. Although we might conclude from this discrepancy that several identified marine landscapes have no ecological meaning, we might also explain this by the potential lack of sufficiently detailed information on the marine biota, used to validate the marine landscapes. In conclusion, the level of detail of our current knowledge on the macrobenthos might be insufficient for an unbiased validation of the marine landscapes. If such detailed information would be available, then these data could help to further unravel the ecological meaning of all eight marine landscapes. At the same time, one will never be able to completely explain the occurrence of certain species and communities on the basis of the abiotic environment alone. A biological or an abiotic point of view will never result in exact the same abstractions of the marine seabed, because both approaches are a different way of looking at the same thing.


### 4.4.4   SWOT analysis

A critical evaluation of the protocol to map marine landscapes is performed using a SWOT analysis (strengths, weaknesses, opportunities and threats).
The main *strengths* of the protocol are the following:
  - the possibility to use all available abiotic variables as input for PCA, a technique that eliminates all redundancy of correlating data;
  - the unnecessity to classify the abiotic variables before the clustering and thus the possibility to use continuous abiotic variables as input for the clustering;
  - the use of the C-H criterion to help defining the optimal number of clusters of marine landscapes.
  - the proposed protocol is repeatable and objective; it forms a good alternative for the currently used methodologies which imply subjective decisions to be made.

The main *weaknesses* are:
- the added value of the defined clusters or marine landscapes map is dependent on the availability of relevant abiotic datasets;
- PCA, cluster analysis and INDVAL requires statistical insight and knowledge of the user.

A possible *opportunity* is:
- the application of this protocol for mapping marine landscapes on an international scale and for a larger area than the BPNS (e.g. as contribution to the European Atlas of the Seas in the context of the future European Marine Strategy Framework Directive).

A possible *threat* is:
- for a mapping exercise over a large area, abiotic and biotic datasets, based on different techniques and with different accuracies will be merged, causing an inpredictable error propagation.

Summarizing, the protocol creates interesting opportunities for a mapping exercise on a European scale, but important considerations have to be made about the accuracy of the final result, when datasets of different qualities and origins are used. Foster-Smith et al. (2007b) describe how the accuracy and confidence of marine habitat maps can be assessed, based on a multi-criteria approach.


## 4.5   Conclusion

This paper proposes an objective statistical method for the definition of ecologically relevant marine landscapes. The zones represent well the natural environment and there are clear relationships with the original abiotic variables and the occurrence of macrobenthic species. The methodology is straightforward and allows an easy application to other areas. Marine spatial planning, environmental protection and management of marine zones can benefit from the definition of ecologically relevant marine landscapes (e.g. definition of most important ecological zones, to be protected from dredging and dumping activities or aggregate extraction).


## 4.6   Acknowledgements

database (Marine Biology Section, Ugent – Belgium, 2008) has been compiled by the Marine Biology Section.