



STADIUS

Center for Dynamical Systems,
Signal Processing and Data Analytics

Citation/Reference	Langone R., Alzate C., De Ketelaere B., Vlasselaer J., Meert W., Suykens J.A.K., " LS-SVM based spectral clustering and regression for predicting maintenance of industrial machines ", <i>Engineering Applications of Artificial Intelligence</i> , vol. 37, Jan. 2015, pp. 268-278
Archived version	Author manuscript: the content is identical to the content of the published paper, but without the final typesetting by the publisher
Published version	insert link to the published version of your paper http://dx.doi.org/10.1016/j.engappai.2014.09.008
Journal homepage	insert link to the journal homepage of your paper http://www.journals.elsevier.com/engineering-applications-of-artificial-intelligence/
Author contact	your email rocco.langone@kuleuven.be Klik hier als u tekst wilt invoeren.
IR	url in Lirias https://lirias.kuleuven.be/handle/123456789/463431

(article begins on next page)



LS-SVM based spectral clustering and regression for predicting maintenance of industrial machines

Rocco Langone^{1,5}, Carlos Alzate^{1,4}, Bart De Ketelaere², Jonas Vlasselaer,³
Wannes Meert,³ Johan A. K. Suykens^{1,5}

¹*Department of Electrical Engineering (ESAT), STADIUS, KU Leuven, B-3001 Leuven Belgium
Email: {rocco.langone, carlos.alzate, johan.suykens}@esat.kuleuven.be*

²*Faculty of Bioscience Engineering, BIOSYST MeBioS Qualimetrics, KU Leuven, B-3001
Leuven Belgium
Email: bart.deketelaere@biw.kuleuven.be*

³*Department of Computer Science, KU Leuven, Celestijnenlaan 200a B-3001 Leuven Belgium
Email: {jonas.vlasselaer, wannes.meert}@cs.kuleuven.be*

⁴*Smarter Cities Technology Center, IBM Research-Ireland
Email: carlos.alzate@ie.ibm.com*

⁵ *iMinds Future Health Department*

Abstract

Accurate prediction of forthcoming faults in modern industrial machines plays a key role in reducing production arrest, increasing the safety of plant operations, and optimizing manufacturing costs. The most effective condition monitoring techniques are based on the analysis of historical process data. In this paper we show how Least Squares Support Vector Machines (LS-SVMs) can be used effectively for early fault detection in an online fashion. Although LS-SVMs are existing artificial intelligence methods, in this paper the novelty is represented by their successful application to a complex industrial use case, where other approaches are commonly used in practice. In particular, in the first part we present an unsupervised approach that uses Kernel Spectral Clustering (KSC) on the sensor data coming from a vertical form seal and fill (VFFS) machine, in order to distinguish between normal operating condition and abnormal situations. Basically, we describe how KSC is able to detect in advance the need of maintenance actions in the analysed machine, due the degradation of the sealing jaws. In the second part we illustrate a nonlinear auto-regressive model (NAR), thus a supervised learning technique, in the LS-SVM framework. We show that we succeed in modelling appropriately the degradation process affecting the machine, and we are capable to accurately predict the evolution of

dirt accumulation in the sealing jaws.

Keywords: kernel spectral clustering; LS-SVMs; NAR; time-series prediction; condition monitoring; maintenance; fault detection, machine degradation, artificial intelligence.

1. Introduction

In industrial processes, the detection and analysis of faults ensure product quality and operational safety [26]. Traditionally, three ways to deal with sensory faults have been used [30],[31],[32]: corrective maintenance, preventive maintenance and predictive maintenance. Corrective maintenance is performed only when the machine fails, it is expensive and safety and environmental issues arise. Preventive maintenance [16] is based on periodic replacement of components, which are then utilized in a non-optimal way. Predictive maintenance [5] can be performed in a manual or automatic fashion. In the first case machines are manually checked with expensive monitoring techniques and the components are replaced according to their real status. In the second case a machine's status is automatically inspected and maintenance is planned accordingly. The continuous monitoring of machine parts leads to reliable and accurate lifetime predictions, and maintenance operations can be fully automated and implemented in a cost effective way.

Nowadays, in many industries several process variables like temperature, pressure etc. can be measured. These measurements give an information on the current status of a machine and can be used to predict the faults due to deterioration of key components [23, 10]. As a consequence, an optimal maintenance strategy can be planned.

Condition monitoring using sensor data has been performed for long time by means of basic methods like exponentially weighted moving average, cumulative sum, principal component analysis (PCA) [9], [6]. Only in the past few years engineers in companies started to get convinced to use more advanced techniques for fault detection, like for instance SVMs approaches in semiconductor manufacturing [25, 15]. In this realm, with the aim of contributing to bridge the gap between academia and industry, we propose to use LS-SVMs for predictive maintenance. LS-SVMs [28, 27] are an artificial intelligence technique characterized by an high quality generalization capability [29], the flexibility in the model design and a clear procedure for model selection. Basically, for a given machine we can construct a reliable model of the degradation process to be able to opti-

32 mize the time of maintenance. For this purpose here we propose two alternative
33 approaches with different associated costs, applied to predictive maintenance in a
34 vertical form seal and fill (VFFS) machine:

- 35 • clustering of vibration signals collected by accelerometers. This technique
36 is cheap because of the relatively low cost of the accelerometers. On the
37 other hand, the clustering model does not directly predict machine's deteri-
38 oration, but it infers the degradation based on vibration data.
- 39 • regression applied to thermal camera data. The number of hot area pixels
40 present in the images acquired by the camera represent a direct measure of
41 degradation. A nonlinear autoregressive model is then used to forecast the
42 future trend of machine's deterioration. This procedure is more reliable than
43 the clustering-based methodology because in this case a direct prediction of
44 the degradation is achieved. On the other hand, it is more expensive because
45 of the higher cost of the thermal camera w.r.t. the accelerometers.

46 In the first part of this work we describe the use of kernel spectral clustering
47 (KSC, [2]) for identifying in advance when the (VFFS) machine enters critical
48 conditions. A similar analysis has been carried out in [12], where only two exper-
49 iments instead of three were considered and no comparisons with other techniques
50 were performed. In contrast to [11] and [13] where an initial clustering model is
51 updated over time, in this framework we assume stationarity. Thus the algorithm
52 is built off-line and then used online (by means of the out-of-sample extension
53 property) in order to distinguish between normal operating condition and abnor-
54 mal situations. Then we use the model in an online fashion via the out-of-sample
55 extension property of KSC to recognize these two regimes. In particular, KSC
56 correctly infers the degradation process from the vibration signals registered by
57 the accelerometers placed on the sealing jaws.

58 In the second part of the paper we utilize a nonlinear auto-regressive model
59 (NAR [17]) to catch the deterioration phenomenon acting on the machine. Like
60 KSC, the model is cast in the LS-SVM framework and, once properly trained,
61 it can be used in an online manner thanks to the out-of-sample extension ability.
62 In this experiment, the machine is equipped with a thermal camera that directly
63 measures the degradation in terms of number of hot area pixels in the acquired
64 images. Then the NAR model is used to predict the evolution of the number of
65 hot area pixels over time.

66 The remainder of this paper is structured as follows: Section 2 summarizes
67 the KSC model and the related model selection scheme, i.e. BLF (Balanced Line

68 Fit). Section 3 explains the basics of the LS-SVM regression method and how it
 69 can be used to formulate a NAR model. Section 4 describes the data-sets used
 70 in the experiments. In Section 5 the simulation results are presented. Regarding
 71 KSC, a comparison with k-means and self organizing maps is also depicted. For
 72 what concerns the LS-SVM NAR model, we show that it outperforms a linear
 73 auto-regressive model (AR) and gives very good predictions up to 10 steps ahead.
 74 Section 6 presents a general discussion. Finally, Section 7 concludes the paper.

75 2. Kernel Spectral Clustering

76 2.1. Model

77 Spectral clustering methods use the eigenvectors of the Laplacian to properly
 78 group the data-points [7, 33, 21]. KSC is a spectral clustering formulation cast
 79 in the LS-SVM primal-dual setting, being the primal problem a constrained op-
 80 timization problem with a quadratic loss function and equality constraints. As
 81 mentioned in the introduction, the KSC method has two main advantages, a sys-
 82 tematic model selection scheme for the correct tuning of the parameters and the
 83 extension of the clustering model to out-of-sample points. In KSC, a clustering
 84 model can be trained on a subset of the data and then applied to the rest of the
 85 data in a learning framework. The out-of-sample extension allows then to predict
 86 the memberships of a new point thanks to the model learned during the training
 87 phase. In this way, once a model of the operation of a machine has been con-
 88 structed, we can use it in an online fashion to discover when the machine enters
 89 critical conditions.

Given a training data set $\mathcal{D} = \{x_i\}_{i=1}^N, x_i \in \mathbb{R}^d$, the multi-cluster KSC model [2] is formulated as a weighted kernel PCA problem [19] decomposed in $l = k - 1$ binary problems, where k is the number of clusters to find:

$$\min_{w^{(l)}, e^{(l)}, b_l} \frac{1}{2} \sum_{l=1}^{k-1} w^{(l)T} w^{(l)} - \frac{1}{2N} \sum_{l=1}^{k-1} \gamma_l e^{(l)T} D^{-1} e^{(l)} \quad (1)$$

$$\text{such that } e^{(l)} = \Phi w^{(l)} + b_l 1_N \quad (2)$$

90 The $e^{(l)} = [e_1^{(l)}, \dots, e_N^{(l)}]^T$ are the projections of the data points $\{x_i\}_{i=1}^N$ mapped
 91 in the feature space along the direction $w^{(l)}$, also called score variables. The opti-
 92 mization problem (1) can then be interpreted as the maximization of the weighted
 93 variances $C_l = e^{(l)T} D^{-1} e^{(l)}$ and the contextual minimization of the squared norm
 94 of the vector $w^{(l)}$, $\forall l$. Through the regularization constants $\gamma_l \in \mathbb{R}^+$ we trade-
 95 off the model complexity expressed by $w^{(l)}$ with the correct representation of the

96 training data. The symbol 1_N represents a column vector with the N components
 97 equal to 1, $D^{-1} \in \mathbb{R}^{N \times N}$ is the inverse of the degree matrix D , Φ is the $N \times d_h$ fea-
 98 ture matrix $\Phi = [\varphi(x_1)^T; \dots; \varphi(x_N)^T]$ and $\gamma_l \in \mathbb{R}^+$ are regularization constants.
 99 The clustering model is formulated as:

$$e_i^{(l)} = w^{(l)T} \varphi(x_i) + b_l, i = 1, \dots, N \quad (3)$$

100 where $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^{d_h}$ is the mapping to a high-dimensional feature space, b_l are
 101 bias terms, with $l = 1, \dots, k-1$. The projections $e_i^{(l)}$ represent the latent variables
 102 of a set of $k-1$ binary clustering indicators given by $\text{sign}(e_i^{(l)})$. The set of binary
 103 indicators form a code-book $\mathcal{CB} = \{c_p\}_{p=1}^k$, where each code-word is a binary
 104 word of length $k-1$ representing a cluster. After constructing the Lagrangian and
 105 solving the Karush-Kuhn-Tucker (KKT) conditions for optimality the following
 106 dual problem is obtained:

$$D^{-1} M_D \Omega \alpha^{(l)} = \lambda_l \alpha^{(l)} \quad (4)$$

107 where Ω is the kernel matrix with ij -th entry $\Omega_{ij} = K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$, D
 108 is the graph degree matrix which is diagonal with positive elements $D_{ii} = \sum_j \Omega_{ij}$,
 109 M_D is a centering matrix defined as $M_D = I_N - \frac{1}{1_N^T D^{-1} 1_N} 1_N 1_N^T D^{-1}$, the $\alpha^{(l)}$ are
 110 dual variables, $\lambda_l = \frac{N}{\gamma_l}$. $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is the kernel function and captures
 111 the similarity between the data-points. In all the experiments we use the RBF
 112 kernel function defined by $K(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2 / \sigma^2)$, with σ meaning
 113 the bandwidth parameter. The out-of-sample extension to new nodes is done by
 114 an Error Correcting Output Codes (ECOC) decoding procedure. In the decoding
 115 process the cluster indicators obtained in the validation/test stage are compared
 116 with the code-book and the nearest code-word in terms of Hamming distance is
 117 selected. The cluster indicators can be obtained by binarizing the score variables
 118 for test points in the following way:

$$\text{sign}(e_{\text{test}}^{(l)}) = \text{sign}(\Omega_{\text{test}} \alpha^{(l)} + b_l 1_{N_{\text{test}}}) \quad (5)$$

119 with $l = 1, \dots, k-1$. Ω_{test} is the $N_{\text{test}} \times N$ kernel matrix evaluated using the test
 120 nodes with entries $\Omega_{\text{test},ri} = K(x_r^{\text{test}}, x_i)$, $r = 1, \dots, N_{\text{test}}$, $i = 1, \dots, N$.

121 2.2. Tuning scheme

122 The number of clusters in which to group the data and the parameters of the
 123 kernel function (if any) have to be selected properly to ensure good performances.

124 The model selection scheme used in the experiments is Balanced Line Fit (BLF).
 125 The BLF criterion reaches its maximum value 1 when the clusters are well sepa-
 126 rated, and in this ideal situation they are represented as lines in the space of the
 127 projections $e^{(l)}$. The analytical formula of BLF is [2]:

$$\text{BLF}(\mathcal{D}^V, k) = \eta \text{linefit}(\mathcal{D}^V, k) + (1 - \eta) \text{balance}(\mathcal{D}^V, k) \quad (6)$$

128 where \mathcal{D}^V indicates the validation set and k the number of clusters. The linefit
 129 index equals 0 when the score variables are distributed spherically and equals 1
 130 when the score variables related to points in the same cluster are collinear. The
 131 balance index equals 1 when the clusters have the same size and tends to 0 in
 132 extremely unbalanced cases. The parameter η controls the importance given to the
 133 linefit with respect to the balance index and takes a value in the range $[0, 1]$. The
 134 BLF can be used to select the number of clusters and the kernel tuning parameters
 135 in the following way:

- 136 1. Define a grid of values for the parameters to select
- 137 2. Train the related KSC model
- 138 3. Compute the memberships of the validation points by means of the out-of-
 139 sample extension
- 140 4. For every partition of the validation set calculate the related score in terms
 141 of BLF
- 142 5. Choose the model with the highest score¹.

143 An example of the model selection procedure on a synthetic toy data-set is de-
 144 picted in Figure 1. It can be noticed that the maximum value of the BLF criterion
 145 is reached when the points in the projections space have a strong line structure,
 146 corresponding to optimal parameters. Although in a real-life problem the cluster
 147 structure can be less clear than in this toy example, BLF has shown to be a useful
 148 model selection criterion in several applications.

149 3. LS-SVMs for regression

150 3.1. Primal-Dual Formulation

151 As already pointed out previously, LS-SVMs are a kind of SVM where a
 152 quadratic loss function is used in the primal objective and equality instead of

¹Instead of considering only the highest value of BLF, also additional local maxima may be taken into account.

inequality constraints are present. This typically leads to a linear system or an eigenvalue problem at the dual level for classification and regression or (weighted) kernel PCA respectively. The main advantage of LS-SVMs with respect to other artificial intelligence methods like artificial neural networks is that a global optimum can be obtained, since the problem formulation is usually convex [28].

Given training data $\{x_i, y_i\}_{i=1}^N$, with $x_i \in \mathbb{R}^{d_x}$ and $y_i \in \mathbb{R}$, the regression problem in the primal space can be expressed as follows:

$$\min_{w, e, b} \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{i=1}^N e_i^2 \quad (7)$$

$$\text{such that } y_i = w^T \varphi(x_i) + b + e_i, i = 1, \dots, N. \quad (8)$$

The expression $\hat{y} = w^T \varphi(x) + b$ indicates the model in the primal space, and the objective function (7) is in fact a ridge regression cost function. With γ we indicate the regularization parameter which controls the trade-off between the model complexity and the minimization of the training error.

$$\left[\begin{array}{c|c} 0 & 1_N^T \\ \hline 1_N^T & \Omega + \mathbf{I}/\gamma \end{array} \right] \left[\begin{array}{c} b \\ \alpha \end{array} \right] = \left[\begin{array}{c} 0 \\ y \end{array} \right] \quad (9)$$

where, as before for KSC, $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^{d_h}$ represents the mapping to a high-dimensional feature space, and $y = [y_1; \dots; y_N]$, 1_N is a column vector of ones, $\alpha = [\alpha_1; \dots; \alpha_N]$. The term Ω means the kernel matrix with entries $\Omega_{ij} = \varphi(x_i)^T \varphi(x_j) = K(x_i, x_j)$, and with $K : \mathbb{R}^{d_x} \times \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ we denote the kernel function. By using a radial basis function (RBF) kernel, expressed by $K(x_i, x_j) = \exp(-\|x_i - x_j\|_2^2 / \sigma^2)$, one is able to construct a model of arbitrary complexity. Finally, after solving the previous linear system, the LS-SVM model for function estimation in the dual form becomes:

$$\hat{y} = \sum_{i=1}^N \alpha_i K(x, x_i) + b. \quad (10)$$

3.2. NAR model

A nonlinear auto-regressive model (NAR) describes time-varying phenomena specifying that the output variable depends non-linearly on its own previous values [17]:

$$\hat{y}_{k+1} = f([y_k; \dots; y_{k-p}]) \quad (11)$$

where \hat{y}_{k+1} is the predicted future value based on the previous p values, $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is the nonlinear mapping. The parameter p is also called order of the NAR model and has to be properly tuned. In the LS-SVM framework the NAR model can be described by:

$$\hat{y}_{k+1} = w^T \varphi([y_k; \dots; y_{k-p}]) + b \quad (12)$$

For multi-step ahead prediction a recursive approach can be used:

$$\begin{aligned} \hat{y}_{k+1} &= f([y_k; \dots; y_{k-p}]) \\ \hat{y}_{k+2} &= f([\hat{y}_{k+1}; \dots; y_{k-p+1}]) \\ &\dots \\ \hat{y}_{k+m} &= f([\hat{y}_{k+m-1}; \dots; \hat{y}_{k-p+m}]) \end{aligned}$$

where m is the number of steps ahead.

3.3. Model selection

In order to build-up a valid LS-SVM NAR model for time-series prediction, we have to choose very carefully the parameters γ , σ and p . The parameters γ and σ are selected using 10-fold cross-validation. The Coupled Simulated Annealing (CSA) is used to minimize the mean square error (MSE) in the cross-validation process. Also the simplex algorithm [20] is used afterwards to further improve the search of the optimum. CSA leads to an improved optimization efficiency because it reduces the sensitivity of the algorithm to the initialization parameters [34]. The optimal order p of the NAR model is found using a grid search approach.

4. Data-sets

The data are collected from a Vertical Form Fill and Seal (VFFS) machine used to fill and seal packages in different industries, but mainly in the food industry. The VFFS machine, illustrated in Figure 2, supplies film from a roll which is formed into a bag over the vertical cylinder. Sealing jaws close the bag at the bottom before it is filled. At the end of the cycle, the bag is sealed and cut off with a knife. From previous experimental studies the dirt accumulation on the sealing jaws was observed to strongly affect the process quality. For this reason, in the experiments described here the jaws were monitored to predict in advance the maintenance actions. Maintenance consists of stopping the machine and cleaning the sealing jaws. A total of three experiments have been performed. In the first

two experiments only accelerometers mounted on the jaws are used in order to indirectly detect the dirt accumulation. In the third experiment the machine is also equipped with a thermal camera, which directly measures the dirt accumulation in terms of number of hot area pixels in the acquired images (with hot we mean that the temperature of the sealing jaws is above a user-defined threshold). This has resulted in three data-sets:

- DS_I: this data-set consists of 771 events and 3 external maintenance actions. An event is related to a particular processed bag and takes place every two seconds (i.e. the sampling frequency is $0.5Hz$). Each event is associated with a 150-dimensional accelerometer signal, that is each signal is a vector of length 150 (see top of Figure 3).
- DS_II: it contains a total of 11 632 processed bags and 15 maintenance actions. Here the vibration signals used to monitor the dirt accumulation in the jaws are 190-dimensional time-series (as shown in the bottom part of Figure 3). This is due to a different setting of the data acquisition system for this experiment.
- DS_III: there are 3 519 processed bags, each one associated with a 150-dimensional vibration signal, and 11 maintenance actions. For this data-set we are given, beyond the the accelerometer signals, also the thermal camera measurements, as depicted in Figure 4.

The accelerometer signals registered just before the maintenance events are depicted in red in Figures 3 and 4 (top). They usually correspond to signals with high amplitude, which is related to the increased amount of dirt accumulated in the sealing jaws.

5. Experimental results

Here we present the analysis of the data described in the previous Section. The experiments have been conducted in Matlab R2013a, on a CORE i7 desktop PC with 16 GB of RAM memory. We used *KSClab* and *LS-SVMlab* available at <http://www.esat.kuleuven.be/stadius/ADB/software.php>

In the first part (unsupervised learning) we illustrate the performance of KSC and we compare it with k-means [18] and self organizing maps (SOMs [8]). K-means is considered among the most popular data clustering methods for industrial applications [4] and will be briefly described later in Section 5.1.4. SOMs

233 are an artificial neural network model able to recognize groups of similar input
234 vectors in an unsupervised fashion, thanks to a competitive learning strategy al-
235 lowing to perform vector quantization. The performance of KSC, k-means and
236 SOMs are evaluated according to a standard internal cluster quality measure, that
237 is the mean silhouette index, and the results are reported in table 1. The mean sil-
238houette criterion represents how well each object lies within its cluster, averaged
239 over every cluster [24].

240 In the second part (supervised learning) we explain the time-series prediction
241 on the thermal camera data by means of an LS-SVM NAR model. The compar-
242 ison with a linear model and a zero and first order extrapolation methods is also
243 discussed.

244 5.1. Unsupervised learning

245 In this Section it is shown how KSC can be used to perform just-in-time main-
246 tenance, not too early to take full advantage of component lifetime but also not
247 too late to avoid catastrophic failures and unplanned downtime. In particular, 2
248 regimes were identified, where one of them can be interpreted as normal behaviour
249 (low degradation) and the other as critical conditions (high degradation inducing
250 the need of maintenance). Moreover, a probabilistic interpretation of the results is
251 also provided, which better describes the degradation process experienced by the
252 sealing jaws of the packing machine.

253 We perform clustering on the raw accelerometer signals. However, since the
254 KSC technique is formulated as a weighted kernel PCA model (see Section 2.1), it
255 automatically extracts features from the vibration signals when performing clus-
256 tering. Another procedure could be related to exploiting the frequency content of
257 the signals by using for instance a wavelet kernel [35], but this approach has not
258 been considered in this work.

259 5.1.1. Model selection

260 In order to catch the ongoing deterioration process of the jaws we need to use
261 historical values of sealing quality in our analysis. For this purpose we apply a
262 windowing operation on the data (Figure 5). If we do not apply this preprocessing
263 step, we would perform clustering of only one signal at each time. In this way, the
264 algorithm would detect if the single bag has been sealed correctly or not. On the
265 other hand, when concatenating a certain number of signals, the clustering method
266 is able to distinguish between good working conditions (many good bags in the
267 current window) and faulty state (many bad bags in the current window).

268 We have a total of 3 parameters to determine: the window size (i.e. the number
 269 of signals to concatenate), the number of clusters k and the RBF kernel parameter
 270 σ . According to the BLF criterion the optimal window size is 40 and the optimal
 271 number of clusters is $k = 2$ for the three data-sets, while σ is data-set dependent.

272 Regarding data-set DS_I we have used 140 data points for training and the
 273 remaining 631 as test set. In the other experiments, we directly applied the model
 274 trained on the dataset DS_I on datasets DS_II and DS_III. Here, we exploited the
 275 generalization capability of the KSC algorithm and the fact that the experiments
 276 are similar to each other.

277 5.1.2. *Hard clustering*

278 In Figure 6 the KSC prediction for the data-set DS_I is shown. We can inter-
 279 pret one of the clusters as normal behaviour or low degradation and the other as
 280 maintenance cluster or high degradation². Notice that the KSC model is able to
 281 predict some minutes in advance the maintenance actions before they are actually
 282 performed by the operator. Concerning the data-set DS_II we can draw the same
 283 comments: KSC is very accurate in predicting the worsening of the packing pro-
 284 cess around the actual maintenance events (see top of Figure 7). Finally, Figure
 285 8 illustrates the results on data-set DS_III. In this case KSC predicts the need of
 286 maintenance also in zones where maintenance has not been really performed. As
 287 we will see later, surprisingly it would have been more logical to perform mainte-
 288 nance as suggested by KSC and not as actually done by the operator³.

289 5.1.3. *Soft clustering*

290 In the previous Section we demonstrated the effectiveness of KSC in predict-
 291 ing in advance the maintenance events. Nevertheless the predicted output is binary
 292 (it goes suddenly from normal operation to maintenance). On the other hand a soft
 293 clustering output is in general preferable, since it can provide more interpretable
 294 results [14]. Furthermore, in our case study the discrete output does not give us
 295 a continuous indicator of the incoming maintenance actions. To solve this is-
 296 sue we can use the latent variable $e(x)$ instead of the binarized clustering output

²Also for $k = 3$ the clustering results are meaningful: the BLF reaches the second highest value and three distinct regimes interpretable as normal behaviour, warning state and maintenance state have been observed. However, in this paper we only consider two clusters, corresponding to the maximum of the BLF model selection criterion.

³This conclusion can be corroborated by observing that some signals in proximity of maintenance events (depicted in red in at the top of Figure 4) have rather low amplitude.

297 $\text{sign}(e(x))$ (see Section 2.1).

298 Furthermore, we can rescale it between 0 and 1 to improve the interpretability.
299 This transformation is based on the structure of the latent variable space. As al-
300 ready mentioned in 2.2, in this space every cluster is ideally represented as a line.
301 The tips of the lines can be considered prototypes of their cluster, since they have
302 more certainty to belong to it because they are further from the decision bound-
303 aries [1]. Thus, the Euclidean distance from every point to the cluster prototype
304 can be seen as a confidence measure of the cluster membership. The transformed
305 latent variable is depicted at the bottom of Figure 6 (data-set DS_I), Figure 7 (data-
306 set DS_II) and Figure 8 (data-set DS_III). The value can be considered as a soft
307 membership or "probability" to maintenance. As explained in [3], this soft mem-
308 bership is given as a subjective probability and it indicates the strength of belief in
309 the clustering assignment. In other words, it is not a probability based on statistic
310 analysis of a large number of samples. Also, it is worth mentioning that since
311 the KSC algorithm does not naturally provide a soft output, we had to come up
312 with this post-processing step. However, other clustering methods like Hidden
313 Markov Models (HMMs, [22]) are designed in a probabilistic setting and can au-
314 tomatically supply moderated outputs (in case of HMMs also the probabilities of
315 transition between the different regimes are given).

316 By looking at the bottom of Figures 6, 7 and 8 we can see how the probability
317 increases as the number of faulty bags in the window increases. The value can
318 decrease since the window can move onto zones with good seals after a period
319 of bad seals. This is probably due to a self-cleaning mechanism. Maintenance is
320 predicted when the probability reaches the value 1.

321 For data-sets DS_I and DS_II it can be noticed how KSC is able to discover
322 from the vibration signals registered by the accelerometers the dirt accumulation
323 in the jaws that leads to the maintenance actions. Since clustering is an unsu-
324 pervised technique, this is achieved by not making use of any information on the
325 location of the maintenance actions (like it occurs for classification). For what
326 concerns data-set DS_III, also regions where no maintenance actions have been
327 really performed appear associated with high probability to maintenance.

328 To understand these unexpected outcome, we can rescale the probability to
329 maintenance in the same range of the the number of hot area pixels in the im-
330 ages captured by the thermal camera. This is illustrated in Figure 9, where we
331 depict in red the number of hot area pixels (measured degradation) and in blue
332 the rescaled soft memberships (modelled degradation). We can recognize similar
333 patterns, meaning that KSC was able to catch the degradation process also in this
334 last experiment. This last observation is quite surprising, since the clustering al-

gorithm is only supposed to divide the data into two main groups, as observed at the top side of Figures 6, 7, 8. The fact that the latent variable $e^{(l)}$ of the clustering model trained on the vibration signals follows the trend of the hot area pixels data (Figure 9) is somehow unexpected. To summarize, in the third experiment probably the operator should have performed the maintenance operations in a different way, being coherent with his behaviour in the first two experiments.

5.1.4. *k-means clustering*

K-means clustering is among the most popular methods for finding clusters in a set of data points. After choosing the desired number of cluster centers, the k-means procedure iteratively moves the centers to minimize the total within cluster variance. k-means has several drawbacks:

- the results are strongly influenced by the initialization
- the number of clusters should be provided by the user. This can be done by using a trial and error approach: the method is applied with different number of clusters and the partition corresponding to the highest value of some quality criterion is chosen.
- it can discover only spherical boundaries.

Despite these disadvantages, together with subtractive and hierarchical clustering, it is still widely used since it works effectively in many scientific and industrial applications. For this reason, here we present the results of k-means applied on the three data-sets described in Section 4. Before discussing the results, it is worth to mention that, thanks to the model selection scheme of KSC described in Section 2.2, we give optimal parameters to k-means (number of clusters = 2 and window size of concatenated accelerometer signals = 40). Figure 10 visualizes the outcomes of the hard and soft clustering. Similarly to KSC, the soft clustering results are based on the distance between the final centroids and the data-points in the input space (see [3]). Concerning data-sets DS_I and DS_III the results of KSC and k-means are very similar, while in the data-set DS_II analysis k-means performs worse than KSC, indicating the need of maintenance where it was not really performed. In this case k-means would suggest too many maintenance actions, which is not cost-effective.

5.2. *Supervised learning*

In this paragraph we present the results of applying the NAR model explained in Section 3.2 for predicting the evolution of the hot area pixels in data-set DS_III.

369 We have used the RBF kernel function $K(u, v) = \exp(-||u - v||_2^2/\sigma^2)$ to build
 370 the nonlinear model. The optimal parameters obtained using the 10 fold cross-
 371 validation procedure summarized in Section 3.3 are: $\gamma = 1356.7$, $\sigma^2 = 301.6$.
 372 Moreover, the optimal model order, tuned using a grid search approach, is $p = 5$,
 373 as depicted in Figure 11. After selecting the optimal parameters, we performed
 374 multi-step ahead recursive prediction. As can be noticed in Figure 12, the NAR
 375 model performed very well up to 10 steps ahead prediction, i.e. about 25 seconds
 376 since a bag is processed every 2 – 3 seconds. Afterwards, the mean absolute
 377 error between outputs and actual values starts increasing progressively, even if
 378 the predictions remains good up to 1 minute ahead prediction. In Figure 13 we
 379 visualize the results related to 10 steps ahead prediction. We can notice how
 380 the modelled degradation follows quite well the trend of the true degradation,
 381 in contrast with what was observed in Figure 9. This is not surprising because in
 382 this case we explicitly constructed a regression model to explain the deterioration,
 383 while the latent variable $e^{(l)}$ of the clustering model trained on the vibration signals
 384 in principle is not supposed to follow the trend of the hot area pixels data. The
 385 performance statistics are summarized in Table 2, where a comparison with two
 386 baseline techniques⁴ and a linear autoregressive model (AR) is shown.

387 6. Discussion

388 In the first part of this work we explained the use of KSC for predictive main-
 389 tenance. KSC has been chosen instead of classification because of the high un-
 390 balance of the data (few maintenance events compared to normal operating con-
 391 dition). After applying a windowing operation on the data in order to catch the
 392 deterioration process affecting the sealing jaws, we showed how KSC is able to
 393 recognize the presence of at least two working regimes of the VFFS machine,
 394 identifiable respectively as normal and critical operating condition. Moreover, we
 395 proposed also a soft clustering output that can be interpreted as "probability" to
 396 maintenance, and it is more directly related to the degradation phenomenon affect-
 397 ing the jaws (see for example Figure 9). In addition, as mentioned in Section 5.1.3,
 398 this probability is not defined in a statistical sense, but it is rather constructed to re-
 399 flect the reliability level of the mechanical equipment at current time. KSC is also
 400 compared with k-means and self organizing maps. While in data-sets DS_I and
 401 DS_III all methods give quite similar results, in case of data-set DS_II k-means

⁴The two baseline models correspond to zero order and first order extrapolation methods.

402 performs worse than KSC and SOMs, since it predicts maintenance in regions
403 where it was not actually performed by the operator.

404 In the second part a LS-SVM NAR model for predicting the evolution of the
405 dirt accumulation in the jaws has been constructed. The dirt accumulation is di-
406 rectly measured by means of the number of hot area pixels present in the images
407 obtained by the thermal camera. First we selected the parameters of the model
408 using 10-fold cross-validation with coupled simulated annealing for γ and σ^2 and
409 a grid search approach for the optimal model order. Then we trained the NAR
410 model and we tested its ability in predicting the evolution of the number of hot
411 area pixels. The performance is very good up to 10 steps ahead prediction (around
412 25 seconds), even if the forecast remains good up to 1 minute ahead. Moreover,
413 the NAR outperforms the linear auto-regressive model (AR) and a zero and first
414 order extrapolation methods.

415 To summarize, both the LS-SVM NAR model and KSC, thanks to their fore-
416 casting capabilities, could help to optimize the timing of maintenance actions for
417 the machine under study. In particular, a company who would like to use the LS-
418 SVM techniques described in this work to maximize production capacity has two
419 choices:

- 420 • option 1 (the cheapest): install on every packing machine only accelerom-
421 eters and use the proposed unsupervised learning method, that is KSC, to
422 predict in advance, in an online fashion, when the machine starts entering
423 critical conditions. In this way the operator can perform maintenance at the
424 right time and avoid long stops of the production.
- 425 • option 2: install the thermal camera and use the LS-SVM NAR model to
426 predict up to 1 minute ahead the dirt accumulation on the jaws. Since the
427 thermal camera costs more than the accelerometers, this option is more ex-
428 pensive. Nevertheless it is more reliable because it is a direct prediction
429 of the deterioration while KSC, as surprisingly effective as it is completely
430 unsupervised, infers the degradation indirectly from the vibration signals.

431 7. Conclusions

432 Predictive maintenance of industrial machines is receiving increasing attention
433 in the last years due to its many advantages, like cost efficiency, reduced ecolog-
434 ical impact etc. It is based on constant monitoring the health of machines joined
435 with advanced signal processing and prognostics techniques for optimal main-
436 tenance management. In this paper we proposed a least squares support vector

machine (LS-SVM) framework for maintenance strategy optimization based on real-time condition monitoring of a packing machine. We presented an unsupervised approach through kernel spectral clustering (KSC), and a supervised learning method, namely nonlinear auto-regression (NAR). In the first case we used the data collected by accelerometers positioned on the jaws of a Vertical Form Fill and Seal (VFFS) machine, from which the degradation process of the machine conditions has been inferred. For the time-series analysis with the NAR model data acquired by a thermal camera, which directly measures the dirt accumulation in the jaws, are processed. We showed that LS-SVM can successfully assess and predict mechanical conditions based on sensor data, thanks to their ability to model the degradation process. Moreover LS-SVM achieved higher performance than basic methods, which are commonly used in practice to predict the forthcoming faults. Finally, we proposed two options to use LS-SVM to schedule maintenance on the packing machine with different associated costs.

Acknowledgements

EU: The research leading to these results has received funding from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013) / ERC AdG A-DATADRIVE-B (290923). This paper reflects only the authors views and the Union is not liable for any use that may be made of the contained information. Research Council KUL: GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC), BIL12/11T; PhD/Postdoc grants Flemish Government: FWO: projects: G.0377.12 (Structured systems), G.088114N (Tensor based data similarity); PhD/Postdoc grants IWT: projects: SBO POM (100031); PhD/Postdoc grants iMinds Medical Information Technologies SBO 2014 Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimization, 2012-2017).

References

- [1] Alzate, C., Suykens, J. A. K., 2010. Highly sparse kernel spectral clustering with predictive out-of-sample extensions. In: Proc. of the 18th European Symposium on Artificial Neural Networks (ESANN 2010). pp. 235–240.
- [2] Alzate, C., Suykens, J. A. K., February 2010. Multiway spectral clustering with out-of-sample extensions through weighted kernel PCA. IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (2), 335–347.
- [3] Ben-Israel, A., Iyigun, C., Jun. 2008. Probabilistic d-clustering. J. Classif. 25 (1), 5–26.

- 470 [4] Berkhin, P., 2006. A survey of clustering data mining techniques. Grouping
471 Multidimensional Data, 25–71.
- 472 [5] Bey-Temsamani, A., Engels, M., Motten, A., Vandenplas, S., Ompusunggu,
473 A. P., 2009. A practical approach to combine data mining and prognostics for
474 improved predictive maintenance. In: Proceedings of the Third International
475 Workshop on Data Mining Case Studies (DMCS). pp. 36–43.
- 476 [6] Choi, S. W., Yoo, C. K., Lee, I.-B., 2003. Overall statistical monitoring of
477 static and dynamic patterns. Ind. Eng. Chem. Res. 42, 108 – 117.
- 478 [7] Chung, F. R. K., 1997. Spectral Graph Theory. American Mathematical So-
479 ciety.
- 480 [8] Kohonen, T., Schroeder, M. R., Huang, T. S. (Eds.), 2001. Self-Organizing
481 Maps, 3rd Edition. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- 482 [9] Kourti, T., MacGregor, J. F., 1995. Process analysis, monitoring and diag-
483 nosis, using multivariate projection methods. Chemometrics and Intelligent
484 Laboratory Systems 28 (1), 3 – 21.
- 485 [10] Kriegel, H.-P., Kroger, P., Zimek, A., 2010. Outlier detection techniques.
486 16th ACM International Conference on Knowledge Discovery and Data
487 Mining (SIGKDD).
- 488 [11] Langone, R., Agudelo, O. M., De Moor, B., Suykens, J. A. K., Septem-
489 ber 2014. Incremental kernel spectral clustering for online learning of non-
490 stationary data. Neurocomputing 139, 246–260.
- 491 [12] Langone, R., Alzate, C., De Ketelaere, B., Suykens, J. A. K., 2013. Ker-
492 nel spectral clustering for predicting maintenance of industrial machines. In:
493 IEEE Symposium Series on Computational Intelligence and Data Mining
494 (CIDM) 2013. pp. 39–45.
- 495 [13] Langone, R., Alzate, C., Suykens, J. A. K., 2013. Kernel spectral clustering
496 with memory effect. Physica A: Statistical Mechanics and its Applications
497 392 (10), 2588–2606.
- 498 [14] Langone, R., Mall, R., Suykens, J. A. K., 2013. Soft kernel spectral clus-
499 tering. In: Proc. of the International Joint Conference on Neural Networks
500 (IJCNN 2013). pp. 1028–1035.

- 501 [15] Li, T.-S., Huang, C.-L., 2009. Defect spatial pattern recognition using a
502 hybrid somsvm approach in semiconductor manufacturing. *Expert Systems*
503 *with Applications* 36 (1), 374 – 385.
- 504 [16] Liao, W., Pan, E., Xi, L., Dec. 2010. Preventive maintenance scheduling for
505 repairable system with deterioration. *J. Intell. Manuf.* 21 (6), 875–884.
- 506 [17] Ljung, L. (Ed.), 1999. *System identification (2nd ed.): theory for the user.*
507 Prentice Hall PTR.
- 508 [18] MacQueen, J., 1967. Some methods for classification and analysis of multi-
509 variate observations. In: of California Press., U. (Ed.), *Fifth Berkeley Sym-*
510 *posium on Mathematical Statistics and Probability*. Vol. 1. pp. 281–297.
- 511 [19] Mika, S., Schölkopf, B., Smola, A., Müller, K.-R., Scholz, M., Rätsch, G.,
512 1999. Kernel pca and de-noising in feature spaces. In: *Proceedings of the*
513 *1998 Conference on Advances in Neural Information Processing Systems*
514 *II*. pp. 536–542.
- 515 [20] Murty, K. G. (Ed.), 1983. *Linear programming*. Wiley.
- 516 [21] Ng, A. Y., Jordan, M. I., Weiss, Y., 2002. On spectral clustering: Analysis
517 and an algorithm. In: Dietterich, T. G., Becker, S., Ghahramani, Z. (Eds.),
518 *Advances in Neural Information Processing Systems 14*. MIT Press, Cam-
519 bridge, MA, pp. 849–856.
- 520 [22] Rabiner, L., Juang, B., 1986. An introduction to hidden markov models.
521 *IEEE Acoutics, Speech and Signal Processing Magazine* 3, 4–16.
- 522 [23] Randall, R. B., 2011. *Vibration-based condition monitoring industrial,*
523 *aerospace and automotive applications*. Wiley.
- 524 [24] Rousseeuw, P., Nov. 1987. Silhouettes: a graphical aid to the interpretation
525 and validation of cluster analysis. *J. Comput. Appl. Math.* 20 (1), 53–65.
- 526 [25] Sarmiento, T., Hong, S., May, G., April 2005. Fault detection in reac-
527 tive ion etching systems using one-class support vector machines. In: *Ad-*
528 *vanced Semiconductor Manufacturing Conference and Workshop, 2005*
529 *IEEE/SEMI*. pp. 139–142.

- 530 [26] Sharma, A., Yadava, G., Deshmukh, S., 2011. A literature review and fu-
531 ture perspectives on maintenance optimization. *Journal of Quality in Main-*
532 *tenance Engineering* 17, 5–25.
- 533 [27] Suykens, J. A. K., Alzate, C., Pelckmans, K., 2010. Primal and dual model
534 representations in kernel-based learning. *Statistics Surveys* 4, 148–183.
- 535 [28] Suykens, J. A. K., Van Gestel, T., De Brabanter, J., De Moor, B., Vande-
536 walle, J., 2002. *Least Squares Support Vector Machines*. World Scientific,
537 Singapore.
- 538 [29] Van Gestel, T., Suykens, J. A., Baesens, B., Viaene, S., Vanthienen, J., De-
539 dene, G., de Moor, B., Vandewalle, J., 2004. Benchmarking least squares
540 support vector machine classifiers. *Machine Learning* 54 (1), 5–32.
- 541 [30] Venkatasubramanian, V., Rengaswamy, R., Kavuri, S., 2003. A review of
542 process fault detection and diagnosis. part i: Quantitative model-based meth-
543 ods. *Computers and chemical engineering* 27 (3), 293–311.
- 544 [31] Venkatasubramanian, V., Rengaswamy, R., Kavuri, S., 2003. A review of
545 process fault detection and diagnosis. part ii: Qualitative models and search
546 strategies. *Computers and chemical engineering* 27 (3), 313–326.
- 547 [32] Venkatasubramanian, V., Rengaswamy, R., Kavuri, S., 2003. A review of
548 process fault detection and diagnosis. part iii: Process history based meth-
549 ods. *Computers and chemical engineering* 27 (3), 327–346.
- 550 [33] von Luxburg, U., 2007. A tutorial on spectral clustering. *Statistics and Com-*
551 *puting* 17 (4), 395–416.
- 552 [34] Xavier-De-Souza, S., Suykens, J. A. K., Vandewalle, J., Bollé, D., Apr. 2010.
553 Coupled simulated annealing. *Trans. Sys. Man Cyber. Part B* 40 (2), 320–
554 335.
- 555 [35] Zhang, L., Zhou, W., Jiao, L., 2004. Wavelet support vector machine. *IEEE*
556 *Transactions on Systems, Man, and Cybernetics, Part B* 34 (1), 34–39.

	DS_I	DS_II	DS_III
KSC	0.29	0.25	0.27
k-means	0.28	0.12	0.27
SOMs	0.29	0.21	0.28

Table 1: **Cluster quality evaluation.** Mean Silhouette index (see end of section 5.1.4), with best performance in bold. In the case of data-set DS_II the low value of Silhouette indicates that k-means does not succeed in correctly separating the normal behaviour and the maintenance cluster.

	LS-SVM NAR	AR	Baseline 1	Baseline 2
MAE	227.03	247.25	371.99	251.34
Percentage error	1.42%	4.11%	6.17%	4.18%
R^2	0.93	0.92	0.77	0.92

Table 2: **Summary NAR performance on test data.** Best performance in bold. By baseline 1 we mean a zero order extrapolation model and baseline 2 is related to a first order extrapolation method.

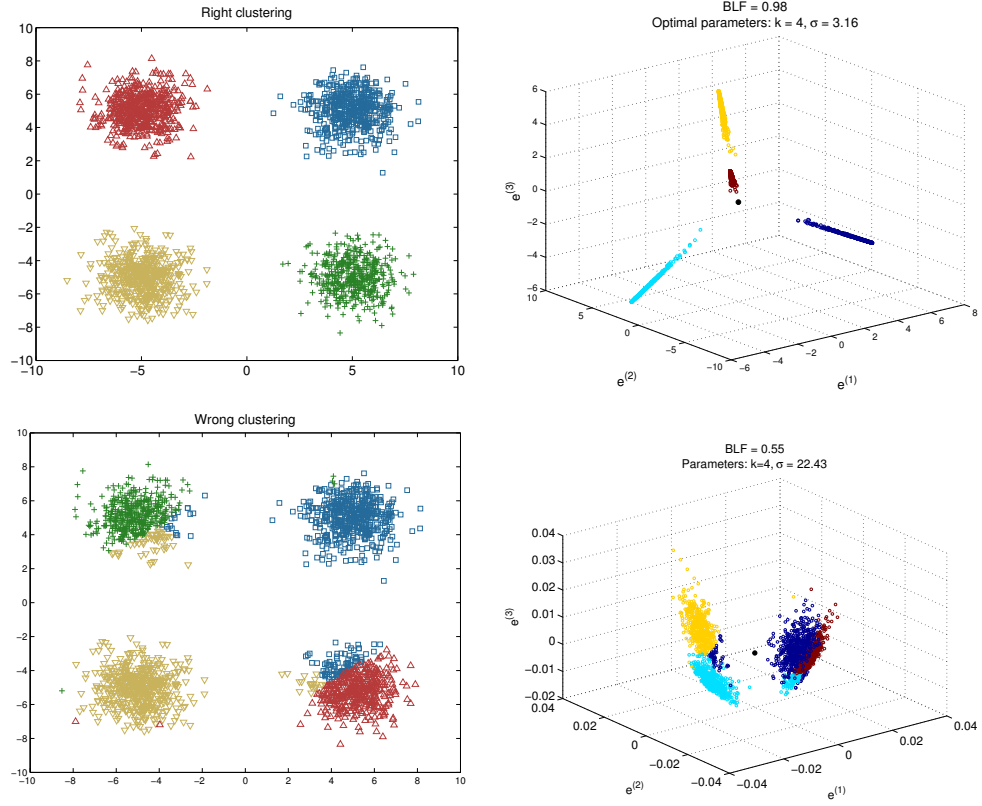


Figure 1: **Model selection illustrative example on synthetic toy data:** for well-chosen kernel parameters, the clusters for out-of-sample data are collinear in the projection space. **Top left:** Correct clustering of 4 Gaussian clouds. **Top right:** Projection space corresponding to an optimal σ value, for which the BLF is maximal. **Bottom left:** Wrong clustering results. **Bottom right:** Projection space for a wrong σ value.

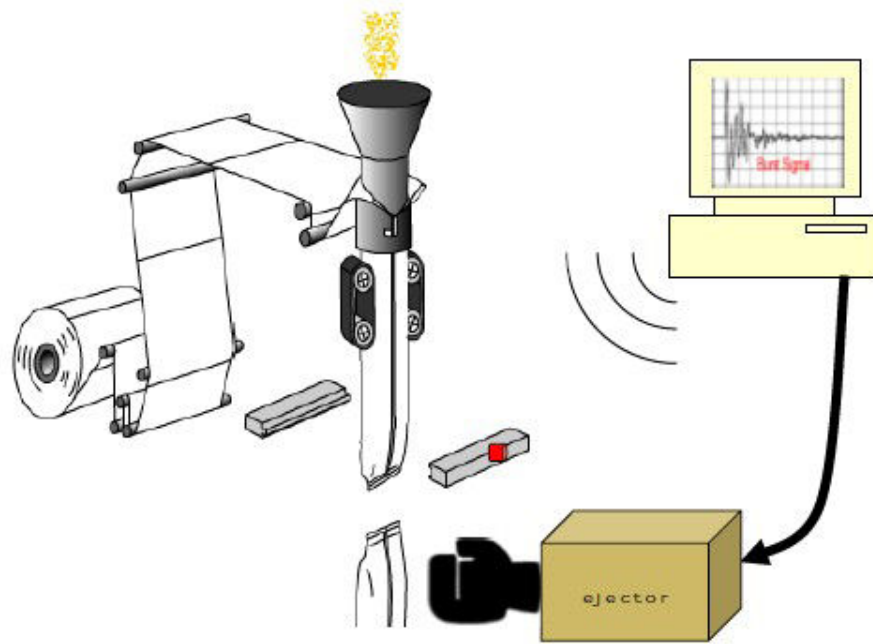


Figure 2: **VFFS machine.** Seal quality monitoring in a packing machine.

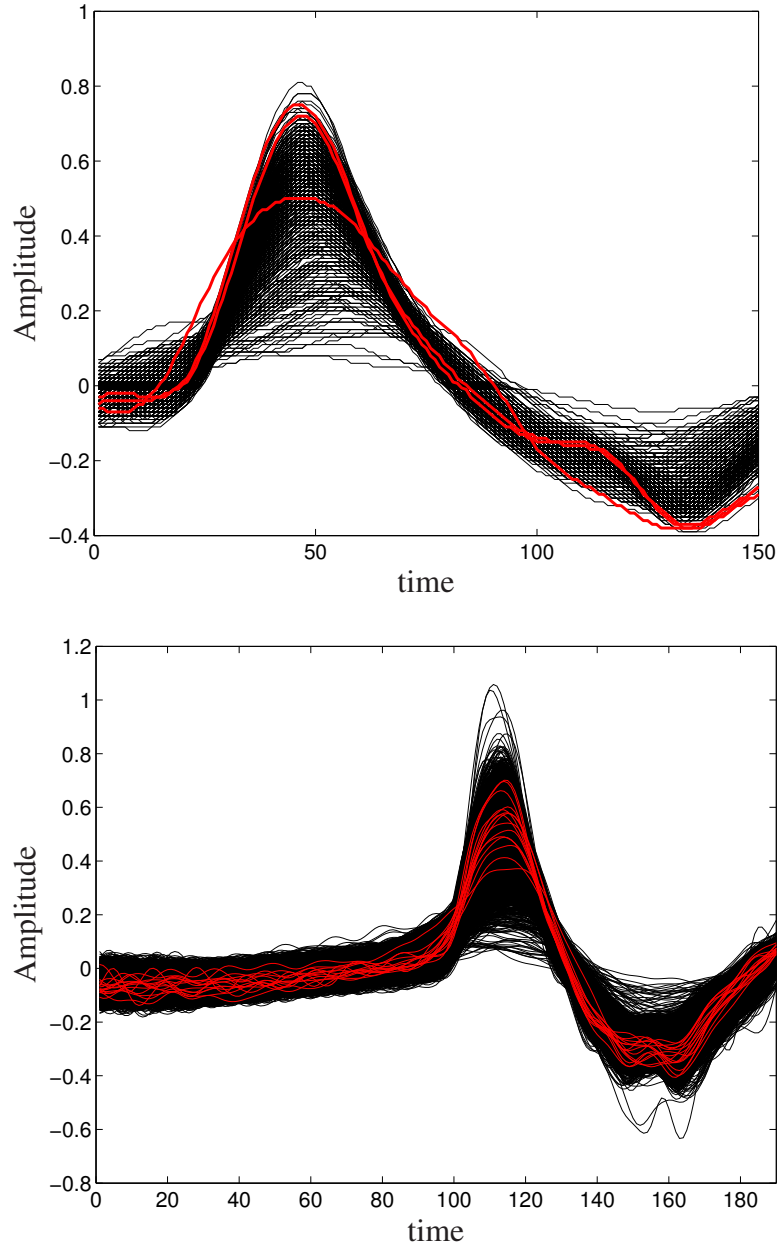


Figure 3: **Data-sets DS_I-DS_II.** **Top:** Accelerometer signals for the data-set DS_I. **Bottom:** Accelerometer signals for the entire data-set DS_II. The signals corresponding to maintenance actions are depicted in red (best visible in colors).

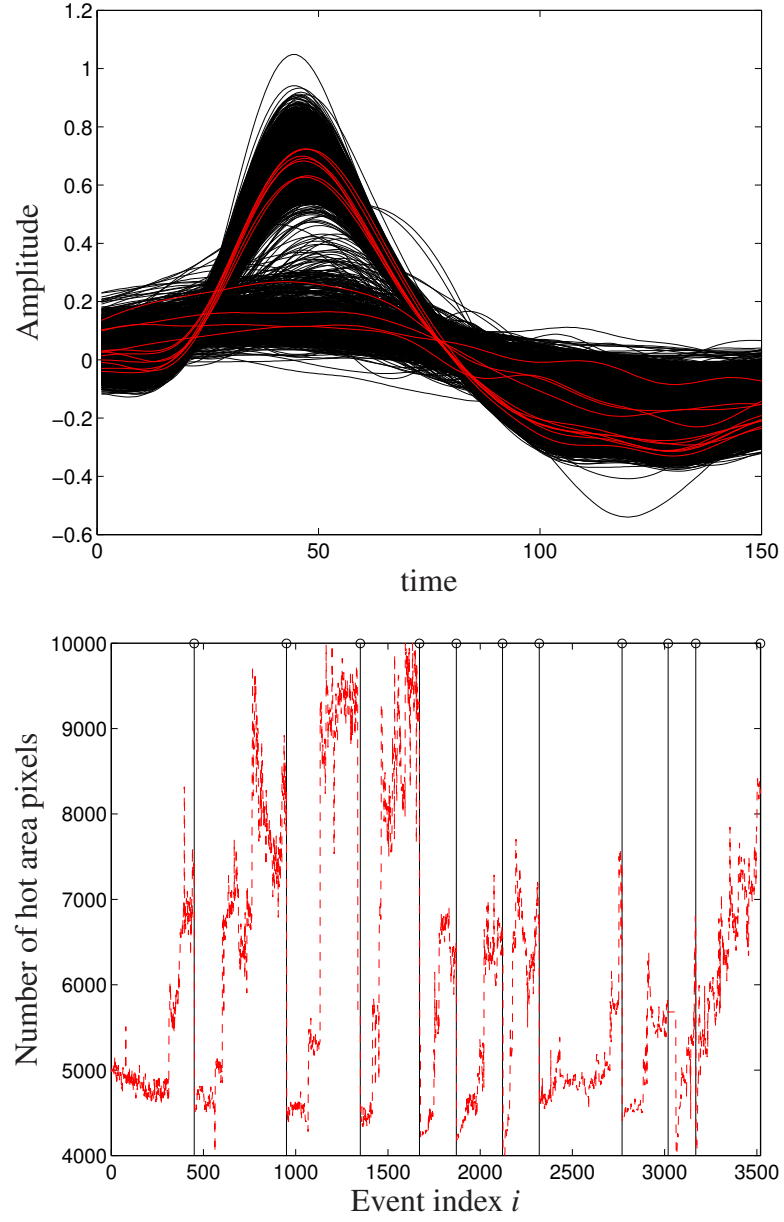


Figure 4: **Data-set DS_III.** **Top:** Accelerometer signals (signals corresponding to maintenance actions are pictured in red). **Bottom:** thermal camera data (the vertical black lines indicates true maintenance actions).

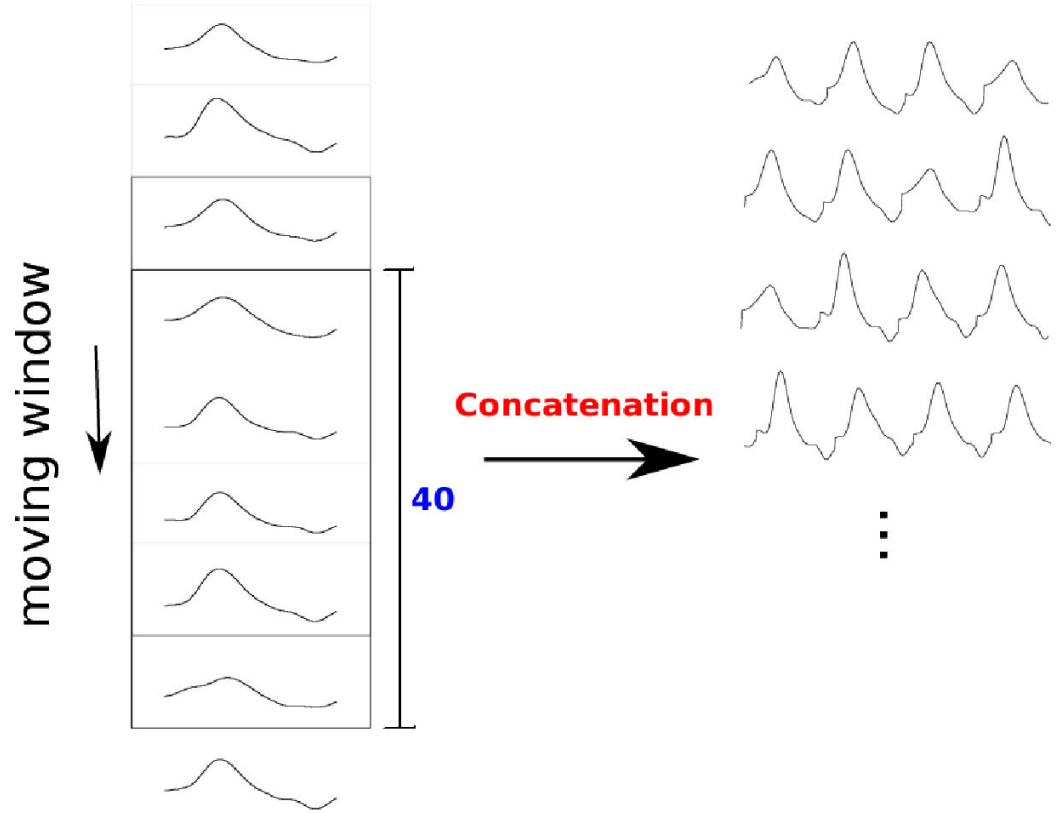


Figure 5: Concatenation of accelerometer signals. After the windowing operation, each data-point is now a time-series of dimension $d = 40 \times 150$ for the first and third data-sets and $d = 40 \times 190$ for data-set DS_II.

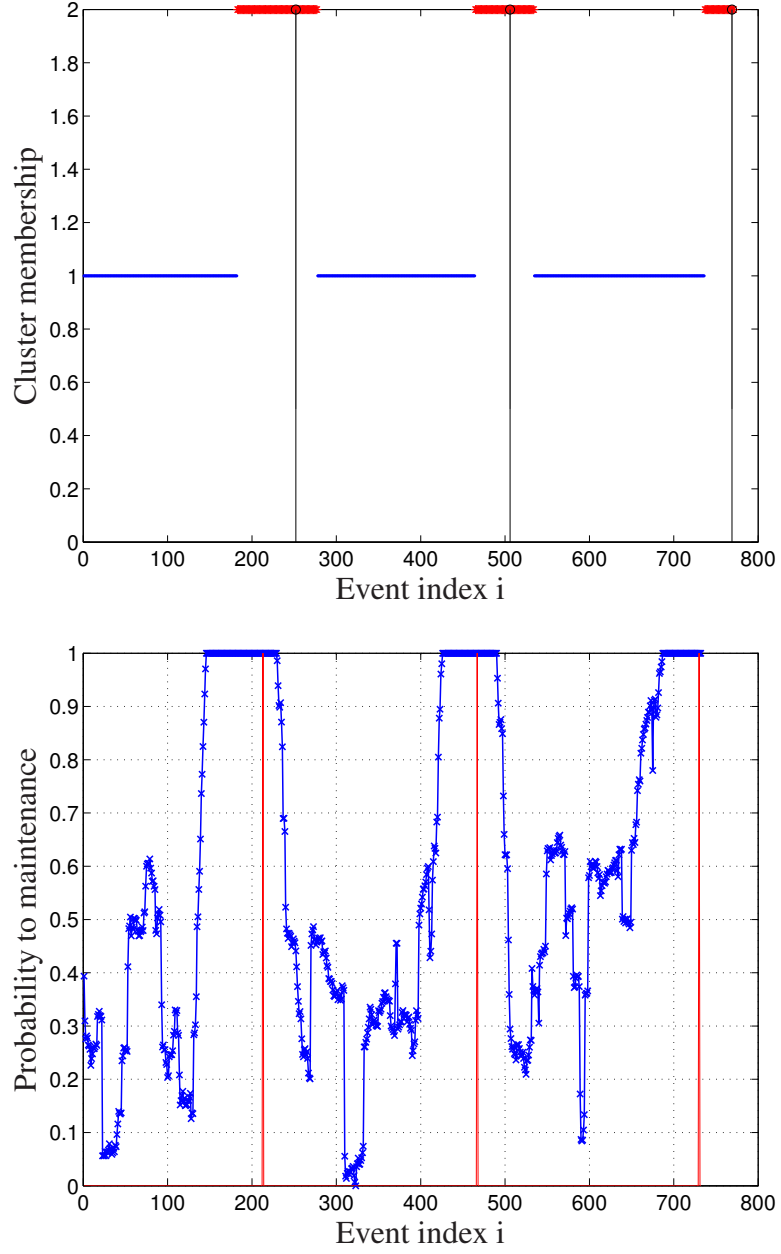


Figure 6: **KSC results data-set DS.I. Top:** Hard clustering results for the whole data-set. Cluster 2 represents predicted maintenance events. The vertical black lines show the true maintenance. **Bottom:** Soft clustering results in terms of probability to maintenance.

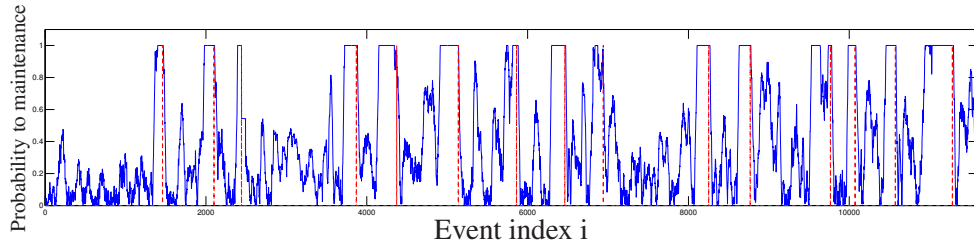
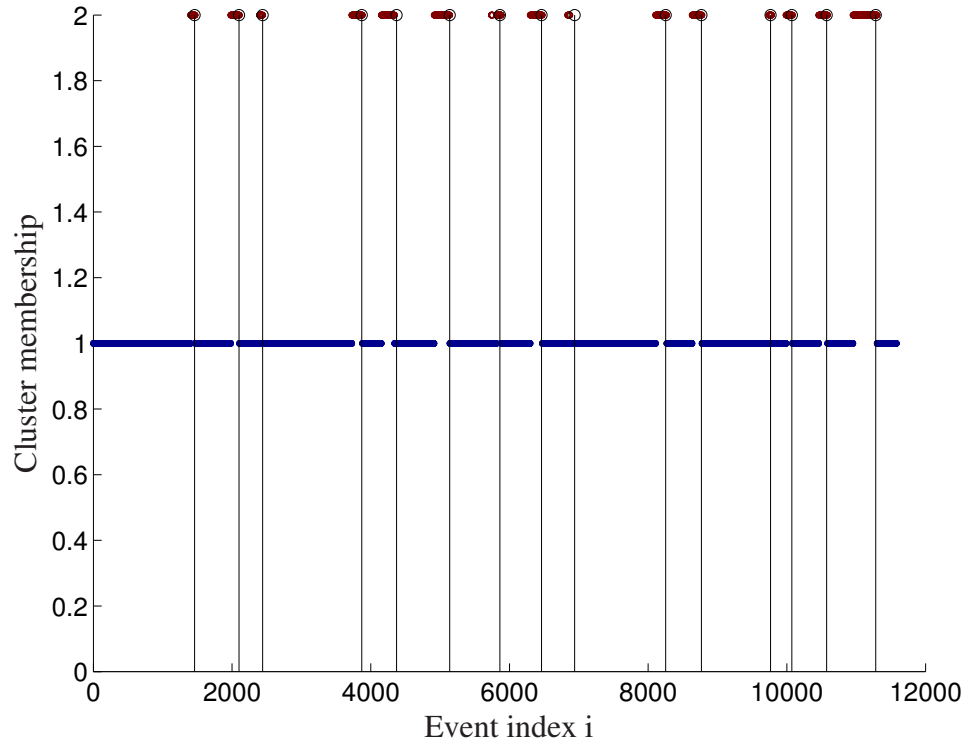


Figure 7: **KSC results data-set DS-II. Top:** Hard clustering, cluster 2 represents predicted maintenance events. **Bottom:** Soft clustering in terms of probability to maintenance.

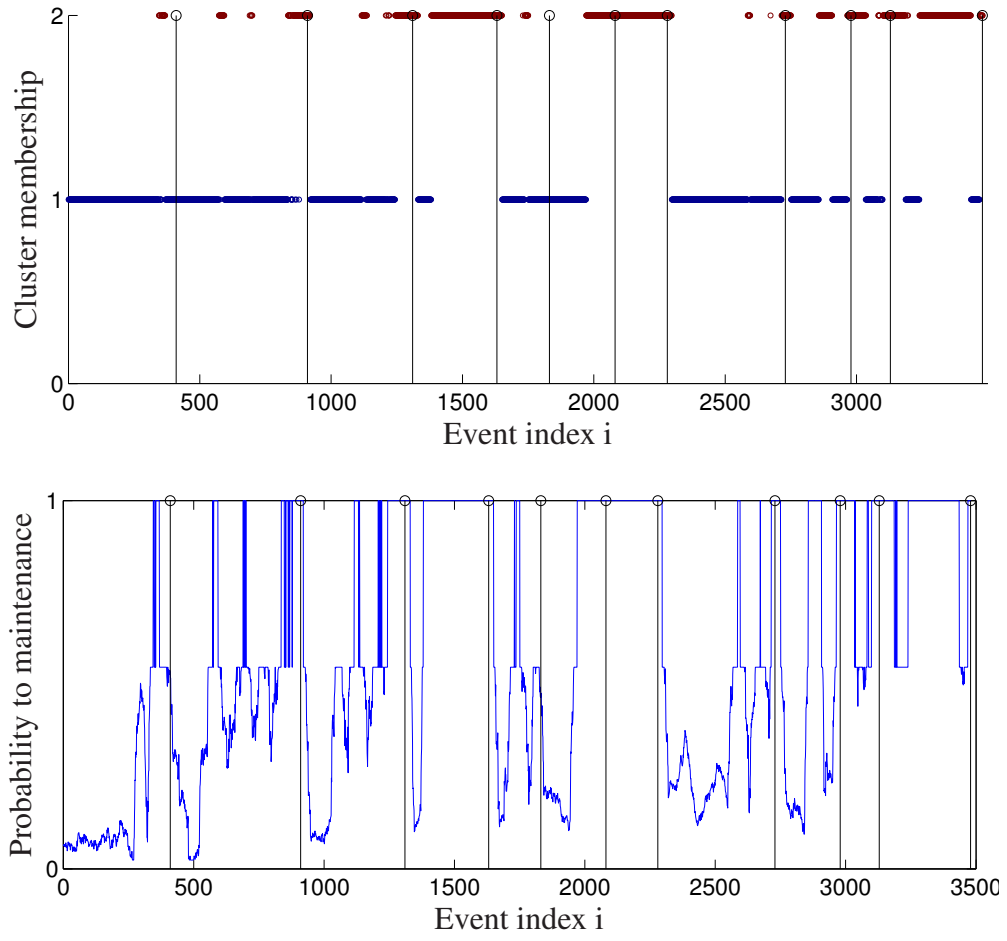


Figure 8: **KSC results data-set DS.III.** **Top:** Hard clustering, cluster 2 represents predicted maintenance events. **Bottom:** Soft clustering in terms of probability to maintenance.

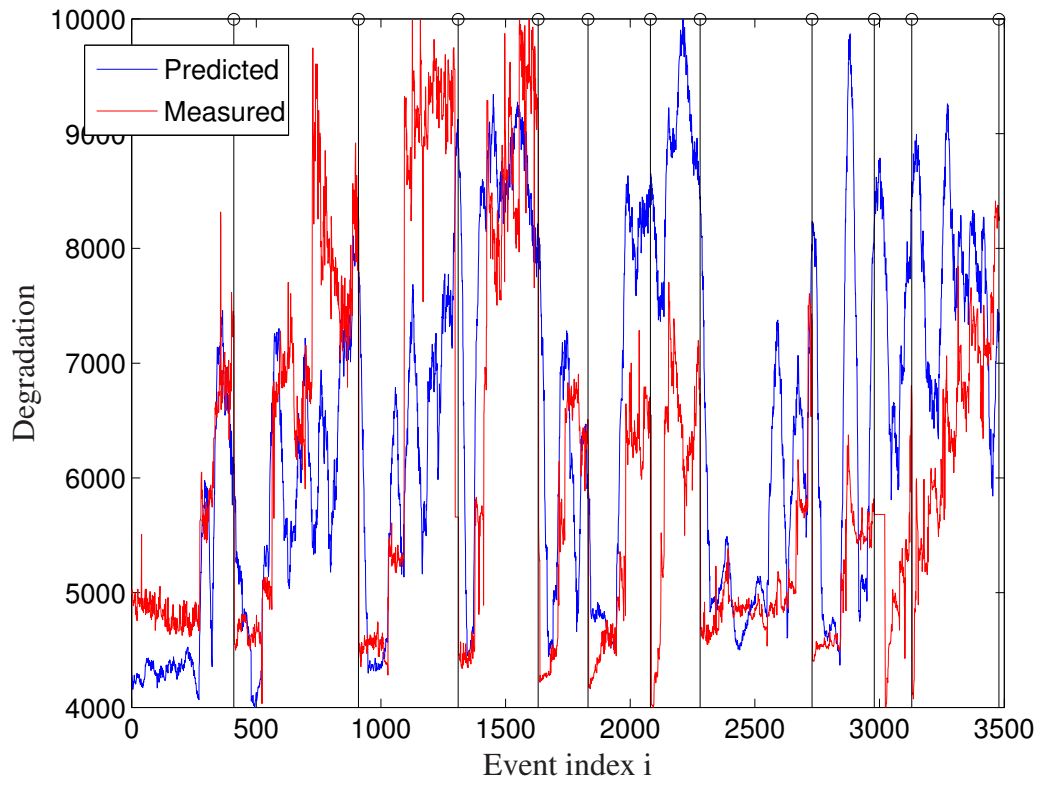


Figure 9: **Degradation dataset DS_III**. Degradation inferred by KSC using as input data the vibration signals (blue) and measured degradation in terms of hot area pixels in the thermal camera images (red).

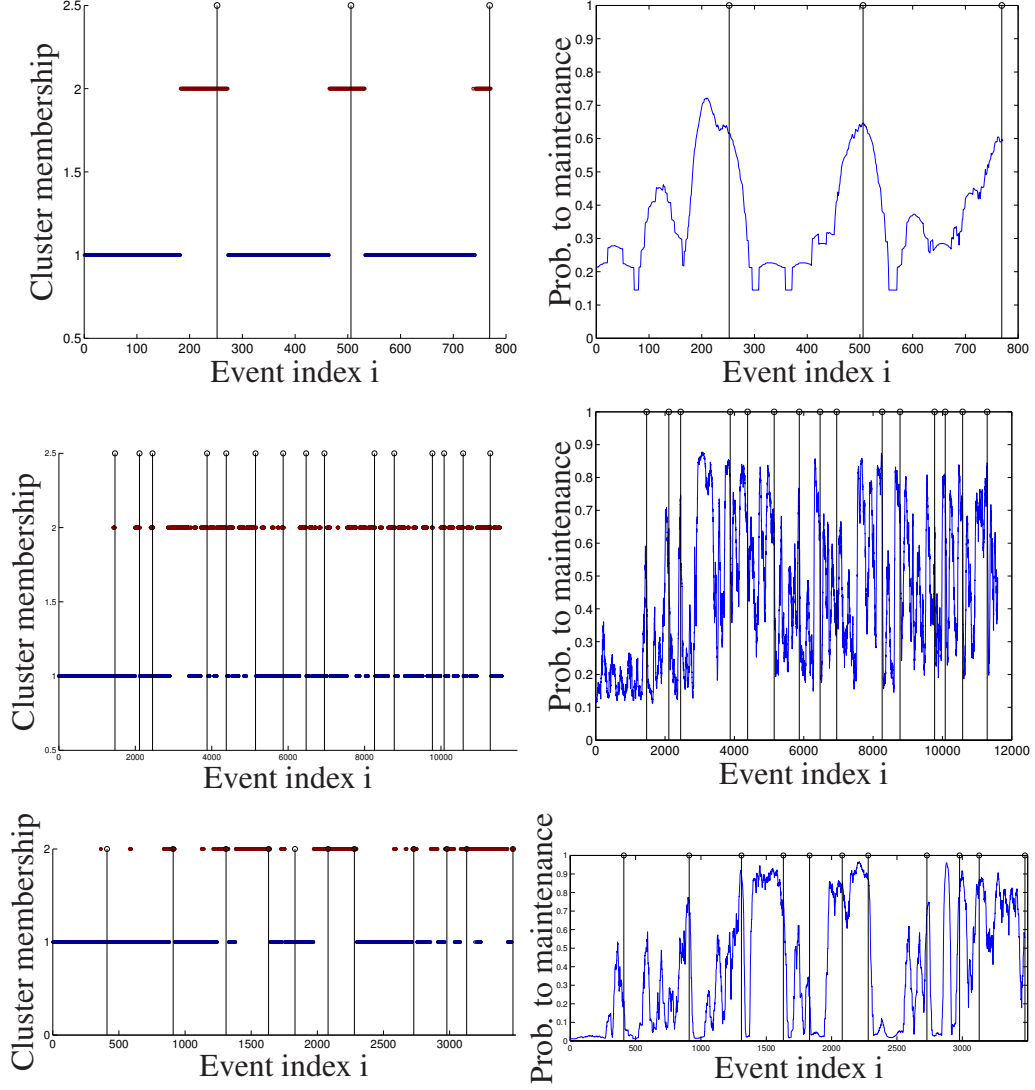


Figure 10: **k-means results.** Cluster 1 symbolizes the normal operating condition, cluster 2 represents predicted maintenance events. The vertical black lines show the true maintenance. **Top:** data-set DS_I. The results are similar to KSC outcomes, even if slightly worse since in the end there is a kind of false alarm (a single prediction of maintenance followed by normal behaviour, before the final maintenance cluster). **Center:** data-set DS_II. In this case k-means performs much worse than KSC, suggesting maintenance in regions not corresponding to actual maintenance events. **Bottom:** data-set DS_III. The results of KSC and k-means are very similar.

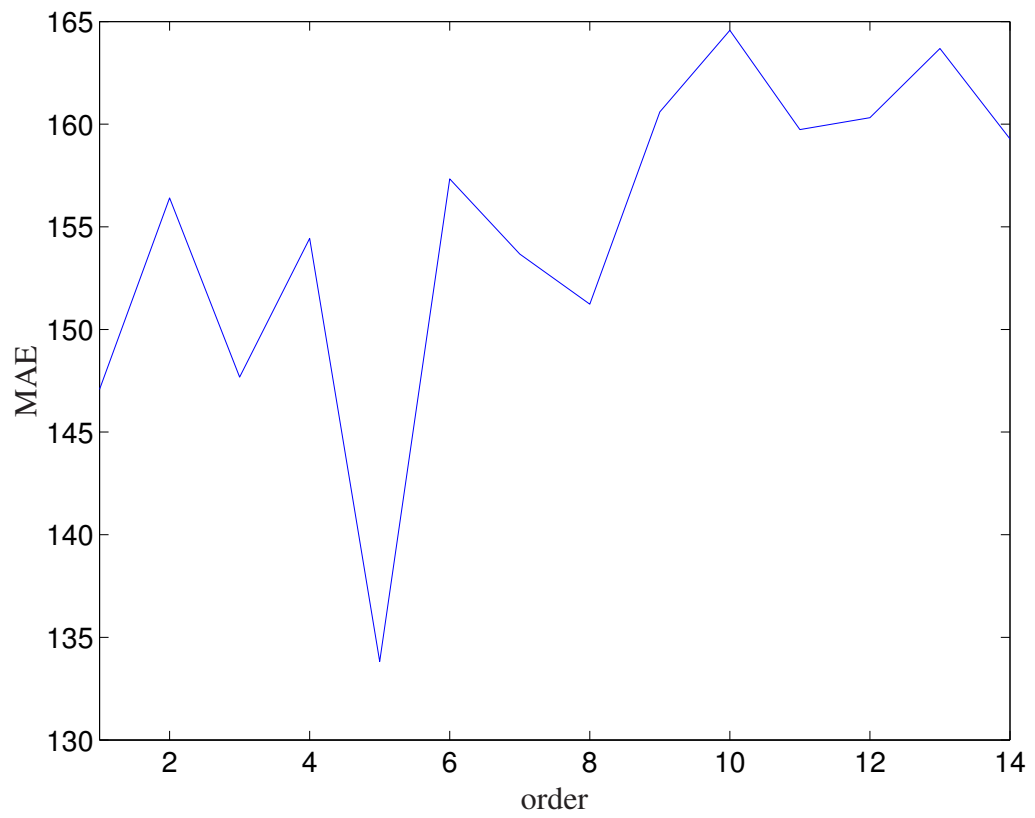


Figure 11: **Tuning order NAR.** Mean absolute error (MAE) between one step ahead prediction and true values with respect to the order of the NAR model. A clear minimum is present at $p = 5$, suggesting to consider this value as the optimal parameter to construct the final NAR model.

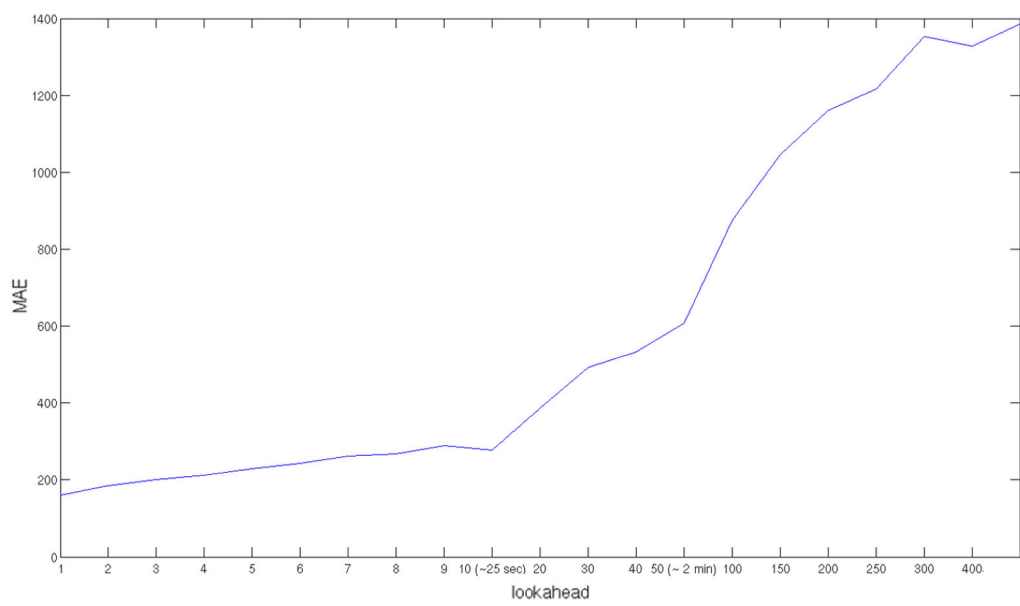


Figure 12: **NAR forecasts versus lookahead.** Prediction performances at the change of the number of steps ahead prediction.

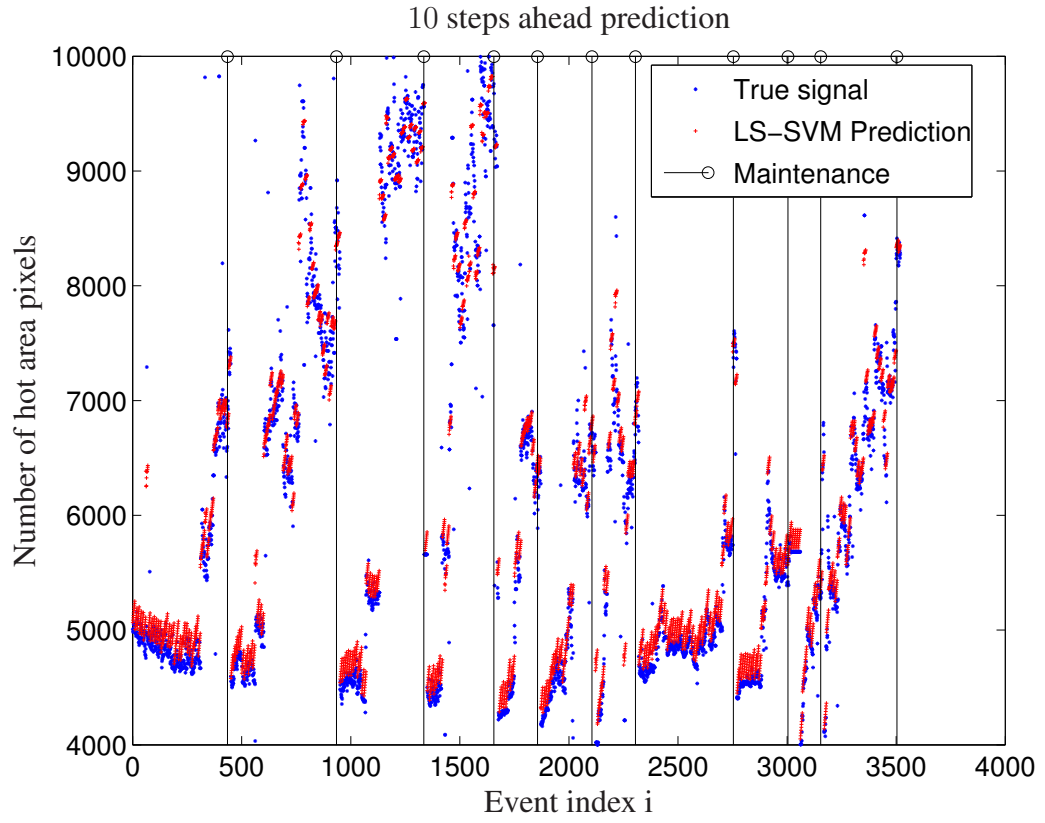


Figure 13: **Prediction NAR.** The NAR model is able to predict very well even 10 steps ahead the future behavior of the thermal camera data based on a window of the 5 past values.