



# Confidence in Signal Reconstruction by the Evolving Clustering Method

Wei Zhao, Piero Baraldi, Enrico Zio

► **To cite this version:**

Wei Zhao, Piero Baraldi, Enrico Zio. Confidence in Signal Reconstruction by the Evolving Clustering Method. PHM-2011, May 2011, Shenzhen, China. pp.199-205, 2011, .

**HAL Id: hal-00658072**

**<https://hal-supelec.archives-ouvertes.fr/hal-00658072>**

Submitted on 12 Jan 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Confidence in Signal Reconstruction by the Evolving Clustering Method

Enrico Zio

1. Ecole Centrale Paris- Supelec  
Paris, France  
[enrico.zio@ecp.fr](mailto:enrico.zio@ecp.fr)  
[enrico.zio@supelec.fr](mailto:enrico.zio@supelec.fr)
2. Dipartimento di Energia  
Politecnico di Milano  
Milano, Italy  
[enrico.zio@polimi.it](mailto:enrico.zio@polimi.it)

Piero Baraldi

Dipartimento di Energia  
Politecnico di Milano  
Milano, Italy  
[piero.baraldi@polimi.it](mailto:piero.baraldi@polimi.it)

Wei Zhao

Group 203, School of Electronics and Information Engineering  
Beihang University  
Beijing, P.R.China  
[zhaowei203@buaa.edu.cn](mailto:zhaowei203@buaa.edu.cn)

**Abstract**—Monitoring the health conditions of equipment allows supplying advanced warning of their incipient failures; this can provide evidence useful to maintenance and replacement practices. However, uncertainties in the signal measurements and incompleteness in the representativeness of the measured data can overshadow the conclusions drawn from condition monitoring, and possibly lead the decision-maker to take wrong actions. In order to reduce the risk of wrong actions, confidence measures on the condition monitoring indications of the state of a component must be provided, so that the decision-maker can know to what degree he or she should trust such indications. As condition monitoring is usually structured in two modules performed in succession, one of reconstruction of the signal values in normal operating conditions and a following one of equipment health state diagnosis, it is reasonable to define confidence measures for the two processes individually, and then integrate the two into a single criterion for the whole condition monitoring. The research presented in this paper focuses on the definition of confidence measures for the signal reconstruction part of condition monitoring. The Evolving Clustering Method (ECM) is adopted to build the empirical model of signal reconstruction. Requirements for the reconstruction confidence are originally defined, and a single confidence measure is proposed to meet all the requirements identified. The confidence measure is analyzed with respect to two-dimensional artificial datasets and a real dataset concerning the Reactor Coolant Pump of a French Pressurized Water Reactor. The results obtained show that the proposed confidence measure meets all requirements and is more informative than the reconstruction error.

**Keywords**—confidence measure; signal reconstruction; condition monitoring; evolving clustering method

## I. INTRODUCTION

Systems for condition monitoring of the health state of an equipment are often based on empirical models of signal regression whose performances may vary depending on the density and information content of the example signal patterns available to train the models [1-7]. Furthermore, the stochastic behavior of the processes and the signal measurements errors can overshadow the health-state conclusions drawn from condition monitoring, and possibly lead the decision-maker to take wrong actions. Given the criticality of these actions, it is important to provide the decision-maker with a measure of confidence on the condition monitoring system outcome [4].

In this respect, the confidence measure should be useful to recognize the two potentially dangerous cases of missing and false alarms. In the first case, the condition monitoring indicates that an unhealthy equipment is operating in normal conditions; this may cause an unexpected failure of the equipment with associated long downtimes, high costs and possible safety and environmental problems. Contrarily, in case of false alarms, a healthy state of the equipment is erroneously recognized as unhealthy; this may lead to an unnecessary stop of the equipment operation with the associated loss of production, and may subtract resources to deal with other actual failures.

A typical scheme of condition monitoring can be described as follows. Historical data of equipment condition under normal operation are used to build a model (often empirical). The model is auto-associative in that it reconstructs the values of the signals measured to characterize the equipment condition. During operation, the actually observed signal values are compared with those reconstructed by the auto-

associative model and the difference is computed. Based on this deviation, a decision about the health state of the equipment is made. The above two actions are often named fault detection [9,10].

Because the condition monitoring is composed by the two phases of signal reconstruction and health state diagnosis, it is reasonable to define the confidence measures for these two phases individually and then to integrate them into a single confidence measure on the whole condition monitoring process.

In this paper, the research is focused on the development of a confidence measure for the signal reconstruction phase. The Evolving Clustering Method (ECM) [3] is adopted to build the empirical reconstruction model of the equipment behavior. On the basis of ECM, the special requirements for the reconstruction confidence are analyzed and a confidence measure is proposed to meet all the requirements.

The remaining parts of the paper are structured as follows. Section II describes the process of condition monitoring in detail and points out the uncertainties existing in the process. Section III presents the ECM algorithm including the training phase and the procedure used for the signal reconstruction. In Section IV, the requirements for a confidence measure of signal reconstruction are discussed and a corresponding measure able to meet the requirement is proposed. In Section V, the proposed confidence measure is applied to both two-dimensional artificial datasets and a real dataset concerning the condition monitoring of a Reactor Coolant Pump of a French Pressurized Water Reactor. Finally, Section VI states some conclusions and draws on potential future steps of the work.

## II. CONDITION MONITORING

Fig. 1 shows the general framework used for condition monitoring of an equipment. This is typically based on:

- 1) a signal reconstruction module,
- 2) a diagnostic decision module.

The former receives in input the vector containing the  $q$  sensor measurements  $\vec{x}^{obs} = (x^{obs}(1), \dots, x^{obs}(q))$  and provides in output the signal values expected in case of normal condition  $\vec{\hat{x}}_{nc} = (\hat{x}_{nc}(1), \dots, \hat{x}_{nc}(q))$ . This module is usually based on an auto-associative model of the component behavior in normal conditions, obtained by techniques such as Principal Component Analysis (PCA) [10], Auto-Associative Kernel Regression Method (AAKR) [1], Auto-associative Neural Networks [11], Evolving Cluster Method [3], trained with data collected during operation in normal conditions.

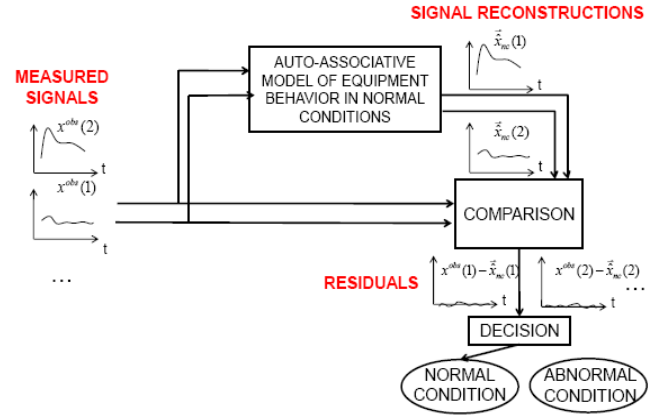


Figure 1 The condition monitoring approach in this work.

The latter module takes the difference between the reconstructed  $(\hat{x}_{nc}(1), \dots, \hat{x}_{nc}(q))$  and observed values  $(x^{obs}(1), \dots, x^{obs}(q))$  to decide whether the system is in normal or abnormal conditions. In case of normal conditions, the measured values are expected to be very similar to the model reconstructed ones, i.e., the residuals are small; on the contrary, under abnormal conditions the measurements tend to deviate from the reconstruction allowing detection of the abnormality.

However, incompleteness of the training data, intrinsic stochasticity of the plant processes and measurement noises may lead to wrong diagnostic decisions on the system health state. For this reason, it is important to develop a confidence measure on the condition monitoring indications of the state of the equipment, so that the decision-maker can know to what degree he or she should trust such indications to take actions.

Proceeding systematically through the two successive phases of condition monitoring, the overall confidence on the condition monitoring,  $Conf_{CM}$ , is sought as a result of the confidence in the signal reconstruction,  $Conf_r$ , and in the diagnostic decision,  $Conf_d$ , properly aggregated.

## III. SIGNAL RECONSTRUCTION BY EVOLVING CLUSTERING METHOD (ECM)

The algorithm considered in this work for reconstructing the equipment behavior in normal conditions is based on a clustering method called Evolving Clustering Method (ECM) [3]. For completeness of the paper, Section III.A briefly describes the basics of ECM, whereas Section III.B is dedicated to the reconstruction procedure based on ECM.

### A. The ECM algorithm

Given a training dataset  $\mathbf{X}^{obs-nc}$  formed by  $n$   $q$ -dimensional patterns  $\vec{x}_i^{obs-nc} = (x_i^{obs-nc}(1), \dots, x_i^{obs-nc}(q))$  recorded during past operation of the equipment in normal conditions, the ECM algorithm provides a procedure to group the training patterns into clusters. The application of the algorithm requires to a priori fix the value of the parameter  $D_{thr}$  which defines the maximum allowed cluster radius. The clusters are then found by performing the following steps:

*Step 0:* Assume the center of the first cluster  $\vec{v}_1$  equal to the first pattern in the dataset  $\vec{x}_1^{obs-nc}$ , the corresponding cluster

radius  $R_1=0$ , and set the pattern counter  $i=1$  and the cluster counter  $m=1$ .

*Step 1:* If all the patterns of the dataset have been processed ( $i=n$ ), exit. Otherwise,  $i=i+1$ , consider the  $i$ -th pattern of the dataset  $\vec{x}_i^{obs-nc}$  and compute its distance  $D_{ik}$  from the centers  $\vec{v}_k$ ,  $k=1, \dots, m$  of all the already formed clusters:

$$D_{ik}^2 = \left\| \vec{x}_i^{obs-nc} - \vec{v}_k \right\|^2, \quad k=1, \dots, m \quad (1)$$

Although the distance computation can be based on weighted distances

$$\left\| \vec{x}_i^{obs-nc} - \vec{v}_k \right\|^2 = \sum_{j=1}^q p(j) * (x_i(j) - v_k(j))^2 \quad (2)$$

with the weights  $p(j)$  related to the importance of signal  $j$  in the equipment monitoring, in this work the Euclidean distance is considered, i.e. all the signals are given equal weights

$$\left\| \vec{x}_i^{obs-nc} - \vec{v}_k \right\|^2 = \sum_{j=1}^q (x_i(j) - v_k(j))^2 \quad (3)$$

*Step 2:* If, among the found distance values  $D_{ik}$ , there are distances lower or equal to the corresponding cluster radii, i.e.

$D_{ik} \leq R_k$ , assign the  $i$ -th pattern to the cluster  $k_0$  with the associated lowest distance  $D_{ik}$ :

$$k_0 = \arg \min_k \left\| \vec{x}_i^{obs-nc} - \vec{v}_k \right\|^2, \quad k=1, \dots, m \quad (4)$$

and go to step 1. Otherwise, if the  $i$ -th pattern does not belong to any cluster, perform steps 3, 4 and 5.

*Step 3:* For all  $m$  existing cluster centers, compute the values  $S_{ik} = D_{ik} + R_k$ ,  $k=1, \dots, m$  and identify the cluster  $k_a$  characterized by the minimum value of  $S_{ik}$ :

$$k_a = \arg \min_k S_{ik} \quad k=1, \dots, m \quad (5)$$

*Step 4:* If  $S_{ik_a}$  is greater than  $2D_{thr}$ , a new cluster  $m=m+1$  is created with center  $\vec{v}_m = \vec{x}_i^{obs-nc}$  and radius  $R_m=0$ ; then, go to Step 1.

*Step 5:* If  $S_{ik_a}$  is less than or equal to  $2D_{thr}$ , the cluster  $k_a$  is updated by moving its center  $\vec{v}_{k_a}$  and modifying its radius  $R_{k_a}$ . The updated radius  $R_{k_a}$  is set to  $S_{ik_a}$  and the new center  $\vec{v}_{k_a}$  is located on the line connecting  $\vec{x}_i^{obs-nc}$  to the old cluster center  $\vec{v}_{k_a}$ , in a position such that the distance between the new center  $\vec{v}_{k_a}$  and the test pattern  $\vec{x}_i^{obs-nc}$  is equal to  $R_{k_a}$ ; then, go to Step 1.

Notice that this procedure guarantees that the maximum possible distance between a pattern of the training dataset and the nearest cluster center is lower or equal to the threshold value  $D_{thr}$ .

#### B. Reconstruction Procedure

Once the  $m$  clusters have been identified, the reconstruction  $\vec{x}_{nc} = (\hat{x}_{nc}(1), \dots, \hat{x}_{nc}(q))$  of a test pattern  $\vec{x}^{obs} = (x^{obs}(1), \dots, x^{obs}(q))$  is based on the following two steps:

1) identification of the cluster with the smallest distance from the test pattern. This is done by computing the distance  $d(\vec{x}^{test}, \vec{v}_k)$  of the test patterns with all the  $m$  cluster centers:

$$d^2(\vec{x}^{test}, \vec{v}_k) = \left\| \vec{x}^{test} - \vec{v}_k \right\|^2, \quad k=1, \dots, m \quad (6)$$

and selecting as nearest cluster  $k_{closest}$  the one with the minimum distance:

$$k_{closest} = \arg \min_k \left\{ d(\vec{x}^{test}, \vec{v}_k) \right\}, \quad k=1, \dots, m \quad (7)$$

2) Reconstruction of the test pattern as the nearest cluster center:

$$\vec{x}^{obs} = \vec{v}_{k_{closest}} \quad (8)$$

Notice that contrary to other algorithms, such as AAKR which requires to perform the reconstruction of a test pattern to have the accessibility to all the training patterns, the ECM reconstruction is based only on the cluster centers and thus it does not require the continuous accessibility of the training set.

#### IV. CONFIDENCE MEASURE

In this Section, a measure of the degree of confidence of the signal reconstruction performed by using the ECM algorithm described in the previous Section III.A is proposed. Basically, the objective is to answer to the questions: how accurate is the reconstruction expected to be? To what degree can we trust the reconstruction?

To this purpose, notice that the metrics ‘‘accuracy’’ and ‘‘robustness’’ proposed in literature [1,2] to estimate the overall performance of the reconstruction algorithm on a set of test patterns different from those used to train the model are not satisfactory for our objective, since the reconstruction performance is expected to vary in different zones of the training space. Thus, the degree of confidence in the reconstruction should not be a fixed quantity independent from the location of the test pattern, but should vary according to the density and information content of the example patterns available to train the model.

##### A. Requirements for a confidence measure

The following four requirements which take into account the position of the test pattern with respect to the clusters found by the ECM algorithm should be considered in order to evaluate the confidence in the reconstruction:

- 1) smaller is the distance between the test pattern and the nearest cluster center, higher should be the reconstruction confidence;
- 2) if the test pattern has nearly the same distance from two or more cluster centers, the reconstruction confidence should be low;
- 3) higher is the number of training patterns in the cluster nearest to the test pattern, higher should be the reconstruction confidence;
- 4) higher is the density of the training patterns in the cluster nearest to the test pattern, higher should be the reconstruction confidence.

Requirement 1) is motivated by the fact that if a pattern does not belong to a cluster or is far away from its center, its reconstruction as the center of the nearest cluster is expected to be not reliable. Requirement 2) considers that we should not be confident in the reconstruction of ambiguous patterns, i.e.

patterns whose position is such that they can belong with the same confidence to two or more clusters. Finally, requirements 3) and 4) are related to the density and information content of the training data: the number of patterns in the cluster is considered since we are less confident in a cluster formed by one or few training patterns which, for example, may correspond to an abnormal equipment condition and be erroneously introduced in the training dataset for the model of normal behavior.

### B. Definition of the confidence measure

Given a test pattern  $\vec{x}^{obs} = (x^{obs}(1), \dots, x^{obs}(q))$ , the degree of confidence of its reconstruction  $\vec{x}_{nc} = (\hat{x}_{nc}(1), \dots, \hat{x}_{nc}(q))$  is defined by:

$$Conf_r = \begin{cases} \frac{n_{k_{nearest}}}{n} \cdot \frac{1}{\sqrt{2\pi}\sigma_{k_{nearest}}} e^{-\frac{d^2(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{2\sigma_{k_{nearest}}^2}} \cdot \left[ 1 - \frac{(m-1)d(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{\sum_{k=1}^m d(\vec{x}^{obs}, \vec{v}_k)} \right] & \text{if } n_{k_{nearest}} \leq n_0 \\ \frac{1}{\sqrt{2\pi}\sigma_{k_{nearest}}} e^{-\frac{d^2(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{2\sigma_{k_{nearest}}^2}} \cdot \left[ 1 - \frac{(m-1)d(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{\sum_{k=1}^m d(\vec{x}^{obs}, \vec{v}_k)} \right] & \text{if } n_{k_{nearest}} > n_0 \end{cases}$$

where:

- $k_{nearest}$  represents the label of the cluster with the center nearest to the test pattern
- $\vec{v}_k$  the center of the  $k$ -th cluster
- $\sigma_k = \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} d^2(\vec{x}_i^{obs-nc}, \vec{v}_k)}$  the mean square distance between the cluster center  $\vec{v}_k$  and all the training patterns belonging to the cluster
- $m$  the total number of clusters
- $n_k$  the number of patterns belonging to the  $k$ -th cluster
- $n$  the total number of training patterns
- $n_0 =$  an integer number representing a threshold for the number of patterns belonging to the cluster: when the number of patterns in the nearest cluster to the test pattern is below  $n_0$ , the confidence in the reconstruction is decreased.

The term  $1 - \frac{(m-1)d(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{\sum_{k=1}^m d(\vec{x}^{obs}, \vec{v}_k)}$  is related to the

distance of the test pattern to all the cluster centers. This term may be rewritten as:

$$\begin{aligned} & 1 - \frac{(m-1)d(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{\sum_{k=1}^m d(\vec{x}^{obs}, \vec{v}_k)} \\ &= 1 - \frac{1}{\frac{\sum_{k=1; k \neq k_{nearest}}^m d(\vec{x}^{obs}, \vec{v}_k)}{m-1} + \frac{1}{d(\vec{x}^{obs}, \vec{v}_{k_{nearest}})} + \frac{1}{m-1}} \\ &= 1 - \frac{1}{\frac{\bar{d}_{others}}{m-1} + \frac{1}{d(\vec{x}^{obs}, \vec{v}_k)}} \\ & \quad \bar{d}_{others} = \frac{\sum_{k=1; k \neq k_{nearest}}^m d(\vec{x}^{obs}, \vec{v}_k)}{m-1} \end{aligned} \quad (10)$$

Since  $\bar{d}_{others} = \frac{\sum_{k=1; k \neq k_{nearest}}^m d(\vec{x}^{obs}, \vec{v}_k)}{m-1}$  can be interpreted as the mean distance between the test pattern and all the cluster centers except the nearest one, the term in (10) results highest for those test patterns very close to the nearest cluster center and far away from all the other clusters, whereas it reaches its minimum value (1/m) in the case in which the test pattern has the same distance to all the cluster centers. These characteristics are in accordance to requirements 1) and 2). This term is derived from [12] where it is used for the estimation of the classification confidence of a  $k$ -nearest neighbor algorithm.

The term  $\frac{1}{\sqrt{2\pi}\sigma_{k_{nearest}}} e^{-\frac{d^2(\vec{x}^{obs}, \vec{v}_{k_{nearest}})}{2\sigma_{k_{nearest}}^2}}$  is introduced in (9)

to take into account the density of the nearest cluster to the test pattern. It depends from the square of the distance between the cluster center and the cluster, and the inverse of the quantity  $\sigma_k$  which is a measure of the cluster dispersion. Thus, once fixed the distance of the test pattern to the nearest cluster center, the degree of confidence in the reconstruction of the test pattern is higher for denser clusters as requested by requirement 4).

Finally, the factor  $\frac{n_k}{n}$  has been introduced in the confidence measure in order to meet requirement 3); basically, if the number of patterns in the nearest cluster is lower than the threshold value the confidence of the reconstruction is decreased of the factor  $\frac{n_k}{n}$ .

## V. VERIFICATION OF THE CONFIDENCE MEASURE

In Section V.A some case studies based on artificial bi-dimensional datasets are designed in order to verify if the proposed confidence measure effectively meets in practice the four requirements of Section IV.A. Then, in Section V.B the proposed confidence measure is applied to a real case study and the obtained results commented.

### A. Application to artificial case studies

In the first case study, two clusters, each one formed by 100 training patterns, have been generated from two bi-dimensional Gaussian distributions centered on (30, 100) and (70, 100), respectively, both with variance equal to 10 (red dots in Fig. 2(a)). The application of the ECM procedure with threshold

parameter  $D_{thr} = 0.15$  has led to the identification of three clusters whose center position along the horizontal axis is reported by a circle in Fig. 2(b), one containing the patterns on the left side of Fig. 2(a) and two containing the patterns on the right side. The obtained clusters have been used for the reconstruction of test patterns on the horizontal straight line passing for the centers of the Gaussian distributions (crosses in Fig. 2(a)). Fig. 2(b) reports the confidence in the reconstruction. Notice that in accordance with the first requirement, the confidence increases as the distance between the test pattern and the closest cluster center decreases. Although the data have been generated from two Gaussian symmetric distributions, a non symmetric behavior of the confidence measure in the two sides of Fig. 2(b) has been found due to the fact that the ECM algorithm results in one cluster for the right Gaussian distribution patterns and two for the left Gaussian distribution patterns.

The training data used for the second case study have been generated by using the same bi-dimensional Gaussian distributions considered in case study 1, whereas the test set is formed by patterns on a circle (crosses in Fig. 3(a)) with center (30, 70) equal to the center of the left Gaussian distribution. Fig. 3(b) shows the variation of the confidence measure as a function of the angle  $\theta$  scanning the circle starting at the dark dotted line in Fig. 3(a) (angle values increase in the counter clockwise direction). The results show that the degree of confidence achieved in the reconstruction of patterns which have the same distance from the center of the Gaussian distribution but different angles  $\theta$ , tends to be different. In particular, pattern A at  $\theta=0^\circ(360^\circ)$  is reconstructed with a degree of confidence equal to 0.33 whereas pattern B at  $\theta = 180^\circ$  with a degree of confidence equal to 0.39. This is in accordance with the second criterion of Section IV.A: the confidence on the reconstruction of pattern A is lower than that of pattern B since A is closer to the cluster on the right and thus more ambiguous. The peak of the confidence measure at  $\theta=270^\circ$  is due to an additional cluster center found by the ECM in proximity of the test pattern at  $\theta=270^\circ$ .

In the third case study, 100 test patterns positioned on a vertical line in the middle between the two clusters (Fig. 4 (a), crosses) are reconstructed by the ECM with threshold  $D_{thr}$  equal to 0.30. Fig. 4 (b) shows that the degree of confidence in the reconstruction of this pattern is lower than 0.3. This result confirms the fulfillment of the second requirement according to which patterns with the same distance to two cluster centers should be reconstructed with a low degree of confidence. Furthermore, according to the first requirement, as the distance to the cluster center increases (values on the vertical axis tend to 50 or 150), the confidence tends to decrease.

The fourth case study is similar to the third one except that the ECM threshold  $D_{thr}$  is decreased to 0.15 in order to allow the generation of small clusters formed by few patterns. Notice that in the case in which the nearest cluster to the test pattern is formed by a single training pattern, the confidence in the reconstruction becomes 0. This is the case of the two test patterns indicated by the circles in Fig. 5 (b).

In the last case study, the training patterns are taken from two Gaussian distributions centered on (20, 100) and (80, 100) and with variance 10 and 100, respectively (Fig 6(a)). The degrees of confidence in the reconstruction of the test patterns belonging to the cluster on the left, characterized by higher density, are higher than those of the less dense cluster on the right.

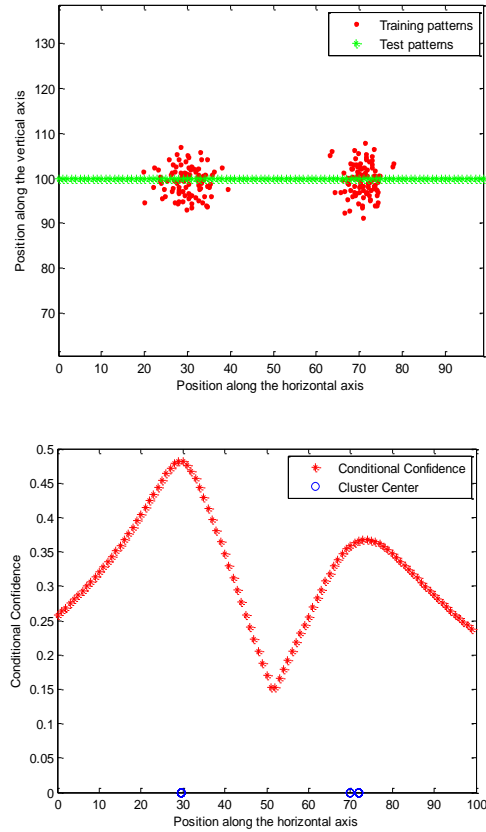


Figure 2 Case study 1: variation of the degree of confidence in the reconstruction as a function of the distance of the test pattern to the cluster centers.

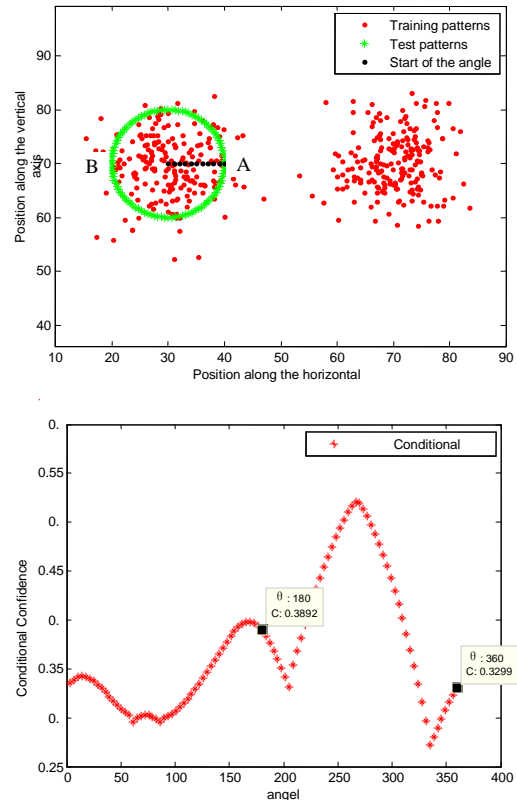


Figure 3 Case study 2: variation of the degree of confidence in the reconstruction depending on the presence of other clusters.

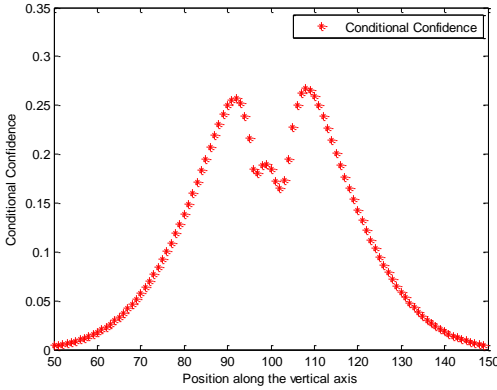
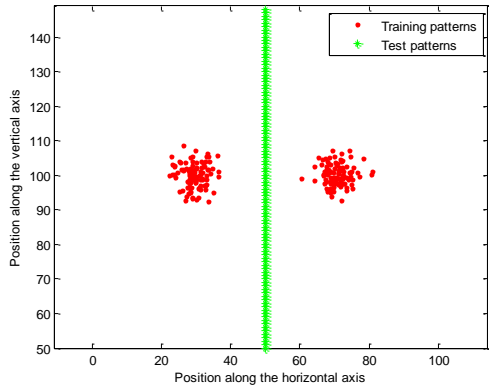


Figure 4 Case study 3: degree of confidence for patterns equally distant from two clusters (ECM with a large radius value).

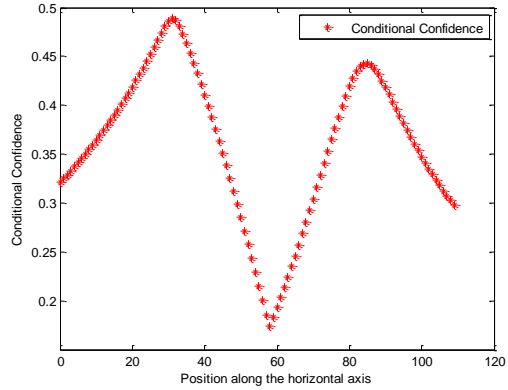
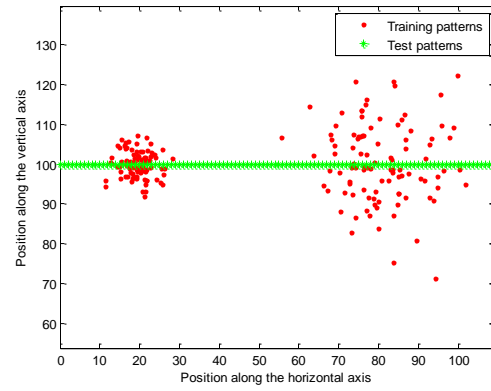


Figure 6 Case study 5: variation of the degree of confidence in function of the cluster density.

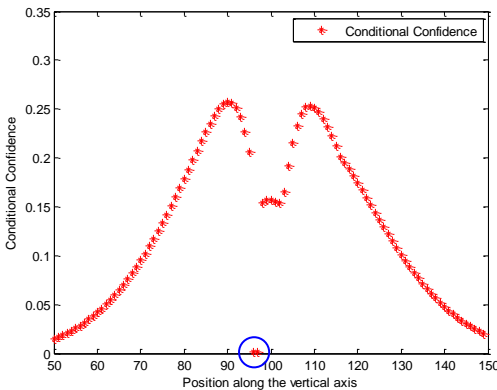
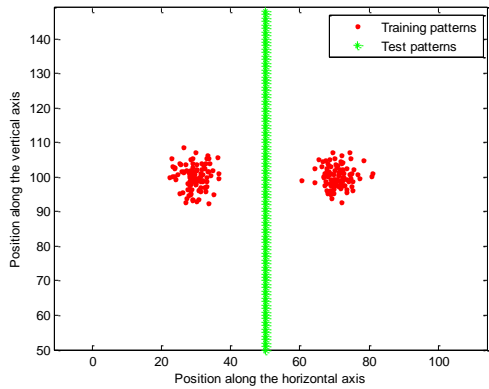


Figure 5 Case study 4: degree of confidence for patterns equally distant from two clusters (ECM with a small radius value).

These five case studies performed with properly designed bi-dimensional Gaussian data have shown that the proposed confidence measure meets the four requirements of Section IV.A.

#### B. Confidence estimation in a real condition monitoring case study

A real case study concerning 48 signals used to monitor the Reactor Coolant Pump (RCP) of a French Pressurized Water Reactor (PWR) is considered in this Section. The signals values have been measured every hour for a period of 11 consecutive months and concern four RCPs, each one on a different line of the primary circuit. The 5768 48-dimensional available patterns have been divided into a training set of 3000 and a test set containing the remaining 2798 test pattern.

The ECM parameter  $D_{thr}$  has been set equal to the value of 0.05 in an attempt to find an optimal compromise between a low value which would generate very accurate but low robust reconstructions (several clusters formed by few patterns) and an high value which would lead to less accurate but more robust reconstructions (few clusters formed by several patterns).

Fig. 7 reports the variation of the degree of confidence with the distance between the test pattern and the nearest cluster center. In accordance to the first requirement, the confidence tends to decrease as the distance between the test pattern and the nearest cluster center increases. Notice that this distance is equal to the reconstruction error since the reconstruction of the test pattern coincides with the nearest cluster center (Section III.B). Furthermore, Fig. 7 shows that there are some patterns reconstructed with confidence 0. These test patterns are assigned to a cluster formed by a single pattern far away from the other patterns; thus, their reconstructions are not believed reliable in accordance to the third criterion. These low degrees

of confidence are justified from the practical point of view by the fact that the training dataset may erroneously contain isolated patterns which do not correspond to normal operation and for this reason should not be considered for the reconstruction of the equipment behavior in normal conditions.

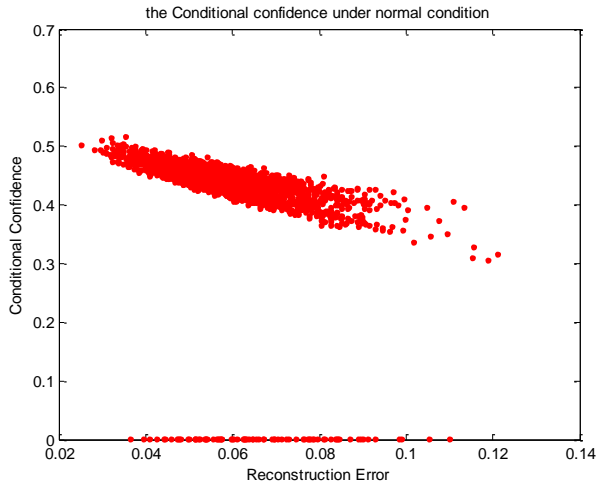


Figure 7 Variation of the degree of confidence of the signal reconstruction with the distance between the test pattern and the nearest cluster center..

A second real case study in which two different training sets are considered has also been performed. Training set 1 is formed by 4000 patterns, training set 2 by 400 patterns randomly sampled from the training patterns of training set 1. Given the composition of the training sets, the clusters obtained by applying the ECM method to training set 2 are expected to be less dense than those obtained based on training set 1. Table 1 reports the mean confidence achieved in the classification of the same test set. In accordance to requirement 4), when the denser training set 1 is used, the mean degree of confidence tends to become higher.

TABLE I. AVERAGE CONFIDENCE IN THE RECONSTRUCTION OF THE TEST PATTERNS IN FUNCTION OF THE CLUSTER DENSITY

Case	Number of training patterns	Average Confidence
raining set 1	4000	0.439
training set 2	400	0.360

## VI. CONCLUSIONS

Due to the criticality of condition monitoring in complex systems, it is important to provide decision makers with not only an estimation of the equipment health state but also a measure of the confidence in the condition monitoring model outcome. To this purpose, we have considered a condition monitoring scheme based on two modules performing (1) signal reconstruction and (2) diagnostics decision on the equipment health state. As reconstruction model, we have employed an auto-associative empirical algorithm based on the

Evolving Clustering Method (ECM). We have, then, introduced a novel measure of confidence on the obtained signal reconstruction, verified it by application to artificial datasets and applied it to a real condition monitoring problem concerning an important component of a nuclear power plant. The results that we have obtained show that the proposed measure meets the requirements which the confidence is expected to depend on, related to the density and information content of the training space.

Future research activity should be devoted to the estimation of the confidence in the diagnostic decision on the equipment health state and to the proper combination of the developed confidence measures in order to obtain an overall measure of confidence in the condition monitoring process.

## REFERENCES

- [1] P. Baraldi, R. Canesi, E. Zio, R. Seraoui, R. Chevalier, Signal Grouping for Condition Monitoring of Nuclear Power Plants Components. Seventh American Nuclear Society International Topical Meeting on Nuclear Plant Instrumentation, Control and Human-Machine Interface Technologies, Las Vegas, USA, 2010.
- [2] J. W. Hines and D. R. Garvey, Development and Application of Fault Detectability Performance Metrics for Instrument Calibration Verification and Anomaly Detection, *Journal of Pattern Recognition Research* 1, pp. 2–15, 2006.
- [3] R. Chevalier, D. Provost, R. Seraoui, Assessment of Statistical and Classification Models For Monitoring EDF's Assets, Sixth American Nuclear Society International Topical Meeting on Nuclear Plant Instrumentation, Control and Human-Machine Interface Technologies, Knoxville, USA, 2009.
- [4] J. Reifman, "Survey of artificial intelligence methods for detection and identification of component faults in nuclear power plants," *Nuclear Technology*, **119**, pp. 76-97. 1997.
- [5] D. Roverso, M. Hoffmann, E. Zio, P. Baraldi, and G. Gola, Solutions for plant-wide on-line calibration monitoring, *Proc. ESREL 2007*, Stavanger, Norway, 2007, Vol. 1, pp.827–832.
- [6] K.C. Gross, K.E. Kumenik, Sequential probability ratio test for nuclear power plant component surveillance, *Nuclear Technology* 93, pp. 131–137. 2007.
- [7] P. Baraldi, A. Cammi, F. Mangili, E. Zio, An ensemble approach to sensor fault detection and signal reconstruction for nuclear system control, *Annals of Nuclear Energy*, Vol 37, No.6, pp. 778-790. 2010.
- [8] M. Na, "Failure detection using a fuzzy neural network with an automatic input selection algorithm", *Intelligent Hybrid Systems. Fuzzy Logic, Neural Network, and Genetic Algorithms*, D. Rua, Springer, New York. 1997.
- [9] P. Baraldi, R. Razavi-Far, E. Zio "A Method for Estimating the Confidence in the Identification of Nuclear Transients by a Bagged Ensemble of FCM Classifiers" Seventh American Nuclear Society International Topical Meeting on Nuclear Plant Instrumentation, Control and Human-Machine Interface Technologies NPIC&HMIT 2010, Las Vegas, Nevada, November 7-11, 2010, on CD-ROM.
- [10] I.T. Jolliffe, "Principal Component Analysis", Springer Eds., 2002.
- [11] M. A. Kramer, Autoassociative Neural Networks, *Computers & Chemical Engineering*, 1992, Vol. 16, No. 4, pp. 313-328
- [12] R.P. Duin, M. Tax, "Classifier conditional posterior probabilities," *Lecture Notes in Computer Science*, LNCS 1451, pp. 611-619, 1998.