



ELSEVIER

Forecasting the probability of finding oil slicks using a CBR system.

Aitor Mata and Juan Manuel Corchado,

Department of Computing Science and Automatic, University of Salamanca, Plaza de la Merced, s/n, Salamanca, Spain

Abstract

A new predicting system is presented in which the aim is to forecast the presence of oil slicks in a certain area of the open sea after an oil spill. Case-Based Reasoning is a computational methodology designed to generate solutions to a certain problem by analysing previous solutions given to previous solved problems. In this case, the system designed to predict the presence of oil slicks wraps other artificial intelligence techniques such as a Radial Basis Function Networks, Growing Cell Structures and Principal Components Analysis in order to develop the different phases of the Case-Based Reasoning cycle. The proposed system uses information such as sea salinity, sea temperature, wind, currents, pressure, number and area of the slicks.... obtained from various satellites. The system has been trained using data obtained after the Prestige oil spill, occurred in the Atlantic waters, in the northwest of Spain. The system developed has been able to accurately predict the presence of oil slicks in the north west of the Galician coast, using historical data.

Keywords: Case-Based Reasoning; oil spill; Growing Cell Structures; Radial Basis Function; Principal Component Analysis.

1. Introduction

Predicting the behaviour of oceanic elements is a quite difficult task. In this case the prediction is related with external elements (oil slicks), what makes the prediction even more difficult. Open ocean is a highly complex system that may be modelled by measuring different variables and structuring them together. Some of those variables are essential to predict the behaviour of oil slicks. In order to predict the future presence of oil slicks in an area, it is obviously necessary to know their previous positions. That knowledge is provided by the analysis of satellite images, obtaining the precise position of the slicks.

The solution proposed in this paper generates, for different geographical areas, a probability (between 0 and 1) of finding oil slicks after an oil spill. The proposed system has been constructed using both historical data and the knowledge generated during the Prestige oil spill, from November 2002 to April 2003. Most of the data used to develop the proposed system has been acquired from the ECCO (*Estimating the Circulation and Climate of the Ocean*) consortium (Menemenlis *et al.*, 2005). Position and size of the slicks has been obtained by treating SAR (*Synthetic Aperture Radar*) satellite images (Palenzuela *et al.*, 2006).

The proposed system is a forecasting Case-Based Reasoning system: the Oil Spill CBR (*OSCBR*). A CBR system has the ability to learn from past situations, and to generate solutions to new problems based in the past solutions given to past problems. Past solutions are stored

in the system, in the *case base*. In *OSCBR* the cases contain information about the oil slicks as long as atmospheric data (wind, salinity, temperature, ocean height and pressure). *OSCBR* combines the efficiency of the CBR systems with artificial intelligence techniques in order to improve the results and to better generalize from past data.

The results obtained with *OSCBR* approximate to the real process occurred in near the ninety per cent of the value of the main variables analyzed, which is a quite important approximation.

In this paper, the CBR technology will be first explained, introducing the specific elements that make this way of predicting work. In second place, the oil spill problem is presented, showing its difficulties and the possibilities of finding solutions to the problem. Finally, *OSCBR* is explained, giving special attention to the techniques applied in the different phases of the CBR cycle. Last, the results are shown and also the future developments that can be achieved with the system.

2. Case-based reasoning systems

Case-Based Reasoning is a technique that has its origin in knowledge based systems. CBR systems learn from previous situations (Aamodt, 1991). The main element of a CBR system is the *case base*; a structure that stores problems, elements (*cases*), and its solutions. So, a case base can be visualized as a database where a collection of problems is stored keeping a relationship with the solutions

to every problem stored, which give the system the ability to generalize in order to solve new problems.

The learning capabilities of the CBR systems are due to its own structure, composed of four main phases (Aamodt and Plaza, 1994): *retrieval*, *reuse*, *revision* and *retention*. The first phase is called *retrieve*, and consists in finding the most similar cases to the proposed problem from the case base. Once a series of cases are extracted from the case base, they must be *reused* by the system. In this second phase, an adaptation of the selected cases is done to fit the current problem. After giving a solution to the problem, that solution is *revised* to check if the proposed alternative is a solution to the problem. If the proposal is confirmed as a solution, then it is *retained* by the system and could eventually serve as a solution to future problems.

Case-Based reasoning is a methodology (Watson, 1999), and so it has been applied to solve different kind of problems. It is a model that can be easily applied to solve soft computing problems (Shiu and Pal, 2004), since the methodology used by CBR is quite easy to assimilate by soft computing approaches. Another interesting application is related with stock market prediction (Chun and Park, 2005), where using different daily values, a CBR system can create a model that may help in stock market investments. Construction is another of the fields of application of CBR, first for the construction of functional databases (Yu and Liu, 2006) to improve the benefits in the usually chaotic organization of the construction projects and also (Chow *et al.*, 2006) to help to choose between different methods and materials, using expert system oriented applications.

Other applications of the CBR methodology cover from health applications (Corchado *et al.*, 2008) to eLearning. CBR has evolved, being transformed so that it can be used to solve new problems, becoming a methodology to plan, or distributed version. Oceanographic problems (Fdez-Riverola and Corchado, 2004), has also been solved with these techniques, helping to predict the value of variable parameters.

But, in most cases, CBR has not been used alone, but combined with various artificial intelligence techniques. Growing Cell Structures has been used with CBR to automatically create the intern structure of the case base from existing data and it has been combined with multi-agent applications (Carrascosa *et al.*, 2007) to improve its results. ART-Kohonen neural networks (Yang *et al.*, 2004), artificial neural networks and fuzzy logic (Fdez-Riverola *et al.*, 2007) has also been used to complement the capabilities of the CBR methodology.

Actual trends in CBR explore the possibility of giving explanations from the very CBR systems (Sørmo *et al.*, 2005). These techniques allow the CBR systems to give the users a better solution, adding extra information to the solution proposed by the system.

3. Oil spill problem

After an oil spill, it is necessary to determine if an area is going to be contaminated or not. To conclude about the presence or not of contamination in an area it is necessary to know how the slicks generated by the spill behave. The most data available; the best solution can be given.

First, position, shape and size of the oil slicks must be identified. The most precise way to acquire that information is by using satellite images. SAR images are the most commonly used to automatically detect this kind of slicks (Solberg *et al.*, 1999). The satellite images show certain areas where it seems to be nothing, like zone with no waves; that are the oil slicks (*figure 1* shows an example of SAR image with oil spills). With these images it is possible to distinguish between normal sea variability and slicks. It is also important to make a distinction between oil slicks and look-alikes. Oil slicks are quite similar to quiet sea areas. If there is not enough wind, the difference between the calmed sea and the surface of a slick is less evident and so, there may be more mistakes when trying to differentiate between an oil slick and something that it is not a slick. This is a crucial aspect in this problem that can also be automatically done by a series of computational tools.

Once the slicks are identified, it is also crucial to know the atmospheric and maritime situation that is affecting the slick in the moment that is being analysed. Information collected from satellites is used to obtain the atmospheric data needed. That is how different variables such as temperature, sea height and salinity are measured in order to obtain a global model that can explain how slicks evolve.

3.1. Previous solutions given to the oil spill problem

There have been different ways to analyze, evaluate and predict situations after an oil spill. One approach is the simulation, where a model of a certain area is created, introducing specific parameters (weather, currents and wind) and working along with a forecasting system. Using this methodology, it is easy to obtain a good solution for a certain area, but it is quite difficult to generalize in order to solve the same problem in new zones. Another way to obtain a trajectory model is to replace the oil spill by drifters comparing the trajectory followed by the drifters with the already known oil slicks trajectories. If the drifters follow a similar trajectory as the one that followed the slicks, then a model can be created and there will be a possibility of creating more models in different areas. A trajectory model has been created to accomplish the NOAA standards, where both the 'best guess' and the 'minimum regret' solutions are generated.

3.2. Models

One step over those solutions previously explained are the systems that, combining a major set of elements, generate response models to solve the oil spill problem.

Table 1
Variables that define a case

Variable	Definition	Unit
Longitude	Geographical longitude	Degree
Latitude	Geographical latitude	Degree
Date	Day, month and year of the analysis	dd/mm/yyyy
Sea Height	Height of the waves in open sea	m
Bottom pressure	Atmospheric pressure in the open sea	Newton/m ²
Salinity	Sea salinity	ppt (parts per thousand)
Temperature	Celsius temperature in the area	°C
Area of the slicks	Surface covered by the slicks present in the analyzed area	Km ²
MeridionalWind	Meridional direction of the wind	m/s
Zonal Wind	Zonal direction of the wind	m/s
Wind Strength	Wind strength	m/s
Meridional Current	Meridional direction of the ocean current	m/s
Zonal Current	Zonal direction of the ocean current	m/s
Current Strength	Ocean current strength	m/s

A different point of view is given by complex systems that analyze large data bases (environmental, ecological, geographical and engineering), using expert systems. This way, an implicit relation between problem and solution is obtained, but with no direct connection between past examples and current decisions. Nevertheless there is a great data mining effort in that kind of solutions.

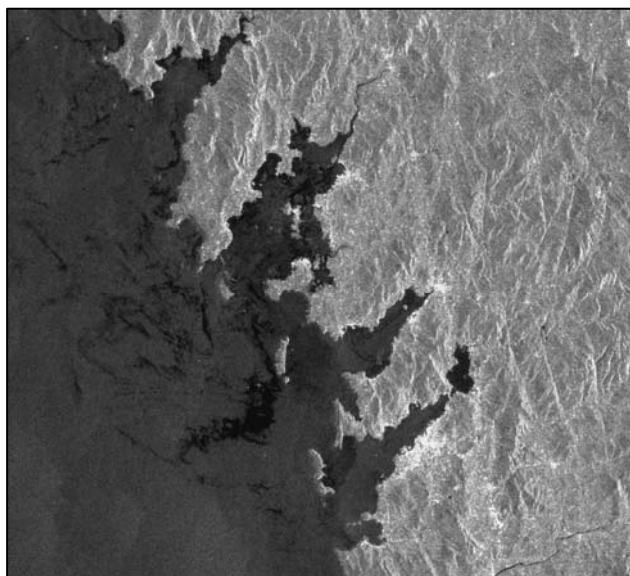


Fig. 1. SAR image showing an oil slick near the Galician coast.

Once the oil spill is produced there should be contingency models that make a fast solution possible (Reed *et al.*, 1999). To get the proper solution expert systems has also been used, using the stored information from past cases, as a repository where future applications will find structured information. Some other complete models have been created, to integrate the different variables affecting the spills, always trying to get better benefits than the possible costs generated by a response to a situation.

The final objective of all these systems is to be decision support systems, in order to help to take all the decisions that need to be taken properly organized. To achieve that great objective, different techniques have been used, from fuzzy logic to negotiation with multi-agent systems (Liu and Wirtz, 2005).

4. Oil spill CBR system - OSCBR

CBR has already been used to solve maritime problems (Corchado and Fdez-Riverola, 2004) in which different oceanic variables were involved. In this case, the data collected from different observations from satellites, is pre-processed, and structured in cases. The created cases are the keys to obtain the solutions to future problems, through the CBR system.

Oil slicks are detected using SAR images. Those images are processed and transformed to be used by the system. In *figure 1* a SAR image is shown. There, a portion of the western Galician coast is shown, as long as some black areas, corresponding to the oil slicks. *Figure 2* shows the interpretation of the previous image after treating the data generated by the SAR image.

OSCBR determines the probability of finding oil slicks

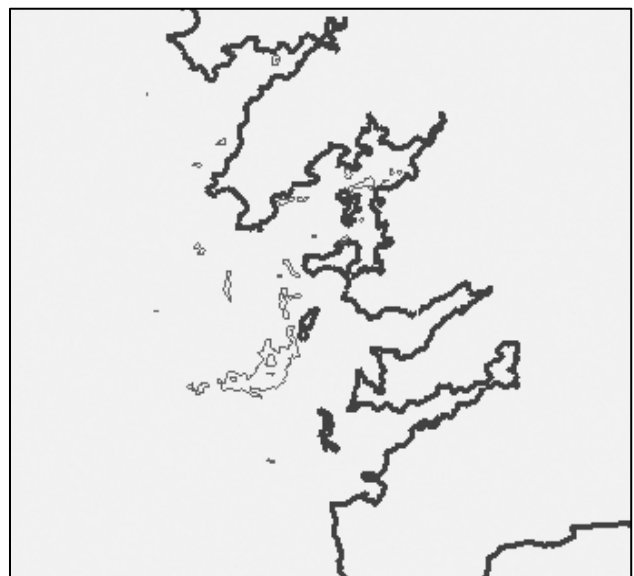


Fig. 2. Oil slicks represented by OSCBR.

in a certain area. To generate the predictions, the system divides the area to be analyzed in squares of approximately half a degree side. The analyzed area and the size of the squares are configurable, been changed to be adapted to the different possible situations. Then the system determines the amount of slicks present in a square.

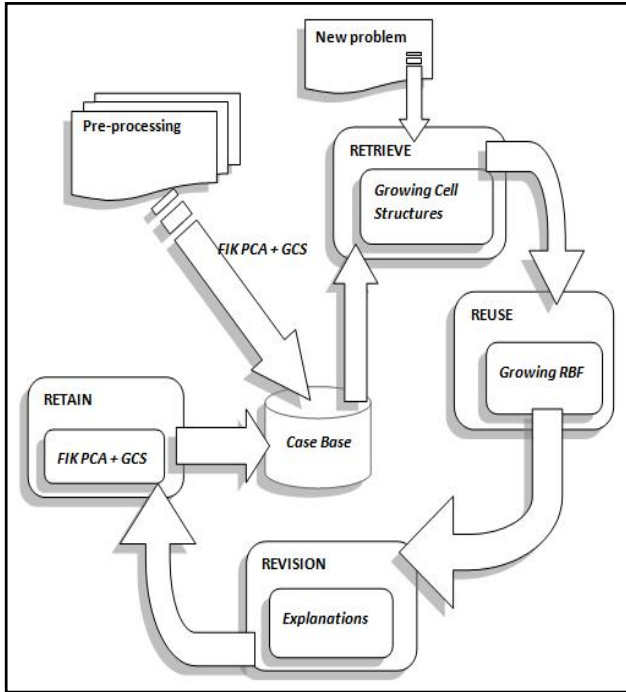


Fig. 3. CBR cycle adapted to the OSCBR system.

A squared zone determines the area that is going to be analyzed independently. The values of the different variables in a square area in a certain moment as long as the is what is called *a case*, which define the problem, the situation that needs to be solved.

In *table 1* the structure of a case is shown. The variables present in a case can be geographical (longitude and latitude), temporal (date of the case), atmospheric (wind, sea height, bottom pressure, salinity and temperature) and variables directly related with the slicks (number and covered area).

Once the data is structured, it is stored in the *case base*. Every case has its temporal situation stored and that relates every case with the next situation in the same position. That temporal relationship is what creates the union between *problem* and *solution*. The problem is the past case, and the solution is the future case, the future state of the square analyzed.

The system has been trained using historical data. The data used to check the system has been obtained after the Prestige accident, between November 2002 and April 2003, in a specific geographical area to the north west of the Galician coast (longitude between 14 and 6 degrees west and latitude between 42 and 46 degrees north). When all that information is stored in the case base, the system is

ready to predict future situations. To generate a prediction, a problem situation must be introduced in the system. Then the most similar cases to the problematic situation are retrieved from the case base. Once a collection of cases are chosen from the case base, they must be used to generate the solution to the current problem.

OSCBR combines the capabilities of the CBR methodology and the power of artificial intelligence techniques. As shown in *figure 3*, every CBR phase uses an artificial intelligence technique in order to obtain its solution. In *figure 3* the four main phases of the CBR cycle are shown as long as the artificial intelligence techniques used in each phase. Those phases with its related techniques are explained next.

Figure 4 shows the graphical user interface of the developed system. In that image the different components of the application can be seen (maps, prediction, slicks, studies...) as well as a visualization of an oceanic area with oil slicks and a squared area to be analyzed.

4.1. Pre-processing

Historical data is used to create the *case base*. As explained before, cases are formed by a series of variables. If the number of those variables is reduced, both the amount of disk space necessary to store the case base and the time to find the most similar cases are reduced. *Principal Components Analysis* (PCA) (Dunteman, 1989) can reduce the number of those variables and then, the system stores the value of the principal components, which are related with the original variables that define a case. PCA has been previously used to analyse oceanographic data and it has proved to be a consistent technique when trying to reduce the number of variables.

In this paper *Fast Iterative Kernel PCA*, an evolution of PCA, has been used (Gunter *et al.*, 2007). This technique reduces the number of variables in a set by eliminating those that are linearly dependent, and it is quite faster than the traditional PCA. To improve the convergence of the Kernel Hebbian Algorithm used by Kernel PCA, *FIK-PCA* set η_i proportional to the reciprocal of the estimated eigenvalues. Let $\lambda_\tau \in \mathfrak{R}_+^r$ denote the vector of eigenvalues associated with the current estimate of the first r eigenvectors. The new Kernel Hebbian Algorithm (KHA) for PCA (Kim *et al.*, 2005) sets de i^{th} component of η_i to:

$$[\eta_t]_i = \frac{1}{|\lambda_{t|} \tau} \eta_0, \quad (1)$$

The final variables are, obviously, linearly independent and are formed by combination of the previous variables. The values of the original variables can be recovered by doing the inverse calculation to the one produced to obtain the new variables. The variables that are less determinant in the final stored variables are those whose values suffer less changes during the periods of time analysed (salinity,

temperature and pressure do not change from one day to another, then, they can be *almost ignored* considering that the final result does not depend on them).

Once applied the *FIK-PCA*, the number of variables is reduced to three, having the following distribution:

$$\begin{aligned} \text{Variable}_1: & -0,560*\text{longitude} - 0,923*\text{latitude} + \\ & 0,991*\text{surface_height} + 0,919*\text{bottom_pressure} + \\ & 0,992*\text{salinity} + 0,990*\text{temperature} - 0,125*\text{area_of_slicks} \\ & + 0,80*\text{meridional_wind} + 0,79*\text{zonal_wind} + \\ & 0,123*\text{wind_strenght} + 0,980*\text{meridional_current} + \\ & 0,980*\text{zonal_current} + 0,980*\text{current_strenght} \end{aligned}$$

$$\begin{aligned} \text{Variable}_2: & 0,292*\text{longitude} - 0,081*\text{latitude} - \\ & 0,010*\text{surface_height} - 0,099*\text{bottom_pressure} - \\ & 0,011*\text{salinity} - 0,013*\text{temperature} - 0,021*\text{area_of_slicks} \\ & + 0,993*\text{meridional_wind} + 0,993*\text{zonal_wind} + \\ & 0,989*\text{wind_strenght} - 0,024*\text{meridional_current} - \\ & 0,024*\text{zonal_current} - 0,024*\text{current_strenght} \end{aligned}$$

$$\begin{aligned} \text{Variable}_3: & 0*\text{longitude} - 0,072*\text{latitude} + \\ & 0,009*\text{surface_height} + 0,009*\text{bottom_pressure} + \end{aligned}$$

$$\begin{aligned} & 0,009*\text{salinity} + 0,009*\text{temperature} + \\ & 0,992*\text{area_of_slicks} + 0,006*\text{meridional_wind} + \\ & 0,005*\text{zonal_wind} + 0,005*\text{wind_strenght} - \\ & 0,007*\text{meridional_current} - 0,007*\text{zonal_current} - \\ & 0,007*\text{current_strenght} \end{aligned}$$

After applying *FIK-PCA*, the historical data is stored in the case base, and is used to solve future problems using the rest of the CBR cycle. Storing the principal components instead of the original variables implies reducing the amount of memory necessary to store the information in about a sixty per cent which is more important as the case base grows. The reduction of the number of variables considered also implies a faster recovery from the case base, about a forty per cent faster than before.

When introducing the data into the case base, *Growing Cell Structures* (Fritzke, 1994) are used. GCS can create a model from a situation organizing the different cases by their similarity. If a 2D representation is chosen to explain this technique, the most similar cells (*cases* in OSCBR) are near one of the other. If there is a relationship between the cells, they are grouped together, and this grouping

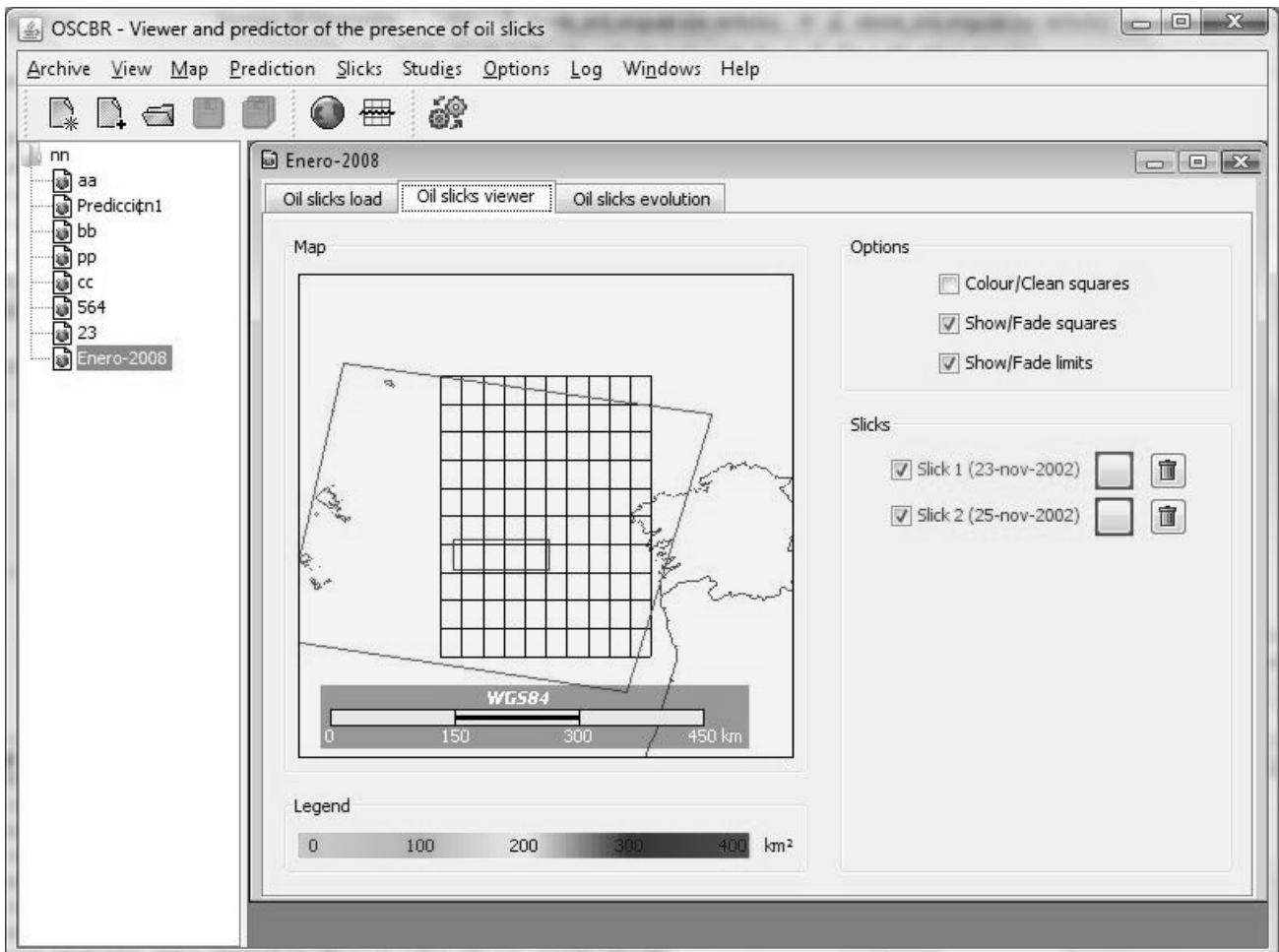


Fig. 4. Graphical user interface of the OSCBR system. The different components of the system can be observed here.

characteristic helps the CBR system to recover the similar cases in the next phase. When a new cell is introduced in the structure, the closest cells move towards the new one, changing the overall structure of the system as shown in (2) and (3). The weights of the winning cell, ω_c , and its neighbours, ω_n , are changed. The changed value is represented by $\omega_c(t+1)$, and $\omega_n(t+1)$ respectively. The terms ε_c and ε_n represent the learning rates for the winner and its neighbours. x represents the value of the input vector.

$$\omega_c(t+1) = \omega_c(t) + \varepsilon_c(x - \omega_c) \quad (2)$$

$$\omega_n(t+1) = \omega_n(t) + \varepsilon_n(x - \omega_n) \quad (3)$$

The pseudocode of the insertion process is shown below:

Growing Cell Structure insertion pseudocode:

1. The most similar cell to the new one is found.
2. The new cell is introduced in the middle of the connection between the most similar cell and the least similar to the new one.
3. Direct neighbours of the closest cell change their values by approximating to the new cell and specified percentage of the distance between them and the new cell.

4.2. Retrieve

Once the case base has stored the historical data, and the GCS has been structured according to the original distribution of the variables, the system is ready to receive a new problem.

When a new problem comes to the system, GCS are used once again. The stored GCS behaves as if the new problem would be stored in the structure, and finds the most similar cells (*cases* in the CBR system) to the problem introduced in the system. In this case the GCS does not change its structure, because it is being used to recover the most similar cases to the introduced problem. Only in the retain phase, the GCS changes again, introducing if it is correct, the proposed solution.

The similarity of the new problem to the stored cases is determined by the GCS calculating the distance between them. Every element in the GCS has a series of values (every value corresponds to one of the principal components created after de FIK-PCA analysis) and then the distance between elements is a multi-dimensional distance, where all the elements are considered to establish the distance between cells.

Then, after obtaining the most similar cases from the case base, they are used in the next phase. The selected cases base will be used to generate an accurate prediction according to the previous solutions related with the introduced problem.

4.3. Reuse

Once the most similar cases to the problem to be solved are recovered from the case base, they are used to generate the solution. The prediction of the future probability of finding oil slicks in an area is generated using an artificial neural network, with a hybrid learning system. An adaptation of *Radial Basis Functions Networks* (Haykin, 1999) are used to obtain that prediction. The chosen cases are used to train the artificial neural network. Radial Basis Function networks have been chosen because of the reduction of the training time comparing with other artificial neural network systems, such as Multilayer Perceptrons. In this case, in every analysis the network is trained, using only the cases selected from the case base, the most similar to the proposed problem.

Growing RBF networks (Ros *et al.*, 2007) are used to obtain the predicted future values corresponding to the proposed problem. This adaptation of the RBF networks allows the system to grow during training gradually increasing the number of elements (prototypes) which play the role of the centres of the radial basis functions. In this case the creation of the Growing RBF must be made automatically, which implies an adaptation of the original GRBF system. The pseudocode of the growing process and the definition of the error for every pattern is shown below:

$$e_i = l/p * \sum_{k=1}^p ||t_{ik} - y_{ik}||, \quad (4)$$

Where t_{ik} is the desired value of the k^{th} output unit of the i^{th} training pattern, y_{ik} the actual values of the k^{th} output unit of the i^{th} training pattern.

Growing RBF pseudocode:

1. Calculate the error, e_i (4) for every new possible prototype.
 - a. If the new candidate does not belong to the chosen ones and the error calculated is less than a threshold error, then the new candidate is added to the set of accepted prototypes.
 - b. If the new candidate belongs to the accepted ones and the error is less than the threshold error, then modify the weights of the neurons in order to adapt them to the new situation.
2. Select the best prototypes from the candidates
 - a. If there are valid candidates, create a new cell centred on it.
 - b. Else, increase the iteration factor. If the iteration factor comes to the 10% of the training population, freeze the process.
3. Calculate global error and update the weights.
 - a. If the results are satisfactory, end the process. If not, go back to step 1.

Once the GRBF network is created, it is used to generate the solution to the proposed problem. The solution will be the output of the network using as input data the introduced problem.

4.4. Revise

After generating the prediction, it is shown to the user in a similar way the slicks are interpreted by OSCBR. A set of squared coloured areas appear. The intensity of the colour corresponds with the possibility of finding oil slicks in that area. The areas coloured with a higher intensity are those with the highest probability of finding oil slicks in them.

In this visual approximation, the user can check if the solution is a good one or not. But the system provides an automatic method of revision that must be, anyway, checked by an expert user.

Explanations (Sørmo *et al.*, 2005) are used to check the correction of the proposed solution, to justify the solution. To obtain a justification to the given solution, the cases selected from the case base are used once again. To create an *explanation*, a comparison between different possibilities has been used. All the selected cases has its own *future situation* associated. If we consider the case and its solution as two vectors, we can establish a *distance* between them, calculating the evolution of the situation in the considered conditions. If the distance between the proposed problem and the solution given is not bigger than the distances obtained from the selected cases, then the solution is a good one, according to the structure of the case base.

Explanation pseudocode:

1. For every selected case in the retrieval phase, the distance between the case and its solution is calculated.
2. The distance between the proposed problem and the proposed solution is also calculated.
3. If the difference between the distance of the proposed solution and those of the selected cases is underneath certain threshold value, then the solution is considered as a valid one.
4. If not, the user is informed and the process goes back to the retrieval phase, where new cases are selected from the case base.

The distances are calculated considering the sign of the values, not using its absolute value. This decision is easily justified by the fact that is not the same to move to the north than to the south, even if the distance between two points is the same.

If the prediction is considered as correct it will be stored in the case base, and it can then be used in next predictions to obtain new solutions.

4.5. Retain

When the proposed prediction is accepted, it is considered as a good solution to the problem. Then, the solution can be stored in the case base in order to serve to solve new problems. It will be used in future situations as the historical data previously stored in the system.

When inserting a new case in the case base, FIK-PCA is used again to reduce the number of variables used and to adapt the data generated by the system. The adaptation is done by changing the original variables into the principal components previously chosen by the system.

Obviously, when introducing a new case in the case base, the GCS formed by the information stored in the case base, also change, to adapt to the new situation generated. When adapting to the new solution introduced in the case base, the GCS system grows and improves its capability of generating good results as new knowledge have been introduced in the system.

5. Results

The historical data used to train the system has been obtained from different satellites. Temperature, salinity, bottom pressure, sea height, wind, currents, number and area of the slicks, as long as the location of the squared area and the date have been used to create a case. All these data define the problem case and also the solution case. The solution to a problem defined by an area and its variables is the same area, but with the values of the variables changed to the prediction obtained from the CBR system.

When the OSCBR system has been used with a subset of the data that has not been previously used to train the system, it has produced quite hopeful results. The predicted situation was contrasted with the actual future situation. The future situation was known, as long as past data was used to train the system and also to test its correction. The proposed solution was, in most of the variables, close to 90% of accuracy.

Table 2
Percentage of good predictions obtained with different techniques.

<i>Number of cases</i>	RBF	CBR	RBF + CBR	OSCBR
100	45 %	39 %	42 %	43 %
500	48 %	43 %	46 %	46 %
1000	51 %	47 %	58 %	64 %
2000	56 %	55 %	65 %	72 %
3000	59 %	58 %	68 %	81 %
4000	60 %	63 %	69 %	84 %
5000	63 %	64 %	72 %	87 %

For every problem, defined by an area and its variables, the system offers nine solutions: the same area, with its proposed variables and the eight closest neighbours. This way of prediction is used in order to clearly observe the direction of the slicks, what can be useful in order to determine the coastal areas that will be affected by the slicks generated after an oil spill.

In *table 2* a summary of the results obtained is shown. In this table different techniques are compared. The table shows the evolution of the results along with the increase of the number of cases stored in the case base. All the techniques analyzed improve its results when increasing the number of cases stored. If the number of cases of the case base increases, then it is easier to find similar cases to the proposed problem and then, the solution can be more accurate. The “*RBF*” column represents a simple Radial Basis Function Network that is trained with all the data available. The network gives an output that is considered a solution to the problem. The “*CBR*” column represents a pure CBR system, with no artificial intelligence techniques included. The cases are stored in the case bases and recovered considering the Euclidean distance. The most similar cases are selected and after applying a weighted mean depending on the similarity, a solution is proposed. It is a *mathematical* CBR. The “*RBF + CBR*” column corresponds to the possibility of using a RBF system combined with CBR. The recovery from the CBR is done using the Manhattan distance to determine the closest cases to the introduced problem. The RBF network works in the reuse phase, adapting the selected cases to obtain the new solution. The results of the “*RBF+CBR*” column are, normally, better than those of the “*CBR*”, mainly because of the elimination of useless data to generate the solution. Finally, the “*OSCBR*” column shows the results obtained by the proposed system, being better than the three previous solutions analyzed.

Table 3
Multiple comparison procedure among different techniques.

	RBF	CBR	RBF + CBR	OSCBR
<i>RBF</i>				
<i>CBR</i>	*			
<i>RBF + CBR</i>	=	=		
<i>OSCBR</i>	*	*	*	

Table 3 shows a multiple comparison procedure (*Mann-Whitney* test) used to determine which models are significantly different from the others.

The asterisk indicates that these pairs show statistically significant differences at the 99.0% confidence level. It can be seen in *Table 3*, that the *OSCBR* system presents statistically significant differences with the rest of the models.

The proposed solution does not generate a trajectory, but a series of probabilities in different areas, what is far more similar to the real behaviour of the oil slicks.

Once the prediction is generated and approved, it can be exported to various formats. First an html file can be generated with the images that represent the prediction, the solution to the problem. The other two output formats are “*Google related*”: the solutions can be exported to *Google Earth* and to *Google Maps*. *Google Maps* needs that the file containing the exported information were in an external server, while *Google Earth* can be used in a local machine.

6. Conclusions and future work

In this paper, the *OSCBR* system has been described. It is a new solution for predicting the presence or not of oil slicks in a certain area after an oil spill.

This system used data acquired from different orbital satellites and with that data the CBR environment was created. The data must be previously classified into the structure required by the CBR system to store it as a case.

OSCBR uses different artificial intelligence techniques in order to obtain a correct prediction. *Fast Iterative Kernel Principal Component Analysis* is used to reduce the number of variables stored in the system, getting about a 60% of reduction in the size of the *case base*. This adaptation of the PCA also implies a faster recovery of cases from the case base (more than 40% faster than storing the original variables).

To obtain a prediction using the cases recovered from the case base, *Growing Radial Basis Function Networks* has been used. This evolution of the RBF networks implies a better adaptation to the structure of the case base, which is organised using *Growing Cell Structures*. The results using *Growing RBF* networks instead of simple RBF networks are about a 6% more accurate, which is a good improvement.

It has been proved that the system can predict in the conditions already known, showing better results than previously used techniques. The use of a combination of techniques integrated in the CBR structure makes possible to obtain better result than using the CBR alone (17% better), and also better than using the techniques isolated, without the integration feature produced by the CBR (11% better).

The next step is generalising the learning, acquiring new data to create a base of cases big enough to have solutions for every season. Another improvement is to create an on-line system that can store the case base in a server and generate the solutions dynamically to different requests. This on-line version will include real time connection to data servers providing weather information of the current situations in order to predict real future situations.

References

- Aamodt, A. (1991) A Knowledge-Intensive, Integrated Approach to Problem Solving and Sustained Learning. *Knowledge Engineering and Image Processing Group. University of Trondheim*.
- Aamodt, A. & Plaza, E. (1994) Case-Based Reasoning: Foundational Issues, Methodological Variations, and System Approaches. *AI Communications*, 7, 39-59.
- Carrascosa, C., Bajo, J., Julian, V., Corchado, J. M. & Botti, V. (2007) Hybrid multi-agent architecture as a real-time problem-solving model. *Expert Systems With Applications*, 34, 2-17.
- Corchado, J. M., Bajo, J. & Abraham, A. (2008) GERAmI: Improving the delivery of health care. *IEEE Intelligent Systems. Special Issue on Ambient Intelligence*.
- Corchado, J. M. & Fdez-Riverola, F. (2004) FSfRT: Forecasting System for Red Tides. *Applied Intelligence*, 21, 251-264.
- Chow, H. K. H., Choy, K. L., Lee, W. B. & Lau, K. C. (2006) Design of a RFID case-based resource management system for warehouse operations. *Expert Systems With Applications*, 30, 561-576.
- Chun, S. H. & Park, Y. J. (2005) Dynamic adaptive ensemble case-based reasoning: application to stock market prediction. *Expert Systems With Applications*, 28, 435-443.
- Dunteman, G. H. (1989) *Principal Components Analysis*, Newbury Park, California.
- Fdez-Riverola, F. & Corchado, J. M. (2004) FSfRT: Forecasting System for Red Tides. *Applied Intelligence*, 21, 251-264.
- Fdez-Riverola, F., Iglesias, E. L., Díaz, F., Méndez, J. R. & Corchado, J. M. (2007) Applying lazy learning algorithms to tackle concept drift in spam filtering. *Expert Systems With Applications*, 33, 36-48.
- Fritzke, B. (1994) Growing cell structures—a self-organizing network for unsupervised and supervised learning. *Neural Networks*, 7, 1441-1460.
- Gunter, S., Schraudolph, N. N. & Vishwanathan, S. V. N. (2007) Fast Iterative Kernel Principal Component Analysis. *Journal of Machine Learning Research*, 8, 1893-1918.
- Haykin, S. (1999) *Neural networks*, Prentice Hall Upper Saddle River, NJ.
- Kim, K. I., Franz, M. O. & Schölkopf, B. (2005) Iterative kernel principal component analysis for image modeling. *IEEE Trans. Pattern Analysis and Machine Intelligence*.
- Liu, X. & Wirtz, K. W. (2005) Sequential negotiation in multiagent systems for oil spill response decision-making. *Marine Pollution Bulletin*, 50, 469-74.
- Menemenlis, D., Hill, C., Adcroft, A., Campin, J. M., Cheng, B., Ciotti, B., Fukumori, I., Heimbach, P., Henze, C. & Köhl, A. (2005) NASA Supercomputer Improves Prospects for Ocean Climate Research. *EOS Transactions*, 86, 89-95.
- Palenzuela, J. M. T., Vilas, L. G. & Cuadrado, M. S. (2006) Use of ASAR images to study the evolution of the Prestige oil spill off the Galician coast. *International Journal of Remote Sensing*, 27, 1931-1950.
- Reed, M., Ekrol, N., Rye, H. & Turner, L. (1999) Oil Spill Contingency and Response (OSCAR) Analysis in Support of Environmental Impact Assessment Offshore Namibia. *Spill Science and Technology Bulletin*, 5, 29-38.
- Ros, F., Pintore, M. & Chrétien, J. R. (2007) Automatic design of growing radial basis function neural networks based on neighborhood concepts. *Chemometrics and Intelligent Laboratory Systems*, 87, 231-240.
- Shiu, S. C. K. & Pal, S. K. (2004) Case-Based Reasoning: Concepts, Features and Soft Computing. *Applied Intelligence*, 21, 233-238.
- Solberg, A. H. S., Storvik, G., Solberg, R. & Volden, E. (1999) Automatic detection of oil spills in ERS SAR images. *IEEE Transactions on Geoscience and Remote Sensing*, 37, 1916-1924.
- Sørmo, F., Cassens, J. & Aamodt, A. (2005) Explanation in Case-Based Reasoning—Perspectives and Goals. *Artificial Intelligence Review*, 24, 109-143.
- Watson, I. (1999) Case-based reasoning is a methodology not a technology. *Knowledge-Based Systems*, 12, 303-308.
- Yang, B. S., Han, T. & Kim, Y. S. (2004) Integration of ART-Kohonen neural network and case-based reasoning for intelligent fault diagnosis. *Expert Systems With Applications*, 26, 387-395.
- Yu, W. & Liu, Y. (2006) Hybridization of CBR and numeric soft computing techniques for mining of scarce construction databases. *Automation in Construction*, 15, 33-46.