

A new CBR approach to the oil spill problem.

Juan Manuel Corchado, Aitor Mata, Francisco De Paz And David Del Pozo

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

Abstract. Oil spills represent one of the most destructing environmental disasters. Predicting the possibility of finding oil slicks in a certain area after an oil spill can be crucial in order to reduce the environmental risks. The system presented here forecasts the presence or not of oil slicks in a certain area of the open sea after an oil spill using Case-Based Reasoning methodology. CBR is a computational methodology designed to generate solutions to a certain problem by analysing previous solutions given to previous solved problems. The proposed system 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 CBR cycle. CBR systems have never been used before to solve oil slicks problems. The proposed system uses information obtained from various satellites such as salinity, temperature, pressure, number and area of the slicks.... OSCBR system has been able to accurately predict the presence of oil slicks in the north west of the Galician coast, using historical data.

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 historical data and checked using the data acquired 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. Position and size of the slicks has been obtained by treating SAR satellite images [1].

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).

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. 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 [2]: *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.

CBR has been applied to different situations like treating problems with no evident solutions. But, in most cases, CBR has not been used alone, but combined with various artificial intelligence techniques. Growing Cell Structures [3] has been used with CBR to automatically create the intern structure of the case base from existing data. Actual trends in CBR explore the possibility of giving explanations from the very CBR systems [4]. 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 [5]. The satellite images show certain areas where it seems to be nothing, like zone with no waves; that are the oil slicks. With these images it is possible to distinguish between normal sea variability and slicks. It is also important to distinguish 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 distinguish 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 essential 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 [6], 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 [7] 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 different way of predicting oil slicks trajectories is to study previous cases to obtain a trajectory model for a certain area with different weather situations [8].

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.

A quite new point of view is given by complex systems [9] 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.

Once the oil spill is produced there should be contingency models that make a fast solution possible. 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.

The final objective of all these systems is to be decision support systems [10], 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.

4 OIL SPILL CBR SYSTEM – OSCBR

CBR has already been used to solve maritime problems [11] 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.

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
Meridional Wind	Meridional direction of the wind	m/s
Zonal Wind	Zonal direction of the wind	m/s
Wind Strenght	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 Strenght	Ocean current strength	m/s

Table 1. Variables that define a case.

OSCBR determines the probability of finding oil slicks in a certain area. To generate the predictions, the system divides the ocean surface in squares of approximately half a degree side. Then the system determines the amount of slicks present in a square. The squares where the slicks are located are coloured with different gradation depending on the quantity of the squared area covered by oil slicks.

The 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 value of the possibility of finding oil slicks in the following day is what is called a *case*, which define the problem and propose the solution.

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 problem (number and area of the slicks).

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 data used to train 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. *Radial Basis Functions* Networks are used to combine the chosen cases in order to obtain the new solution.

OSCBR includes different artificial intelligence techniques to achieve the objectives of every CBR phase. As shown in *figure 1*, every CBR phase uses an artificial intelligence technique in order to obtain its solution. Those phases with its related techniques are going to be explained next.

4.1 Pre-processing

Historical data collected from November 2002 to april 2003 is used to create the *case base*. As explained before, cases are formed by a series of variables. *Principal Components Analysis* (PCA) 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.

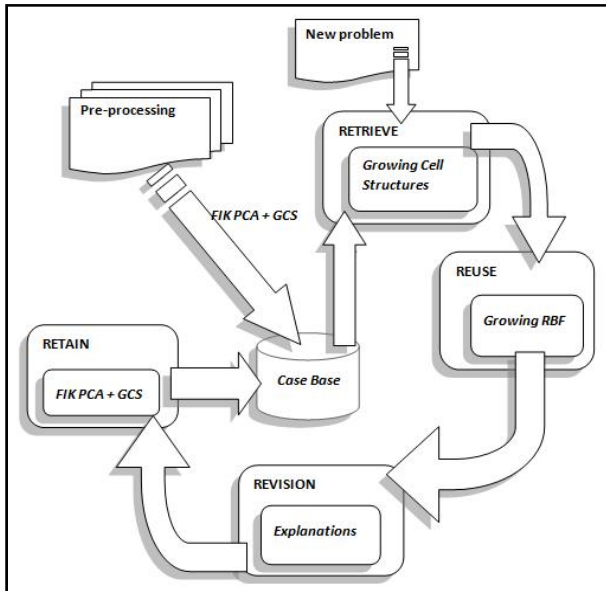


Figure 1. CBR cycle adapted to the OSCBR system.

In this paper *Fast Iterative Kernel PCA*, an evolution of PCA, has been used [12]. 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_t proportional to the reciprocal of the estimated eigenvalues. Let $\lambda_t \in \mathfrak{R}_+^r$ denote the vector of eigenvalues associated with the current estimate of the first r eigenvectors. The new KHA algorithm sets de i^{th} component of η_t to:

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

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 used 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 *ignored* considering that the final result does not depend on them).

After applying FIKPCA, 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 forty 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.

When introducing the data into the case base, *Growing Cell Structures* [13] 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 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. The weights of the winning cell, ω_c , and its neighbours, ω_n , are changed. The terms ε_c and ε_n represent the learning rates for the winner and its neighbours, respectively. x represents the value of the input vector.

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

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

4.2 Retrieve

Once the case base has stored the historical data, and the GCS has learned from 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 once used. 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 obtain 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 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 most similar cases stored in the case base will be used to obtain an accurate prediction according to the previous solutions related with the selected cases.

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* are used to obtain that prediction [14; 15]. 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 [16] 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 centers 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 [17]. 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.

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 selected cases from the case base.

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 are used to check the correction of the proposed solution, to justify the solution [18]. 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.

5 RESULTS

The data used to train the system has been obtained from different satellites. Temperature, salinity, bottom pressure, sea height, 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.

Number of cases	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 %

Table 2. Percentage of good predictions obtained with different techniques.

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 the correction of it. The proposed solution was, in most of the variables, close to 90% of accuracy.

In 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. 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.

The proposed solution do 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.

6 CONCLUSIONS AND FUTURE WORK

In this paper, the OSCBR system has been explained. 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 40% 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 7% 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 4% 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

- [1] Palenzuela, J.M.T., Vilas, L.G. and 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 (10), 1931-1950.
- [2] Aamodt, A. and Plaza, E. (1994) Case-Based Reasoning: Foundational Issues, Methodological Variations, and System Approaches, *AI Communications*, 7 (1), 39-59.
- [3] Azuaje, F., Dubitzky, W., Black, N. and Adamson, K. (2000) Discovering relevance knowledge in data: a growing cell structures approach, *IEEE Transactions on Systems, Man and Cybernetics, Part B*, 30 (3), 448-460.
- [4] Roth-Berghofer, T.R., Cassens, J. and Sørmo, F. (2005) Goals and Kinds of Explanations in Case-Based Reasoning, *Proceedings of WM*, 264-268.
- [5] Solberg, A.H.S., Storvik, G., Solberg, R. and Volden, E. (1999) Automatic detection of oil spills in ERS SAR images, *IEEE Transactions on Geoscience and Remote Sensing*, 37 (4), 1916-1924.
- [6] Brovchenko, I., Kuschan, A., Maderich, V. and Zheleznyak, M. (2002) The modelling system for simulation of the oil spills in the Black Sea, *3rd EuroGOOS Conference: Building the European capacity in operational oceanography*, 192.
- [7] Price, J.M., Ji, Z.G., Reed, M., Marshall, C.F., et al. (2003) Evaluation of an oil spill trajectory model using satellite-tracked, oil-spill-simulating drifters, *OCEANS 2003. Proceedings*, 3.
- [8] Vethamony, P., Sudheesh, K., Babu, M.T., Jayakumar, S., et al. (2007) Trajectory of an oil spill off Goa, eastern Arabian Sea: Field observations and simulations, *Environmental Pollution*.
- [9] Douligieris, C., Collins, J., Iakovou, E., Sun, P., et al. (1995) Development of OSIMS: An oil spill information management system, *Spill Science & Technology Bulletin*, 2 (4), 255-263.
- [10] Keramitsoglou, I., Cartalis, C. and Kassomenos, P. (2003) Decision Support System for Managing Oil Spill Events, *Environmental Management*, 32 (2), 290-298.
- [11] Corchado, J.M. and Fdez-Riverola, F. (2004) FSfRT: Forecasting System for Red Tides, *Applied Intelligence*, 21, 251-264.
- [12] Gunter, S., Schraudolph, N.N. and Vishwanathan, S.V.N. (2007) Fast Iterative Kernel Principal Component Analysis, *Journal of Machine Learning Research*, 8, 1893-1918.
- [13] Fritzke, B. (1994) Growing cell structures—a self-organizing network for unsupervised and supervised learning, *Neural Networks*, 7 (9), 1441-1460.
- [14] Martin, B. and Sanz, A. (1997) Redes neuronales y sistemas borrosos, *Zaragoza: Editorial Ra-Ma*.
- [15] Haykin, S. (1999) *Neural networks*. Prentice Hall Upper Saddle River, NJ.
- [16] Karayiannis, N.B. and Mi, G.W. (1997) Growing radial basis neural networks: merging supervised and unsupervised learning with network growth techniques, *Neural Networks, IEEE Transactions on*, 8 (6), 1492-1506.
- [17] Ros, F., Pintore, M. and Chrétien, J.R. (2007) Automatic design of growing radial basis function neural networks based on neighborhood concepts, *Chemometrics and Intelligent Laboratory Systems*, 87 (2), 231-240.
- [18] Sørmo, F., Cassens, J. and Aamodt, A. (2005) Explanation in Case-Based Reasoning—Perspectives and Goals, *Artificial Intelligence Review*, 24 (2), 109-143.