

Accepted Manuscript

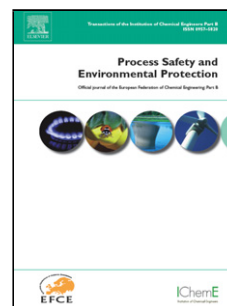
Title: Ozone Measurements Monitoring Using Data-Based Approach

Author: Fouzi Harrou Farid Kadri Sofiane Khadraoui Yin Sun

PII: S0957-5820(16)00020-3
DOI: <http://dx.doi.org/doi:10.1016/j.psep.2016.01.015>
Reference: PSEP 686

To appear in: *Process Safety and Environment Protection*

Received date: 4-5-2015
Revised date: 20-1-2016
Accepted date: 22-1-2016



Please cite this article as: Fouzi Harrou, Farid Kadri, Sofiane Khadraoui, Yin Sun, Ozone Measurements Monitoring Using Data-Based Approach, *Process Safety and Environmental Protection* (2016), <http://dx.doi.org/10.1016/j.psep.2016.01.015>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Ozone Measurements Monitoring Using Data-Based Approach

Fouzi Harrou

CEMSE Division, King Abdullah University of Science and Technology, Thuwal 23955-6900 Saudi Arabia,

Tel.: +966 546326240; fax: +974 012 8080602

E-mail: fouzi.harrou@kaust.edu.sa

Farid Kadri

PIMM Laboratory, UMR CNRS 800, Arts et Métiers ParisTech, Paris, France

E-mail: farid.kadri@ensam.eu

Sofiane Khadraoui

University of Sharjah, Department of Electrical and Computer Engineering, Sharjah, United Arab Emirates

E-mail: sofiane.khadraoui@qatar.tamu.edu

Yin Sun

CEMSE Division, King Abdullah University of Science and Technology, Thuwal 23955-6900 Saudi Arabia,

E-mail: ying.sun@kaust.edu.sa

Ozone Measurements Monitoring Using Data-Based Approach

Highlights:

1. An improved anomaly detection algorithm based on PCA developed
2. The proposed algorithm is applied to monitor ozone measurements
3. The detection results show effectiveness of the proposed method

Accepted Manuscript

Ozone Measurements Monitoring Using Data-Based Approach

Abstract

The complexity of ozone (O_3) formation mechanisms in the troposphere make the fast and accurate modeling of ozone very challenging. In the absence of a process model, principal component analysis (PCA) has been extensively used as a data-based monitoring technique for highly correlated process variables; however conventional PCA-based detection indices often fail to detect small or moderate anomalies. In this work, we propose an innovative method for detecting small anomalies in highly correlated multivariate data. The developed method combine the multivariate exponentially weighted moving average (MEWMA) monitoring scheme with PCA modelling in order to enhance anomaly detection performance. Such a choice is mainly motivated by the greater ability of the MEWMA monitoring scheme to detect small changes in the process mean. The proposed PCA-based MEWMA monitoring scheme is successfully applied to ozone measurements data collected from Upper Normandy region, France, via the network of air quality monitoring stations. The detection results of the proposed method are compared to that declared by Air Normand air monitoring association.

Keywords: Anomaly detection; MEWMA statistic; MSPC; Principal components analysis; Ozone pollution; Data-driven strategy.

1. Introduction

Atmospheric pollution is one of the most serious problems confronting our modern world. The impact of atmospheric pollution on human health is now forefront of population concerns [1]. Numerous epidemiological studies highlight the influence on the health of certain chemical compounds such as sulfur dioxide (SO_2), nitrogen dioxide (NO_2), ozone (O_3) or dust particle in the air [1]. The influence of this pollution is noticeable on sensitive populations such as asthmatics, children, and elderly. Currently, among the monitored compounds, ozone is one of the greatest concern. Ozone is one of the most important photochemical oxidant that exerts adverse effects on human health as well as damages ecosystems, agricultural crops and materials at certain concentration levels [2, 3, 4]. France, like most European countries, has often known during the last summer seasons (2003 especially) episodes of ozone pollution, affecting a large part of the territory. The detection of abnormal pollution in the measured concentrations of these compounds is therefore an important issue for health.

The acceptable concentrations of these pollutants, harmful for human health and the environment, are defined by European standards. Air quality monitoring networks have the following main missions: the measurement network management (recording of pollutant concentrations and a range of meteorological parameters related to pollution

events) and the diffusion of data for permanent information of population and public authorities in reference to norms. The objective of this work is to propose a statistical detection method able to detect abnormal ozone measurements caused by air pollution or any incoherence between the different network sensors or sensor dysfunction. The complexity of ozone (O_3) formation mechanisms in the troposphere [5], the complexity of meteorological conditions in urban areas and the uncertainty in the measurements of all the parameters involved, make the fast and accurate modeling of O_3 very challenging. As an alternative, implicit modelling approaches, which are data-based techniques (like principal component analysis), are particularly well adapted to reveal linear relationships among the process variables without formulating them explicitly. To overcome this difficulty, the principal component analysis (PCA) (a basic method in the framework of multivariate analysis techniques) can be used because they need no prior knowledge about the process model [6]. PCA is one of the most popular multivariate statistical technique used in extracting information from data and is widely used by scientists and engineers in various disciplines, such as in face recognition, data compression, image analysis, visualization, as well as in anomaly detection [7, 8, 9]. In the absence of a process model, principal component analysis (PCA) has been successfully used as a data-based anomaly detection technique for highly correlated process variables [7]. Due to its simplicity and efficiency in processing huge amount of process data, it is recognized as a powerful tool of statistical process monitoring [10, 11]. PCA and its extensions has been successfully applied in a wide range of applications, such as in chemical processes [12], water treatment [13] and hospital management [14].

Generally, in PCA based process monitoring, PCA develop a reference model using the normal data collected from the normal process. The new process behavior can thus be compared with the predefined one by the monitoring system to ensure whether it remain under normal operating conditions or not. When anomaly occur, the process moves out of the normal operation regions indicating that the change in the process behaviors has occurred. Typically, Hotelling T^2 statistic [15] and the sum of squared residuals SPE [16] which is also known as the Q statistic [17] are used in PCA-based monitoring to elucidate the pattern variations in the model and residual subspaces, respectively. The T^2 statistic is defined by the Mahalanobis distance whereas the Q statistic is defined by the Euclidean distance to avoid ill-conditioning due to small eigenvalues [18, 19, 20, 7, 21]. In other words, the T^2 statistic is a measure of the variation in the PCA model and the Q statistic is a measure of the amount of variation not captured by the PCA model. The main disadvantage of using PCs in process monitoring is the lack of physical interpretation [22, 23]. In addition, in a previous study, [17] have shown that the T^2 statistic can result in false negatives (missed detection) due to the latent space sometimes being insensitive to moderate process upsets, which is because each latent variable is a combination of all process variables. Additionally, the disadvantage of T^2 statistic is that anomalies in the process mean that are orthogonal to the first PCs cannot be detected by using the T^2 [24]. The Q statistic, however, is more sensitive to additive anomalies than the T^2 statistic because additive anomalies propagate to the model error. However, the Q statistic can better detect changes in the correlations between the process variables than T^2 [25], and is also more sensitive than T^2 to modeling errors [25]. Nevertheless, the major disadvantage of the conventional

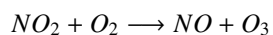
PCA-based detection indices, is that use only the information enclosed about the process in the last observation and they ignore information given by the sequence of all observations. Consequently, this make these detection indices relatively insensitive to small changes in the process variables [26]. These shortcomings of the T^2 and Q statistics motivate the use of other alternatives in order to mitigate these disadvantages. To overcome the previous shortcomings, an alternative approach is proposed in this paper, in which PCA is used as a modeling framework in a model-based anomaly detection method. In this approach, PCA is used to express a process data matrix as the sum of approximate and residual matrices. After a model is obtained using PCA, various methods for anomaly detection can be applied, such as the multivariate exponentially weighted moving average (MEWMA) monitoring scheme, which is utilized in this work to improve anomaly detection. Therefore, the main contribution of the paper is to exploit the greater ability of the MEWMA monitoring scheme to detect small shifts in the process mean for improved anomaly detection of conventional PCA. More specifically, this paper proposes PCA based-MEWMA anomaly detection methodology for detecting abnormal ozone measurements.

The rest of this paper is organized as follows. Section 2 provides a brief overview of ground-level ozone (i.e., tropospheric ozone) pollution. The used data sets and study site are described in Section 3. Then, PCA and a description of how it can be used in anomaly detection is presented in Section 4. Next, the multivariate EWMA which is commonly used in quality control is described in Section 5. Then, the proposed PCA-based MEWMA anomaly detection approach, that integrates PCA modeling and MEWMA monitoring scheme, is presented in Section 6. In Section 7, we present the application of the PCA-based MEWMA anomaly detection approach to detect abnormal ozone measurements of an air quality monitoring network in Upper Normandy, France. Conclusions and future works are finally presented in the last Section.

2. Ozone pollution

Generally, two types of ozone are distinguished: 1) Stratospheric or good ozone, present at around 13 to 30 kilometers of altitude, is a natural filter that protects life on earth from the harmful (ultraviolet) rays of the sun [3]. The ozone hole is a partial disappearance of this filter, linked to the ozone destroying effects of certain pollutants emitted into the troposphere and that move slowly into the stratosphere. 2) Tropospheric ozone or ground-level ozone, present in the air we breathe, is bad: it causes eye irritation, bronchial, and can cause respiratory problems, especially among vulnerable persons (children, elderly) or asthma. The Tropospheric ozone (O_3) is a pollutant that has attracted growing interest in recent years [27, 28]. Unlike other pollutants, ozone is not directly emitted to the atmosphere. It is a pollutant called secondary formed as a result of complex chemical reactions involving two large families of pollutants known as primary: volatile organic compounds (VOC) and industrial emissions release a family of nitrogen oxides (NO_x) [29]. It is formed gradually under the action of solar radiation (NO_x and VOC combine chemically with oxygen to form ozone during sunny) and ozone important peaks can be seen in the summer. High levels of ozone are usually formed in the heat of the afternoon and early evening, dissipating during the cooler nights. The tropospheric

ozone is a pollutant that must be monitored. Ozone, O_3 , is produced by the reaction represented by the following equation:



where NO_2 is the nitrogen dioxide, NO is nitrogen monoxide and O_2 is the oxygen. The nitrogen oxides (NO_2) result from the combination of oxygen (O_2) with nitric oxide (NO) induced by human activities (combustion of hydrocarbons, for transportation or heating...) and volatile organic compounds (VOCs) mainly coming from industries. Solar radiation of wavelengths less than $430nm$ are capable of dissociating NO_2 into a molecule of nitric oxide (NO) and oxygen (O). This last is combined with the oxygen to form the molecule of ozone (O_3).

This reaction provides two essential information: (i) Ozone photochemical pollutant is formed only during daylight hours under appropriate conditions, but is destroyed throughout the day and night. Ozone concentrations are higher on hot, sunny, calm days. Generally, ozone concentration are highest in the rural sites than the urban sites. Higher concentrations in rural areas can be result from nitrogen oxides and volatile organic compounds being transported from upwind urban or industrial areas, by natural ozone being transported to ground-level from the upper atmosphere, or from natural volatile organic compounds emitted from vegetation [30, 31]. (ii) At night, ozone produced in the light of day (due direct solar radiation), disappears. This is due to the destruction of ozone by nitric oxide, which is emitted by vehicles. Nitric oxide can remove ozone by reacting with it to form nitrogen dioxide ($3NO + O_3 \rightarrow 3NO_2$). Ironically, the concentrations of nitric oxide are very low in most rural areas to completely destroy ozone, so ozone remains in the atmosphere for a longer period. Ozone levels tend to be higher in rural areas where there are less local emissions of nitrogen dioxides to destroy any ozone that has formed in the atmosphere [32].

2.1. Diurnal variation of ground-level ozone

Diurnal variations of ozone concentrations follow a typical cycle, with a minimum in late night and a maximum around mid afternoon [33], as shown in Figure 1. This figure shows the measurements of 7 different stations (located in the same network) for the same day. The 7 curves have a daily behavior very similar. The ozone concentration begins to increase just after sunrise, and attains its maximum level in the afternoon due to photochemical production of O_3 mainly from oxidation of natural and anthropogenic hydrocarbons, carbon monoxide (CO), and methane (CH_4) by hydroxyl (OH) radical in the presence of a sufficient amount of NO_x .

2.2. Anomalies in ozone measurements

Two types of anomalies in ozone measurements (atypical ozone peaks) can be distinguished: true and false anomalies. True anomalies correspond to peaks in the ozone levels due to the production of photochemical ozone. The formation of a true peak of ozone requires certain conditions, such as sunny days under stagnant and humid air conditions, high humidity and high temperatures to promote the formation of ozone, and low wind speeds to accumulate

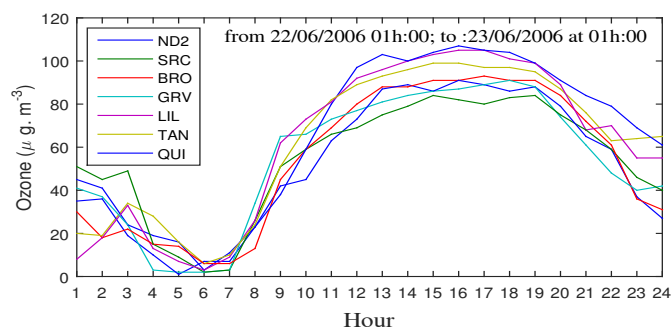


Figure 1: Example of daily ozone concentrations.

high pollution levels. These peaks are usually large with a duration of several hours (due to long reaction times needed for a gradual formation of the photochemical ozone). Therefore, this type of anomalies usually exhibit bell shaped curves. Furthermore, false anomaly are usually observed outside the summer period, where the ozone concentration abruptly increases with very high ozone concentrations (to be in the range of $150\mu\text{g}/\text{m}^3$ to $600\mu\text{g}/\text{m}^3$) for short periods of time (around one hour). These abnormal measurements are sharply pointed, which are different from those observed in the case of photochemical ozone. The presence of this type of anomalies can be due to different phenomena: (a) malfunctioning sensor(s), (b) transported ozone produced elsewhere in the region, (c) transported ozone produced elsewhere in the region, and others [34].

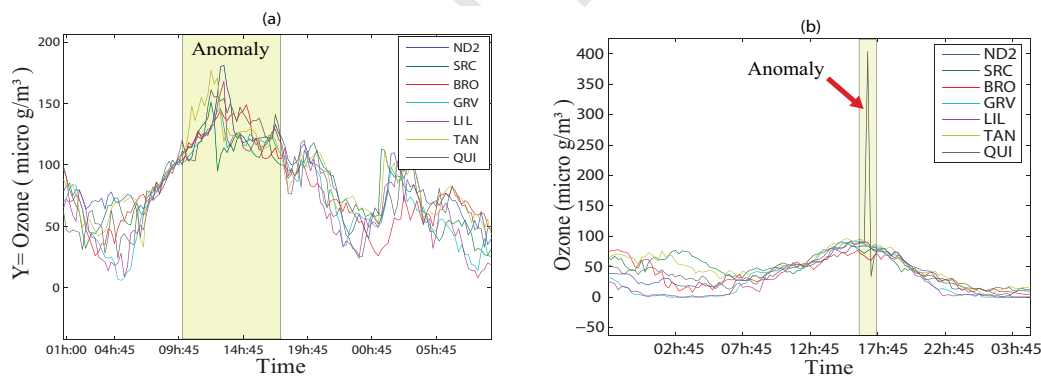


Figure 2: Types of ozone anomalies: (a) true anomaly, (b) false anomaly.

3. Air quality monitoring in French using network of measurement stations

Pollution of the lower atmosphere by ozone is a growing problem in industrialized countries. In France, the law on air quality and rational use of energy (LAURE, law n° 96-1236, 30th December 1996) provides a set of measures to guarantee for citizens the best air quality. Hence the fight against air pollution become a priority. Today, according to this law, all cities in France, with more than 100 000 inhabitants, have an air quality monitoring network. Actually in France, we have 40 networks where each of them is managed by a local association. Fourteen air quality monitoring

associations (AASQA) have been created and approved by the ministry of environment to monitor air quality in France. Atmo federation groups all these forty approved associations. These associations measure, collect, monitor and observe air quality. AASQA continuously monitor the presence in the ambient air of 13 pollutants regulated by European directives and national legislation. Ozone (O_3) forms part of the pollutants which are measured by monitoring associations, because it can cause a number of respiratory health effects. Monitoring networks for air quality generally consist of several measuring stations spread over the geographical area concerned. When the air pollutant concentration exceed a certain threshold (defined by decree in air quality regulations) or there is a risk to exceed it, the association is in charge to inform general public with information on the measured values and to give advices/recommendations for the exposed populations.

The heat wave of summer 2003 in France, was linked with an exceptional ozone pollution, that affected the whole European community. These levels were specially high and related to the weather conditions and exceptional temperatures. The consequences of this heat wave demonstrated the importance to dispose of reliable warning systems for detection of unexpected pollution and unforeseeable events. Considerable efforts have been deployed (and still are) to equip AASQA by descriptive models of ozone dispersion. However, we can notice that these so-called deterministic models are sometimes far from reality. Hence, it is important to propose new optimal descriptive models and statistical methodology for the detection of peak ozone levels. This will be the principal objective of this study. In the next subsection the ozone data set used in this study will be briefly described.

3.1. Data sets and study region

In this study, the Upper Normandy region was selected for data collection. Upper Normandy is located at northwest of Paris, near the south side of Manche sea and is one of the most highly industrialized areas in France. This city, like most large European cities, faces air pollution problems. The association Air Normand is the official association responsible for monitoring air quality over Upper Normandy region, and providing with information on the results.

Generally, there are different types of air quality monitoring stations: local, urban, rural and industrial. The local stations, directly exposed to industrial locations or positioned close to traffic, convey the concentration of pollutants emanating from an identified source. The urban stations measure the ambient air pollution to which the majority of the population is exposed. Finally, the rural stations are representative of the levels observed in the sparsely populated areas and enable the long distance consequences to be assessed. Each station consists of a set of sensors, dedicated to the acquisition of pollutants (ozone O_3 , nitrogen oxides NO, sulfur dioxide SO_2 ,...).

In order to measure and control tropospheric ozone pollution the Air Normand association consists of 7 stations placed in industrial, peri-urban and urban sites, across the region. Ozone concentrations have been measured every fifteen minutes by Air Normand network. Figure 3 shows a map of France and the location of study sites (Champagne-Ardenne and Upper normandy).

The aim of this study is to apply the proposed PCA-based MEWMA anomaly detection algorithm in order to



Figure 3: Location of study sites in France (<http://education.francetv.fr/CartesInteractives>).

detect abnormal measurements of ozone, both of anthropogenic origin (pollution peaks caused by human activity) or the result of dysfunction of sensors (anomalies, interference,...). A brief introduction to the principles of PCA, and how it can be used in anomaly detection is presented next.

4. Principal component analysis (PCA)

PCA is a linear dimensionality reduction modeling method, which can be helpful when handling data with a high degree of cross correlation among the variables. The main idea behind PCA is briefly introduced in this section, and more details can be found in [35, 20].

4.1. PCA modeling

Let us consider the following raw data matrix $\mathbf{X} = [\mathbf{x}_1^T, \dots, \mathbf{x}_n^T]^T \in R^{n \times m}$ consisting of n observations and m correlated variables. The data are collected when the monitored process is under normal operating condition so that the PCA's model that will be built represents a reference of the normal process behavior. Before computing the PCA model, the raw data matrix \mathbf{X} is usually pre-processed by scaling every variable to have zero mean and unit variance. This is because variables are measured with various means and standard deviations in different units. This pre-processing step puts all variables on an equal basis for analysis [36]. Let \mathbf{X}_s denote the autoscaled matrix of \mathbf{X} . By using singular value decomposition (SVD), PCA transforms the data matrix \mathbf{X}_s into a new matrix $\mathbf{T} = [t_1 \ t_2 \ \dots \ t_m] \in R^{n \times m}$ of uncorrelated variable called score or principal components (PCs). Indeed, PCs are just mathematical constructs chosen to represent the variance as efficiently as possible, even if their physical meaning is obscure. Each principal component is a linear combination of the original variables, so that \mathbf{T} is obtained from \mathbf{X}_s by

an orthogonal transformations (rotations) designed by $\mathbf{P} = [p_1 \ p_2 \ \dots \ p_m] \in R^{m \times m}$ which is given as following:

$$\mathbf{T} = \mathbf{X}_s \mathbf{P} \quad \text{and} \quad \mathbf{X}_s = \mathbf{T} \mathbf{P}^T = \sum_{i=1}^m t_i p_i^T. \quad (1)$$

where the column vectors $p_i \in R^m$ of the matrix $\mathbf{P} \in R^{m \times m}$ (also known as the loading vectors) are formed by the eigenvectors associated with the covariance matrix of \mathbf{X}_s , i.e., Σ . The covariance matrix, Σ , is defined as follows:

$$\Sigma = \frac{1}{n-1} \mathbf{X}_s^T \mathbf{X}_s = \mathbf{P} \Lambda \mathbf{P}^T \quad \text{with} \quad \mathbf{P} \mathbf{P}^T = \mathbf{P}^T \mathbf{P} = \mathbf{I}_n, \quad (2)$$

where, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ is a diagonal matrix containing the eigenvalues in a decreasing order ($\lambda_1 > \lambda_2 > \dots > \lambda_m$), \mathbf{I}_n is the identity matrix, and the i^{th} eigenvalue equals the square of the i^{th} singular value (i.e. $\lambda_i = \sigma_i^2$) [37].

Note that the PCA model results in the same number of principal components as the number of original variables (m). In the case of collinear process variables, however, a smaller number of principal components (l) are needed to capture most of the variations in the data. Often, a small subset of the principal components (corresponding to the largest eigenvalues) can extract most of the important information in a data set, and thus simplify its analysis. The first PC indicates the direction of largest variation in data, the second PC indicates the largest variation unexplained by the first PC in a direction orthogonal to the first PC, and so on. Figure 4 shows how a 3-dimensional collinear data set can be represented in a reduced 2-dimensional space using only 2 principal components. The number of the retained PCs is usually less than the number of measured variables.

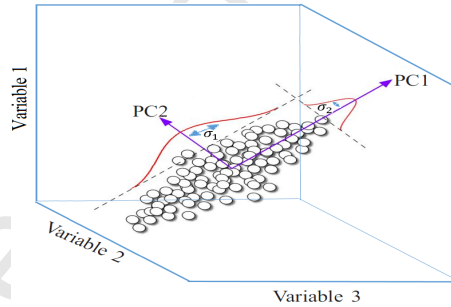


Figure 4: Principle of PCA.

A key step in the building of PCA model is to determine the number of PCs, l , that are required to adequately capture the major variability in the data sets. The goodness of the PCA model depends on a good choice of how many PCs are retained [38]. The first (l) largest principal components normally describe the most of the variance of the data. On the other hand, the smallest principal components are considered as a noise contributor. Too few components implies that there are not enough dimensions to represent the process variability, which degrades the prediction quality of the PCA model. While too many components implies that one can introduce noise and the model fails to capture some of the information. A number of techniques have been proposed to determine the number of PCs to be retained in a PCA model including cross validation [39], Scree plot [40], and cumulative percent variance (CPV). In this study,

the CPV technique will be used to determine the number of PCs for PCA model. The CPV is defined as follows: $CPV(l) = \frac{\sum_{i=1}^l \lambda_i}{\text{trace}(\Sigma)} \times 100$. Once the number of principal components l is determined, the PCA algorithm decomposes \mathbf{X}_s into two orthogonal parts: an approximated data matrix $\widehat{\mathbf{X}}$ and a residual data matrix \mathbf{E} , i.e.,

$$\mathbf{X}_s = \sum_{i=1}^l t_i p_i^T + \sum_{i=l+1}^m t_i p_i^T = \widehat{\mathbf{X}} + \mathbf{E} \quad (3)$$

Of course, if some of the variables in the data set are collinear or highly correlated, then a smaller number of principal components l are required to explain the majority of the variance in the data. In practice, the variance left unexplained by the PCs is captured by the residual subspace, which are often associated with the instrument or process noise.

4.2. PCA-based detection indices

As shown in equation (3), any measured vector \mathbf{x} can be expressed using PCA as the sum of two orthogonal parts, approximated vector $\widehat{\mathbf{x}}$ and residual vector \mathbf{e} (see Figure 5), corresponding to the projection onto the PC subspace S_p and residual subspace S_r , respectively. In anomaly-detection using PCA, a PCA model is constructed using fault-free data, and then the model is used to detect faults using one of the detection indices, such as the Hotelling's T^2 and Q statistics, which are described next.

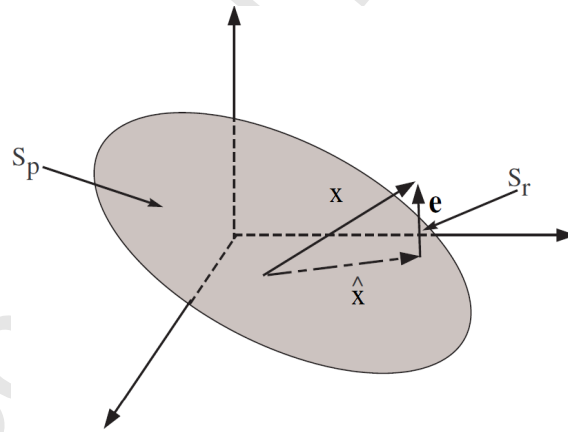


Figure 5: Geometric principle of PCA

4.2.1. Hotelling's T^2 statistic

The T^2 statistic measures the variations in the principal components or score vectors at different time samples. T^2 at any instance of time is defined as follows [15]:

$$T^2 = \mathbf{x}^T \widehat{\mathbf{P}} \widehat{\Lambda}^{-1} \widehat{\mathbf{P}}^T \mathbf{x} = \sum_{i=1}^l \frac{t_i^2}{\lambda_i} \quad (4)$$

where the matrix $\widehat{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_l)$, is a diagonal matrix containing the eigenvalues associated with the l retained principal components. The threshold value used for the T^2 statistic can be computed as follows [15]:

$$T_{l,n,\alpha}^2 = \frac{l(n-1)}{n-l} F_{l,n-l,\alpha} \quad (5)$$

where n is the number of samples in the data, l is the number of retained PCs, α is the level of significance (α usually takes values between 1% and 5%), and $F_{l,n-l}$ is the Fisher F distribution with l and $n-l$ degrees of freedom. When the number of observations, n , is rather large, the T^2 statistic threshold can be approximated with a χ^2 distribution with l degrees of freedom, i.e., $T_{l,n,\alpha}^2 = \chi_{l,\alpha}^2$. These threshold values are computed using fault-free data. For new testing data, when the value of T^2 exceeds the value of the threshold, $T_{l,n,\alpha}^2$ or T_{α}^2 , a fault is declared.

4.2.2. Q statistic or squared prediction error (SPE)

The Q statistic or Rao-statistic (also referred to as the squared prediction error, SPE) measures the projection of a data sample on the residual subspace, which provides an overall measure of how a data sample fits the PCA model. Q is defined as the sum of squares of the residuals obtained from the PCA model, i.e., [7]:

$$Q = \mathbf{e}^T \mathbf{e} \quad (6)$$

The upper control limit of this statistic is defined as [37]:

$$Q_{\alpha} = \varphi_1 \left[\frac{h_0 c_{\alpha} \sqrt{2\varphi_2}}{\varphi_1} + 1 + \frac{\varphi_2 h_0 (h_0 - 1)}{\varphi_1^2} \right] \quad (7)$$

where, c_{α} is the value of the normal distribution with α level of significance, $\varphi_i = \sum_{j=l+1}^m \lambda_j^i$ for $i = 1, 2, 3$, and $h_0 = 1 - \frac{2\varphi_1\varphi_3}{3\varphi_2^2}$. This value of threshold is calculated based on the assumptions that the measurements are time-independent and multivariate normally distributed. The Q fault detection index is very sensitive to modeling errors and its performance largely depends on the choice of the number of retained principal components, l , [7]. The PCA fault detection algorithm is summarized next.

1) **Given:**

- A training fault-free data set that represents the normal process operations and a testing data set (possibly faulty data),

2) **Data preprocessing**

- Scale the data to zero mean and unit variance,

3) **Build the PCA model using the training fault-free data**

- Compute the covariance matrix, Σ , using equation (2),
- Calculate the eigenvalues and eigenvectors of Σ and sort the eigenvalues in decreasing order,

- Determine how many principal components to be used. Many techniques can be used in this regards. In this work, the *CVP* criterion is used,
- Express the data matrix as a sum of approximate and residual matrices as shown in equation (3),
- Compute the control limits for the statistical model (e.g., the Q_α statistic limits)

4) Test the new data

- Scale the new data,
- Generate a residual vector, \mathbf{e} , using PCA,
- Compute the monitoring statistic (Q or T^2 statistics) for the new data using equation (4) or (6),

5) Check for anomalies

- Declare an anomaly when new data exceeds the control limits (e.g., $Q \geq Q_\alpha$).

Unfortunately, the T^2 and Q statistics use only the observed data at the current time point alone for making decision about the process performance at the current time point. They take into account only the present information of the process thus they have a short memory. For this reason the T^2 and Q statistics are also called detection indices without memory. Consequently, these detection indices are relatively insensitive to small changes in the process variables, and thus may result in missed detections [26]. These drawback of the T^2 and Q statistics motivate the use of other alternatives in order to surmount these disadvantages. Note that the ability to detect smaller parameter shifts can be improved by using a chart based on a statistic that corporate information from past samples in addition to current samples. In this study, anomaly detection technique which is based on PCA model and MEWMA control scheme will be developed in order to surmount these drawbacks and improve detection performance compared to the conventional PCA based anomaly detection method. A succinct introduction to the basic ideas behind MEWMA monitoring scheme is exposed in the subsequent section.

5. Multivariate EWMA statistical control scheme

Control charts are one of the most frequently used procedures in statistical process control (SPC), and have been widely used as a monitoring tool in quality engineering to detect the existence of possible anomalies in the mean or variance of process measurements. Many control charts are referenced in the bibliography, and they can be broadly categorized into main classes: univariate and multivariate techniques [26, 41]. The univariate control charts such as Shewhart, cumulative summation (CUSUM) [42], and EMWA [26] have been designed to essentially to monitor only one process variable. However, modern industrial processes often present a large number of highly correlated process variables. This is the area where univariate control charts are unable to explain different aspects of the process and, therefore, it is not appropriate for modern day processes. Moreover, to monitor several different process variables in the same time multivariate statistical monitoring charts such as Multivariate Shewhart [26], Multivariate EWMA

(MEWMA) [43] and Multivariate CUSUM (MCUSUM) [26] were developed in analogy with the univariate charts. In fact, most commonly used multivariate control charts are the natural extension of the univariate charts, e.g. the Hotelling's T^2 charts [44], MEWMA charts and MCUSUM charts [26, 43]. A multivariate SPC charts take into account the additional information due to the correlation between a process variables while univariate SPC charts do not. These concepts may be used to develop more efficient control charts than the simultaneous operation of several univariate control charts.

The MEWMA chart was first proposed by Lowry et al [43] to monitor mean shifts of a multivariate process. This is a multivariate extension of the univariate EWMA chart proposed by Roberts [45]. This monitoring chart is constructed based on a weighted moving average of all observed data and available at the current time point. The MEWMA is utilized when there are several correlated process variables to be monitored simultaneously where detecting faults with small magnitudes is of interest. Suppose that we observe $\mathbf{X}_t = (X_1, X_2, \dots, X_m)^T$, a m -dimensional set of observations at time t . A MEWMA control chart is proposed by Lowry et al [43] as follows:

$$\mathbf{Z}_t = \mathbf{R}\mathbf{X}_t + (\mathbf{I}_{m \times m} - \mathbf{R})\mathbf{Z}_{t-1}. \quad (8)$$

Where $\mathbf{R} = \text{diag}(r_1, r_2, \dots, r_m)$ which is a diagonal matrix with r_1, r_2, \dots, r_m on the main diagonal, and m is the number of variables; $0 < r_j \leq 1$ is a weighting parameter for j -th component of \mathbf{X} , for $j = 1, 2, \dots, m$, $\mathbf{I}_{m \times m}$ is the identity matrix, \mathbf{Z}_i is the i^{th} EWMA vector, and X_i is the i^{th} observation vector $i = 1, 2, \dots, n$. The initial value \mathbf{Z}_0 is usually obtained as equal to the in-control mean vector of the process. Generally, in quality control, a smaller value of r leads to quicker detection of smaller shifts [46]. Indeed, r should be adjusted to a value appropriate for the characteristic of the monitored process. Usually, the larger the shift is, the greater the, r is. The value of r is usually set between 0.2 and 0.3 [47]. It can be noticed that if $\mathbf{R} = \mathbf{I}$, then the MEWMA control chart is equivalent to the T^2 Chart. In this case, a MEWMA chart has been automatically changed into T^2 chart.

In practice, if there is no priori reason to weight different components differently, then we can simply choose $r_1 = r_2 = \dots = r_m = r$. In this case the equation 8 can be written as follows:

$$\mathbf{Z}_t = r\mathbf{X}_t + (1 - r)\mathbf{Z}_{t-1}. \quad (9)$$

The MEWMA decision function, \mathbf{V}_t^2 , can be calculated recursively as follows [43]:

$$\mathbf{V}_t^2 = \mathbf{Z}_t^T \boldsymbol{\Sigma}_{\mathbf{Z}_t}^{-1} \mathbf{Z}_t \quad (10)$$

where $\boldsymbol{\Sigma}_{\mathbf{Z}_t}$ is the variance-covariance matrix of \mathbf{Z}_t . When $r_1 = r_2 = \dots = r_p = r$, the variance-covariance matrix of \mathbf{Z}_t can be simplified to:

$$\boldsymbol{\Sigma}_{\mathbf{Z}_t} = \frac{r}{(2-r)} \left[1 - (1-r)^{2n} \right] \boldsymbol{\Sigma} \quad (11)$$

where Σ is the covariance matrix of the input data. The MEWMA chart statistic is usually constructed in terms of the asymptotic covariance matrix. When t becomes large, the covariance matrix converges to: $\Sigma_{Z_t} = (\frac{r}{2-r})\Sigma$.

Under nominal conditions, the statistic Z is distributed according to the gaussian law with zero mean and variance-covariance matrix Σ_{Z_t} , $Z \sim \mathcal{N}(0, \Sigma_{Z_t})$. The distribution of the statistic Z in the presence of additive mean shift μ_1 is given as: $Z \sim \mathcal{N}(r \sum_{j=1}^n [(1-r)^{n-j}\theta], \Sigma_{Z_t})$. The MEWMA chart declares the presence of anomaly when $V_t^2 > h$, where h is the control limit. The distribution of V_t^2 under in-control condition is χ_p^2 . However, because the variables in the time series V_t^2 , $t = 1, 2, \dots$ are correlated, the control limit h cannot simply be chosen to be $(1 - \alpha)$ -th quantile $\chi_{1-\alpha, p}^2$ of the χ_p^2 distribution. One of the main troubles on this chart is the selection of the h . The value of h can be calculated by simulation to achieve a specific control limits. Various authors have used theoretical derivation, Markov chain approximation, integral equation approximation, and monte Carlo simulation, or combinations of the three techniques to compute the control limit h according to the parameters r , p , and α [48, 49]. Bodden [50] proposed an algorithm to find the control limit h in order to respect a given number of false alarm and a given r .

6. Anomaly detection using a PCA-based MEWMA control scheme

In this section, PCA is integrated with MEWMA to develop a new anomaly detection scheme with a higher sensitivity to small or moderate anomalies in the data. Towards this end, PCA is used to represent a matrix of the process measurements as the sum of two orthogonal parts (an approximated data matrix and a residual data matrix) as shown in equation (3). In PCA model, the principal components associated with large eigenvalues capture most of the variations in the data, where, ones associated with small eigenvalues mostly represent noise and are sensitive to the observations that are inconsistent with the correlation among the variables [51, 52]. Therefore, the smallest principal components (i.e., associated with small eigenvalues) should be useful in anomaly detection. The smallest ignored PCs can be used as an indicator about the existence or absence of faults. When the monitored process is under healthy conditions (no anomaly), the least important principal components are close to zero. However, when a anomaly occurs, then they tend to largely deviate from zero indicating the presence of a new condition that is significantly distinguishable from the normal healthy mode. In this paper, MEWMA is used to enhance process monitoring through its integration with PCA. Because of the ability of the MEWMA control scheme to detect small/moderate changes in the data, this technique is appropriate to improve the detection of moderate anomalies. Thus, this work exploits the advantages of the MEWMA control scheme to improve anomaly detection over the conventional PCA-based methods. Towards this end, the MEWMA control scheme is used to monitor the ignored principal components, which correspond to the small eigenvalues of the PCA model.

6.1. PCA-based MEWMA process monitoring algorithm:

In this approach, the MEWMA monitoring scheme is applied using the principal components ignored (which have smallest variances) from the PCA model. If the matrix of ignored principal components is defined as $\tilde{\mathbf{T}} =$

$[\mathbf{t}^{l+1}, \dots, \mathbf{t}^j, \dots, \mathbf{t}^m]$, where $\mathbf{t}^j \in R^n$, i.e., $\mathbf{t}^j = [\mathbf{t}_1^j, \dots, \mathbf{t}_l^j, \dots, \mathbf{t}_n^j]$, then the MEWMA function can be computed using the residuals of the j^{th} principal component as follows:

$$z_t^j = r\mathbf{t}_t^j + (1-r)z_{t-1}^j, \quad j \in [1, m-l]. \quad (12)$$

The MEWMA decision function, V_t^2 , can be calculated recursively as follows [43]:

$$V_t^2 = \mathbf{Z}_t^T \Sigma_{\mathbf{Z}_t}^{-1} \mathbf{Z}_t \quad (13)$$

where $\Sigma_{\mathbf{Z}_t}$ is the variance-covariance matrix of \mathbf{Z}_t .

In this case, since the MEWMA control scheme is applied on the ignored $m-l$ principal components, one MEWMA decision function will be computed to monitor the process. However, this approach can only detect the presence of anomalies, i.e., it can not determine their locations. This approach is summarized in Table 1.

Table 1: PCA-based MEWMA fault detection algorithm.

Step	Action
1.	<p>Given:</p> <ul style="list-style-type: none"> • A training fault-free data set that represents the normal process operations and a testing data set (possibly faulty data), • The parameters of the MEWMA control scheme: smoothing parameter r and the probability of false alarm α,
2.	<p>Data preprocessing</p> <ul style="list-style-type: none"> • Scale the data to zero mean and unit variance,
3.	<p>Build the PCA model using the training fault-free data</p> <ul style="list-style-type: none"> • Express the data matrix as a sum of approximate and residual matrices as shown in equation (3), • Compute the ignored principal components $\tilde{\mathbf{t}}^j$, using PCA, • Compute the MEWMA control limits,
4.	<p>Test the new data</p> <ul style="list-style-type: none"> • Scale the new data, • Compute the principal components $\tilde{\mathbf{t}}^j$, using PCA, • Compute the MEWMA decision function, V_t^2,
5.	<p>Check for anomalies</p> <ul style="list-style-type: none"> • Declare a fault when the MEWMA decision function, V_t^2, exceeds the control limits.

In the next section, the performance of the proposed PCA-based MEWMA fault detection method will be evaluated and compared to that of the conventional PCA anomaly detection scheme through their application to monitor wind rotor induction machines.

7. Results and discussion

In this section, the proposed PCA-based MEWMA anomaly detection scheme is applied in order to detect abnormalities in ozone measurements caused by air pollution or any incoherence between the different network sensors or sensor faults in the framework of regional ozone surveillance network in Upper Normandy. The performance of the proposed method is compared to that obtained with the conventional PCA approach and to that declared by Air Normand air monitoring association.

7.1. Problem setting

In this study, the data that we use was extracted from the Upper Normandy region. The ozone concentrations data are measured each fifteen minutes in order to limit spatial and temporal sampling problems. The data series of ozone concentrations measured from 11 August to 19 August, 2006 with a total number of 773 observations were used to develop a PCA model without faults. Plots of the original ozone concentration times series and of the corresponding auto-correlation functions (ACF) are shown in Figures 6. Only the curves of the three stations 'SRC', 'QUI' and 'ND2' are plotted for better readability of the figures. These three stations behave like the others network stations.

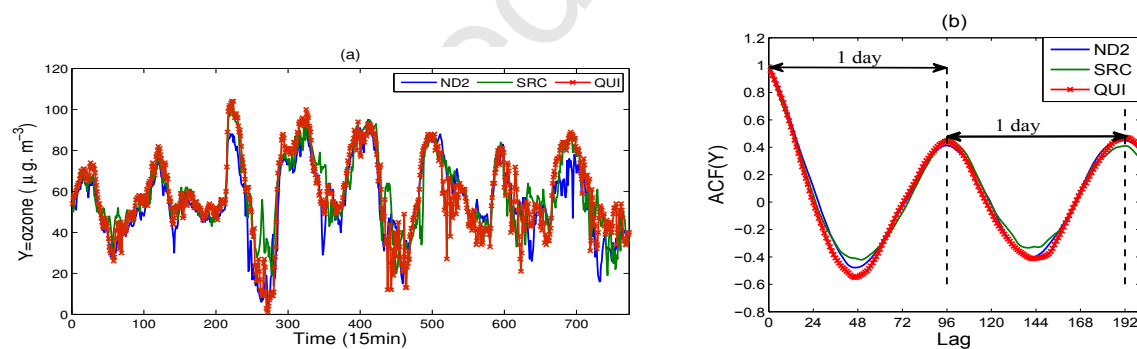


Figure 6: (a) Quarter-hourly ozone time series and (b) ACF of ozone time series.

From Figures 6, the ACF graphics shows an apparent periodicity of 24 hours. It is well known that the distance between extremum points in the autocorrelation functions gives the period of the time series. We suspect that this periodicity is related to the diurnal cycle of ozone which is primarily caused by the diurnal temperature cycle. This periodic variation is due to the cycle of solar radiation (day/night) which is closely related to the mechanism of formation of this pollutant. We also can see the similarity between the autocorrelation functions of ozone concentrations

of the majority of network stations. Monitoring such data therefore requires an initial processing step where such explainable patterns and seasonality are removed. PCA can handle the high dimension of the measurement network and the high degree of correlation among some variables. The purpose is to detect abnormalities in ozone measurements.

7.2. PCA modelling

Firstly, a PCA model is build using training data set. The fault-free data used to develop the model was arranged in a matrix \mathbf{X} with 773 rows (samples) and 7 columns (ozone concentration variables). These data matrix are scaled (to be zero mean with a unit variance), and then used to construct a PCA model.

The scaled fault-free data matrix is used to construct a PCA model, and the computed principal components are shown in Figure 7. Indeed, the principal components (PCs) are linear combinations of the original ones and are uncorrelated. Although PCs represent directions (or patterns) that explain most of the observed variability, their interpretation is, however, not always simple. More specifically, they are just mathematical constructs chosen to represent the variance as efficiently as possible and to be orthogonal to each other. It can be noticed from figure 7 that the principal components t_3, \dots, t_7 represent mainly noise while the first two principal components t_1 and t_2 capture most of the important variations in the data. More specifically, the first principal component, t_1 , is the direction of greatest variability in the data (capture 86:88% of the total variations in the data). The second, t_2 , is the next orthogonal (uncorrelated) direction of greatest variability (capture 4:34% of the total variations in the data). In this case study, t_1 and t_2 capture most of the important variations in the data.

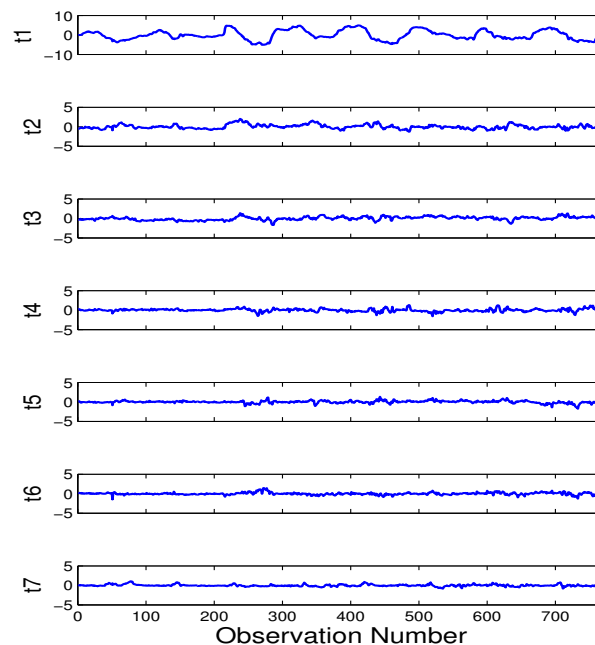


Figure 7: The principal components of the fault-free data.

In PCA, most of the important variations in the data are usually captured in few principal components corresponding to the largest eigenvalues. In this work, the cumulative percent variance (CPV) method is used to determine the optimum number of retained principal components. Using a CPV threshold value of 90%, only the first two principal components will be retained since they capture 86.88% and 4.34% of the total variations in the data.

Indeed, the principal components are linear combination of the original ones, and are uncorrelated with one another. To determine whether principal components are uncorrelated, the scatter plot of PC1 and PC2 is examined. If there were a noticeable relationship in this plot, it would be attributed to non-linear relationships in the data. The PC technique removes all linear correlations and results in a scatter plot when the non-linear relationships are small or nonexistent. Figure 8 shows the bivariate scores plot of PC1 versus PC2 and shows that PC1 and PC2 are uncorrelated. The PCA technique removes all linear correlations.

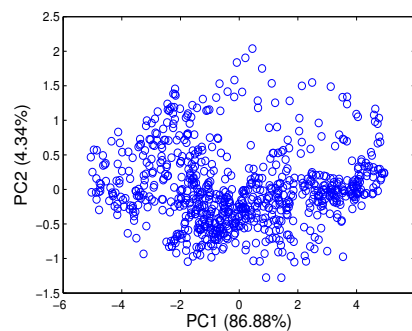


Figure 8: PC1 versus PC2.

Figure 9 present standardized measurements and estimation for of the whole measurements network , the estimations being given by the PCA model. By taking into account the nature of considered process, the results are very satisfactory. With this PCA model based on the first two PCs, the ozone concentrations is generally correctly estimated. However, for some variables we can have modelling errors as shown in Figure 9 (stations ND2, TAN and QUI). In conclusion, the linear PCA was able to model the relations between the various variables. However as we could not it, certain variables being less better estimated than others, we now will examine the effect of the medelling errors on the fault detection phase.

7.3. Detection results

In this section, the anomaly detection abilities of the developed PCA-based MEWMA anomaly detection approach will be assessed using the Upper Normandy ozone data which are completely independent from the training data used to construct the reference PCA model. To evaluate the performance of the developed method, the detection results of the proposed method are compared to that declared by Air Normand, and to that of conventional PCA. Three different testing data sets have been used to evaluate the performance of the PCA-based MEWMA anomaly detection scheme. The first sample covers the period from 11 June 2006 to 09 July 2006, a period of 27 days. The second sample

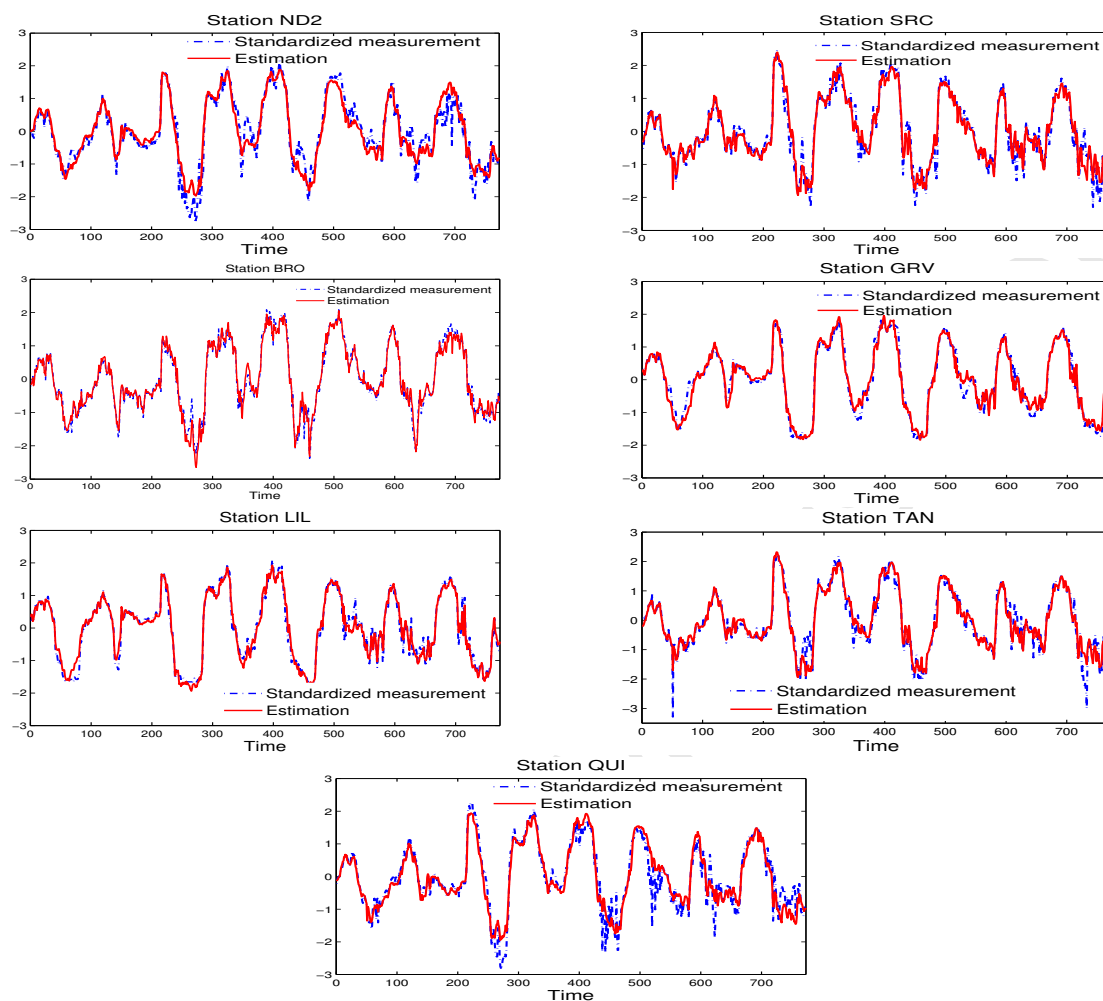


Figure 9: Measurements and estimation of ozone level for the three station.

covers the period from 19 August 2006 to 8 September 2006, a period of 21 days. The latter covers the period from 9 September 2006 to 10 October is a period of 29 days. When the developed PCA-based MEWMA anomaly detection scheme is applied using the fault-free data, the MEWMA threshold value is found to be $h(\alpha) = 9.65$ for a smoothing parameter $r = 0.25$ and a false alarm probability of $\alpha = 0.005$. The detection results are given in Table 2 and are visually illustrated in Figure 10.

In Table 2, the first seven columns present the results of analysis given by Air Normand experts. The first column presents the date of an anomaly observed by experts of Air Normand. The second and third column present the time and the maximum peak intensity. The column 4 presents the station name where the anomaly has occurred and columns 5, 6 and 7 show the beginning, the end and the duration of this anomaly. The column 8 shows the results of detection given by PCA-based MEWMA anomaly detection scheme. The columns 9 and 10 show the results of detection given by the conventional PCA detection indices, T^2 , and Q respectively. If the result is yes, then it is a

correct detection. If the result is *no*, then it is a missed detection. For example take the first two lines to describe how to read this table. The first line indicates that the station 'LIL' has measured abnormal level ozone 12/06/2006 between 11:30 and 12:45 for a total duration of 0:45 minutes and the anomaly peak has occurred at 11:45 with a maximum intensity level in $141.3\mu\text{g}$. The developed PCA-based MEWMA anomaly detection scheme does not detect this anomaly (see column 8 detection). The results of the T^2 and Q statistics shown in columns 9 and 10, respectively, show that the conventional PLS was unable to detect this anomaly. In the second line, ND2 and LIL stations have presented abnormalities on 13/06/2006. The PCA-based MEWMA scheme has correctly detected these anomalies. The results using the Q statistic given in column 10 show that it could successfully detect this anomaly. However, Hotelling's T^2 statistic was unable to detect this anomaly. This result may be explained by the fact that the T^2 statistic provides a measure of the deviation in the PCs that are of greatest importance to the normal process condition. Thus, the normal operating region defined by the T^2 control limits is usually larger than that defined by the Q control limits. Therefore, anomalies with moderate magnitudes can easily exceed the Q threshold, but not the T^2 threshold, which makes the Q statistic usually more sensitive than T^2 for this anomaly. By comparing the results obtained by the PCA-based MEWMA detector and results declared by Air Normand, we note that the PCA-based MEWMA detector has detected almost the totality of anomalies (see Table 2 and Figure 10). For our application, the proposed fault anomaly algorithm improves the anomaly detection compared to classical detection indices Q and T^2 . The developed PCA-based MEWMA anomaly detection algorithm takes very little time to give its verdict. Hence, the proposed algorithm can be used as an automatic tool of abnormal ozone peaks (or sensors faults) detection in the framework of regional air quality monitoring networks.

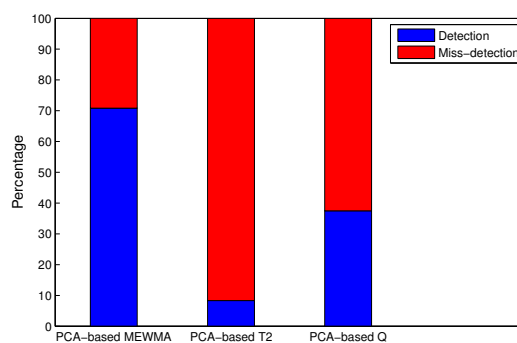


Figure 10: Compare detection results.

8. Conclusion

In this paper, an anomaly detection scheme based on principal component analysis is proposed to monitor the ozone concentrations in the Upper Normandy region, France. To enhance anomaly detection a new PCA-based monitoring strategy combining PCA with the multivariate exponentially weighted moving average (MEWMA) monitoring

Table 2: detection results.

Date	Air Normand detection						PCA-MEWMA	PCA- T^2	PCA- Q
	Hour	Intensity	Places	Beginning	End	Duration			
12/06/2006	11:45	141	LIL	11 :30	12 :15	0 :45	no	no	no
13/06/2006	13 :15	168	LIL	12 :30	13 :45	1 :15	yes	no	yes
		181	ND2	12 :15	14 :15	2 :00	yes	no	no
17/06/2006	08 :00	132	SRC	7 :15	10:15	3 :00	no	no	no
	08 :30	141	TAN	8 :00	9 :00	1 :00	yes	no	no
23 /06/2006	14 :15	137	LIL	13 :00	15 :00	2 :00	no	no	no
	14 :30	126	ND2	13 :00	15 :15	2 :15	no	no	no
	14 :45	127	QUI	13 :15	15 :15	2 :00	no	no	no
30/06/2006	08 :00	144	TAN	7 :15	8 :15	1 :00	yes	no	no
03/07/2006	08 :15	244	TAN	8:15	9 :15	1:00	yes	yes	yes
	10 :15	242	TAN	9 :15	11 :15	2 :00	yes	yes	yes
	9 :30	179	LIL	9 :00	10 :15	1 :15	yes	no	yes
	10 :00	166	QUI	9 :15	10 :15	1 :00	yes	no	no
04/07/2006	07 :45	201	ND2	6 :30	10 :00	3 :00	yes	no	yes
05/09/2006	09 :45	180	LIL	7 :45	10:45	3:00	yes	no	no
	09 :45	115	TAN	8 :15	11 :00	2 :45	no	no	no
06/09/2006	09 :45	182	LIL	8 :15	10 :30	2 :15	yes	no	yes
	11 :15	168	LIL	10 :30	13 :15	2 :45	yes	no	no
	14 :00	168	ND2	13 :15	15 :00	1 :45	yes	no	no
	14 :30	168	GRV	13 :45	15 :00	1 :15	yes	no	no
10/09/2006	09 :30	167	QUI	7 :30	10 :00	2 :30	yes	no	no
	09 :45	146	LIL	8 :45	10 :30	1 :45	yes	no	no
	11 :00	180	TAN	10 :15	11 :30	1 :15	yes	no	no
	12 :00	166	GRV	11 :30	12 :45	1 :15	no	no	no

scheme is proposed. In the proposed approach, MEWMA control scheme is applied on the ignored principal components (which have smallest variances) to detect the presence of anomalies. The proposed PCA-based MEWMA anomaly detection scheme is successfully applied to data of the ozone concentrations collected from the Upper Normandy region, France. For this application, the PCA-based MEWMA scheme improves the anomaly detection compared to that of the conventional PCA-based monitoring charts. The results indicate that the PCA-based MEWMA test can be used as an automatic tool to detect abnormal ozone measurements.

References

- [1] H. Moshhammer, "Communicating health impact of air pollution," 2010.
- [2] A. Nawahda, "An assessment of adding value of traffic information and other attributes as part of its classifiers in a data mining tool set for predicting surface ozone levels," *Process Safety and Environmental Protection*, vol. 99, pp. 149–158, 2016.
- [3] S. Sillman, "Tropospheric ozone and photochemical smog," in B. Sherwood Lollar, ed., *Treatise on Geochemistry, Environmental Geochemistry, Ch. 11, Elsevier*, vol. 9, pp. 407–431, 2003.
- [4] M. Chiogna and F. Pauli, "Modelling short-term effects of ozone on morbidity: an application to the city of milano, italy, 1995–2003," *Environmental and Ecological Statistics*, vol. 18, no. 1, pp. 169–184, 2011.
- [5] J. Seinfeld and S. Pandis, "Atmospheric chemistry and physics: from air pollution to climate change." *John Wiley & Sons, Inc , New York*, 2006.
- [6] S. Yin, G. Wang, and H. Karimi, "Data-driven design of robust fault detection system for wind turbines," *Mechatronics*, vol. 24, no. 4, pp. 298–306, 2014.
- [7] S. Qin, "Statistical process monitoring: Basics and beyond," *Journal of Chemometrics*, vol. 17, no. 8/9, pp. 480–502, 2003.
- [8] A. Herve and J. Lynne, "Principal component analysis," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, pp. 433–459, 2010.
- [9] S. Yin, S. Ding, X. Xie, and H. Luo, "A review on basic data-driven approaches for industrial process monitoring," *IEEE Transactions on Industrial Electronics*, vol. 61, no. 11, pp. 6418–6428, 2014.
- [10] S. Qin, "Survey on data-driven industrial process monitoring and diagnosis," *Annual Reviews in Control*, vol. 36, no. 2, pp. 220–234, 2012.
- [11] F. Khan, S. Rathnayaka, and S. Ahmed, "Methods and models in process safety and risk management: Past, present and future," *Process Safety and Environmental Protection*, vol. 98, pp. 116–147, 2015.
- [12] A. Banimostafa, S. Papadokonstantakis, and K. Hungerbühler, "Evaluation of EHS hazard and sustainability metrics during early process design stages using principal component analysis," *Process safety and environmental protection*, vol. 90, no. 1, pp. 8–26, 2012.
- [13] J. George, Z. Chen, and P. Shaw, "Fault detection of drinking water treatment process using PCA and hotellingAes T^2 chart," *World Academy of Science, Engineering and Technology*, vol. 50, pp. 970–975, 2009.
- [14] F. Harrou, F. Kadri, S. Chaabane, C. Tahon, and Y. Sun, "Improved principal component analysis for anomaly detection: Application to an emergency department," *Computers & Industrial Engineering*, vol. 88, pp. 63–77, 2015.
- [15] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol. 24, pp. 417–441, 1933.
- [16] G. E. P. Box, "Some theorems on quadratic forms applied in the study of analysis of variance problems: Effect of inequality of variance in one-way classification," *The Annals of Mathematical Statistics*, vol. 25, pp. 290–302, 1954.
- [17] J. Romagnoli and A. Palazoglu, "Introduction to process control," *CRC Press, United States of America*, 2006.
- [18] P. Geladi and B. Kowalski, "Partial least-squares regression: a tutorial," *Analytica chimica acta*, vol. 185, pp. 1–17, 1986.
- [19] T. Kourti and J. MacGregor, "Process analysis, monitoring and diagnosis using multivariate projection methods: A tutorial," *Chemometrics and Intelligent Laboratory Systems*, vol. 28, no. 3, pp. 3–21, 1995.
- [20] J. MacGregor and T. Kourti, "Statistical process control of multivariate processes," *Control Engineering Practice*, vol. 3, no. 3, 1995.
- [21] Q. Chen, U. Kruger, M. Meronk, and A. Leung, "Synthesis of T^2 and Q statistics for process monitoring," *Control engineering practice*, vol. 12, no. 6, pp. 745–755, 2004.
- [22] G. Ranger and F. Alt, "Choosing principal components for multivariate statistical process control," *Communications in Statistics-Theory and Methods*, vol. 25, no. 5, pp. 909–922, 1996.
- [23] T. Kourti and J. MacGregor, "Multivariate spc methods for process and product monitoring," *Journal of Quality Technology*, vol. 28, no. 4, 1996.
- [24] C. Mastrangelo, G. Runger, and D. Montgomery, "Statistical process monitoring with principal components," *Quality and Reliability Engineering International*, vol. 12, no. 3, pp. 203–210, 1996.

- [25] M. Harkat, G. Mourot, and J. Ragot, "An improved PCA scheme for sensor FDI: application to an air quality monitoring network," *Journal of Process Control*, vol. 16, no. 6, pp. 625–634, 2006.
- [26] D. C. Montgomery, "Introduction to statistical quality control," *John Wiley & Sons, New York*, 2005.
- [27] C. Vlachokostas, S. Nastis, C. Achillas, K. Kalogeropoulos, I. Karmiris, N. Moussiopoulos, E. Chourdakis, G. Baniyas, and N. Limperi, "Economic damages of ozone air pollution to crops using combined air quality and GIS modelling," *Atmospheric Environment*, vol. 44, pp. 3352–3361, 2010.
- [28] A. Detournay, S. L. Meur, and V. Delmas, "Understanding of the atypical ozone peaks phenomenon observed around the petrochemical industrial zone of port-jerome in upper normandy, france," *Pollution atmosphérique*, vol. 196, pp. 405–422, 2007.
- [29] G. Brulfert, O. Galvez, F. Yang, and J. Sloan, "A regional modelling study of the high ozone episode of June 2001 in southern Ontario," *Atmospheric Environment*, vol. 41, pp. 3777–3788, 2007.
- [30] C. Dueñas, M. Fernández, S. Cañete, J. Carretero, and E. Liger, "Analyses of ozone in urban and rural sites in Málaga (Spain)," *Chemosphere*, vol. 56, pp. 631–639, 2004.
- [31] A. Proyou, G. Toupance, and P. Perros, "A two year study of ozone behaviour at rural and forested sites in eastern France," *Atmospheric Environment*, vol. 25A, no. 10, pp. 2145–2153, 1991.
- [32] E. Brankov, R. Henry, K. Civerolo, W. Hao, S. Rao, P. Misra, R. Bloxam, and N. Reid, "Assessing the effects of transboundary ozone pollution between Ontario, Canada and New York, USA," *Environmental Pollution (Oxford, United Kingdom)*, vol. 123, no. 3, pp. 403–411, 2003.
- [33] W. Chen, H. Tang, and H. Zhao, "Diurnal, weekly and monthly spatial variations of air pollutants and air quality of Beijing," *Atmospheric Environment*, vol. 119, pp. 21–34, 2015.
- [34] I. Zdanevitch, "Etude d'épisodes inexplicables d'ozone," *Rapport LCSQA, conversion 41/2000. INERIS, Paris.*, 2001.
- [35] R. J. Patton and J. Chen, "A review of parity space approaches to fault diagnosis," in *Proceedings of SAFEPROCESS'91*, 1991, pp. 239–255.
- [36] P. Ralston, G. DePuy, and J. Graham, "Computer-based monitoring and fault diagnosis: a chemical process case study," *ISA Transactions*, vol. 40, no. 1, 2001.
- [37] J. Jackson and G. Mudholkar, "Control procedures for residuals associated with principal component analysis," *Technometrics*, vol. 21, pp. 341–349, 1979.
- [38] S. Qin and R. Dunia, "Determining the number of principal components for best reconstruction," *Journal of Process Control*, vol. 10, no. 2, pp. 245–250, 2000.
- [39] B. Li, J. Morris, and E. Martin, "Model selection for partial least squares regression," *Chemometrics and Intelligent Laboratory Systems*, vol. 64, no. 1, pp. 79–89, 2002.
- [40] M. Zhu and A. Ghodsi, "Automatic dimensionality selection from the scree plot via the use of profile likelihood," *Computational Statistics & Data Analysis*, vol. 51, pp. 918–930, 2006.
- [41] D. Bissell, *Statistical Methods for SPC and TQM*. CRC Press, 1994, vol. 26.
- [42] E. S. Page, "Continuous inspection schemes," *Biometrika*, pp. 100–115, 1954.
- [43] C. A. Lowry, W. H. Woodall, C. W. Champ, and S. E. Rigdon, "A multivariate exponentially weighted moving average control chart," *Technometrics*, vol. 34, no. 1, pp. 46–53, 1992.
- [44] H. Hotelling, "Multivariate quality control illustrated by the air testing of sample bomb sights, techniques of statistical analysis, ch. ii," 1947.
- [45] S. W. Roberts, "Control chart tests based on geometric moving averages," *Technometrics*, vol. 1, no. 3, pp. 239–250, 1959.
- [46] J. Lucas and M. Saccucci, "Exponentially weighted moving average control schemes: properties and enhancements," *Technometrics*, vol. 32, no. 1, pp. 1–12, 1990.
- [47] J. S. Hunter, "The exponentially weighted moving average," *Journal of Quality Technology*, vol. 18, no. 4, pp. 203–210, 1986.
- [48] G. Runger and S. Prabhu, "A Markov chain model for the multivariate exponentially weighted moving averages control chart," *Journal of the American Statistical Association*, vol. 91, no. 436, pp. 1701–1706, 1996.
- [49] S. Rigdon, "An integral equation for the in-control average run length of a multivariate exponentially weighted moving average control chart," *Journal of Statistical Computation and Simulation*, vol. 52, no. 4, pp. 351–365, 1995.

- [50] K. M. Bodden and S. E. Rigdon, "A program for approximating the in-control ARL for the MEWMA chart," *Journal of Quality Technology*, vol. 31, no. 1, pp. 120–123, 1999.
- [51] J. Jobson, *Applied multivariate data analysis*. Springer Heidelberg, 1992, vol. 2.
- [52] D. Donnell, A. Buja, and W. Stuetzle, "Analysis of additive dependencies and concavities using smallest additive principal components," *The Annals of Statistics*, pp. 1635–1668, 1994.

Accepted Manuscript