# Statistical Measurement Validation with Application to Electronic Nose Technology

Mina Mirshahi, Vahid Partovi Nia and Luc Adjengue

*Department of Mathematics and Industrial Engineering, Polytechnique Montreal, Montreal, Quebec, Canada*

Keywords:     Artificial Olfaction, Electronic Nose, Gas Sensor, Odor, Outlier, Robust Covariance Estimation.

Abstract:     An artificial olfaction called electronic nose (e-nose) relies on an array of gas sensors with the capability of mimicking the human sense of smell. Applying an appropriate pattern recognition on the sensor's output returns odor concentration and odor classification. Odor concentration plays a key role in analyzing odors. Assuring the validity of measurements in each stage of sampling is a critical issue in sampling odors. An accurate prediction for odor concentration demands for careful monitoring of the gas sensor array measurements through time. The existing e-noses capture all odor changes in its environment with possibly varying range of error. Consequently, some measurements may distort the pattern recognition results. We explore e-nose data and provide a statistical algorithm to assess the data validity. Our online algorithm is computationally efficient and treats data as being sampled.

## 1 INTRODUCTION

The ability to recognize the chemicals in the environment is a very basic and essential need for the living organisms; from a single-cell amoebae to human beings, all species are provided with a chemical awareness system. Human beings have three sensory systems to detect odors: sense of taste, sense of smell, and chemical feel with receptors all over the body. All species employ their chemical senses to approach and being attracted to possibly safe conditions, as well as avoiding and being resisted to the harmful ones. As for human beings, in every breath, the sense of smell collects a sample from its environment and forwards it to the brain for further analyses. Unlike the sense of taste, smell can be captured from a distance and assist the brain in producing a warning. Unfortunately, the human sense of smell does not respond to all harmful air pollutants. Additionally, sensitivity of humans to many air pollutants varies — one can be accustomed to a toxic smell. In the last decade, great attention has been paid to the subject of air quality because it directly influences the environmental and human health. A crucial element in assessment of indoor and outdoor air quality is auditing the odorants. There exists various odor measurement techniques such as dilution-to-threshold, olfactometers, and referencing techniques (McGinley and Inc, 2002). The performance of these approaches depend on human evaluation. Due to the high variability of individual's sensitivity, the common methods mostly lack accuracy. In 1982, the first gas multisensor array was invented as primary artificial olfaction (Persaud and Dodd, 1982). The term electronic nose (e-nose) was introduced by Gardner and Bartlett (1994). E-nose is an artificial olfactory system which consists of an array of gas sensors. The e-nose is designed for recognizing complex odors in its surrounding environment. The gas sensor array receives chemical information about gaseous mixtures as input and converts it to measurable signals. Sensors act independently and simultaneously in this device. Cross-sensitivity of gas sensors is inevitable in sensor array structure. The cross-sensitivity is the interaction among chemicals that leads to a different signal from the component in a mixture compared to the single component. Gas sensor's performance is affected by different elements which make it unstable and less sensitive to odors. One of the most serious deterioration in sensors is owing to a phenomenon called *drift*. Drift is a temporal change in sensor's response while all other external conditions are kept constant. The majority of manufactured sensor arrays are subject to drift, and several methods have been introduced to overcome this problem (Carlo and Falasconi, 2012; Artursson et al., 2000; Padilla et al., 2010; Zuppa et al., 2007). The behavior of a sensor is directly influenced by the surrounding chemical and

physical conditions. For instance, the sensor response may depend on the temperature of the gas under examination. Therefore, thermal conditions around the sensing elements need to be supervised. The multivariate response of gas sensor arrays undergoes different pre-processing procedures before the prediction is performed using statistical tools such as regression, classification, or clustering. Numerous methods have been developed for analyzing the gas sensor array data, including Gutierrez-Osuna (2002); Kermiti and Tomic (2003); Bermak et al. (2006).

## 2 PROBLEM STATEMENT

The e-nose has partially addressed the human sense of smell in diverse industrial sites. Unwanted variability may occur in sensor's output data. This happens due to environmental factors or physical impairment of the system, since e-noses are installed in outdoor fields where the conditions can dramatically fluctuate. This demands for monitoring the critical factors through adding extra sensors and temperature compensation in sensor pre-processing. The sensor's output is used to quantify odor concentration. Transferring the data to olfactometry is both time consuming and costly. Only small portions of data are appointed for further analyses of its concentration in olfactometry. Pattern recognition methods are employed in order to predict the odor concentration for each set of sensor values. To assess the accuracy of predictions, the validity of sensor values must be ensured. Sensors in the e-nose structure may report incorrect values or some stop functioning for a short period of time. These anomalies are ought to be diagnosed and reported in real time using a computationally efficient algorithm.

## 3 DATA DESCRIPTION

The data under the study include 11 distinct attributes, each representing sensor values of the e-nose. Sensors react to almost all gases in the air, but they are designed so that each sensor is more sensitive to a specific type of gas. Some of the sensors are highly positively correlated with each other, see Figure 1 and Figure 2 (left panel).

Suppose that $\mathbf{x}_{p \times 1}^{\top}$ is a random vector of $p = 11$ attributes, in which $\mathbf{a}^{\top}$ illustrates the transpose of vector $\mathbf{a}$, and its $n$ independent realization are stored in the rows of data matrix $\mathbf{X}_{n \times p}$. The covariance matrix of
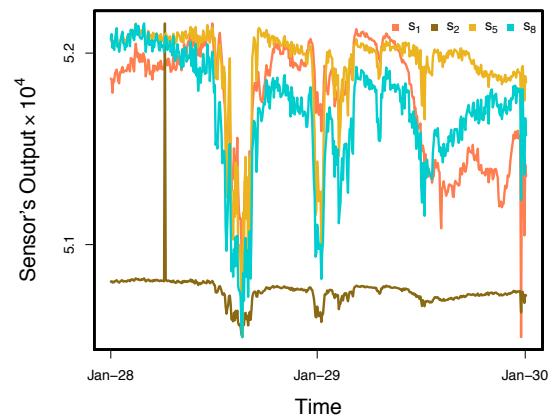


Figure 1: Senor's output during three days of sampling for 4 randomly selected sensors.

$\mathbf{x}_{p \times 1}$ , say $\Sigma = [\sigma_{ij}]_{i,j=1,2,...,p}$, is defined as

$$\Sigma_{p \times p} = \mathrm{Cov}(\mathbf{x}) = \mathrm{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\top}\},$$

where $\boldsymbol{\mu}$ represents the mean of $\mathbf{x}$, E is the expectation operator. The covariance, $\sigma_{ij}$, measures the degree to which two attributes are linearly associated. It is well-known that the inverse of covariance matrix, commonly known as precision matrix, yields the partial correlation between the attributes. The partial correlation is the correlation between two attributes conditioning on the effect of other attributes. Non-zero elements of $\Sigma^{-1}$ implies the conditional dependence. Therefore, the sparse estimation of $\Sigma^{-1}$ pinpoints the block dependent structure of attributes. The sparse estimation of $\Sigma^{-1}$ set some of the $\Sigma^{-1}$ entries exactly to zero. Investigation of the inherent dependence between the sensor values is then performed by means of the partial correlation. In order to obtain a clear image of sensors which are potentially grouped together, the *graphical lasso* (Friedman et al., 2008) is used. Friedman et al. (2008) considered estimating the inverse of covariance matrix, $\Sigma^{-1}$, sparsely by applying a *lasso penalty* (Tibshirani, 1996). In Figure 2 (right panel), the undirected graph connects two variables which are conditionally correlated given all other attributes. For instance, the sensors 9, 10, and 11 are conditionally correlated with each other. This also agrees with the heatmap of the correlation matrix Figure 2 (left panel). Thus, this dependence must be taken into account while modeling data. Another vital assumption that should be verified is the Gaussianity of the data. The non-Gaussianity of the sensor values is established using various methods such as analyzing the distribution of individual sensor values, scatter plot of the linear projection of data using principal components, estimating the multivariate kurtosis and skewness, and also multivariate Mardia test, see Figure 3.
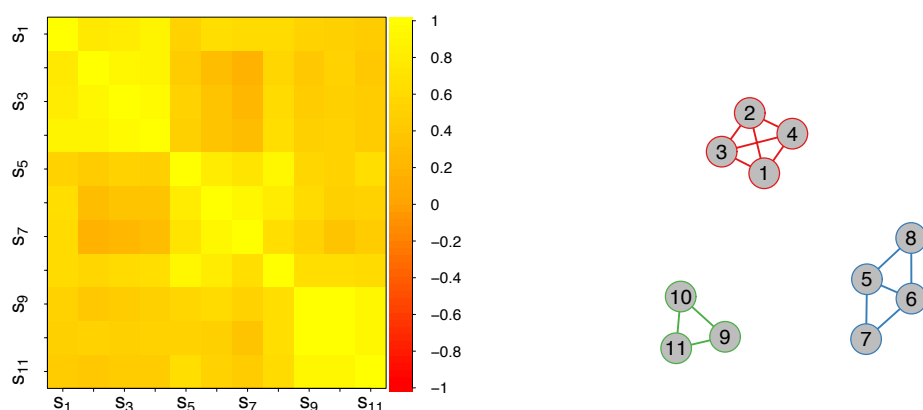
Figure 2: Left panel, heatmap of the correlation matrix of the sensor values ($s_1$–$s_{11}$). Right panel, the undirected graph of partial correlation using the graphical lasso. The undirected graph of the right panel approves the block structure of the heatmap of the left panel.
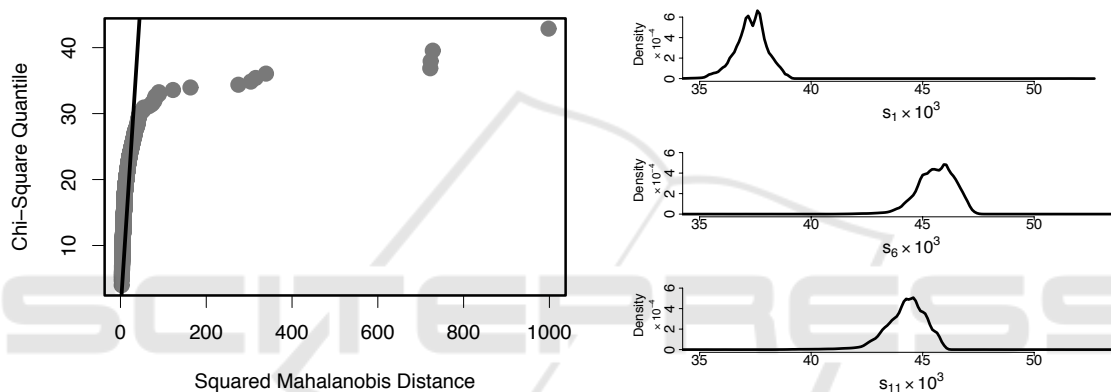


Figure 3: Left panel, the Q-Q plot of squared Mahalanobis distance supposed to follow chi-square distribution for Gaussian data. Right panel, the marginal density for some randomly chosen sensor values. Both graphs confirm the non-Gaussianity of data.

# 4 METHODOLOGY

In order to demonstrate the validity of the e-nose measurements, we aim to allocate each sample to different zones. To be able to verify the validity of the measurements, it is necessary to have some reference samples for the purpose of comparison. These reference samples are collected while the e-nose is at its best performance, and the conditions are fully under control. For the data set under the study, there are two distinct reference sets. *Reference* 1 is constituted of data in a period of sampling defined by an expert after installation of the e-nose. We call the data in this period of sampling as *proposed set*. *Reference* 2, upon its availability, is manually gathered samples from the field and brought to the laboratory to quantify the odor concentration. We call the latter data, *calibration set* to emphasize that it can be used for data modeling using supervised learning. If new data diverge greatly

from the overall pattern of data previously seen, then it is marked as an outlier and is allocated to the red zone. This zone represents a dramatic change in the pattern of samples and refer to "risky" samples. If new data is non-outlier and it is also located within the data polytope of the Reference 1 or the Reference 2, it is assigned to green or blue zone respectively. These zones represent the "safe" samples. If new data is non-outlier, but outside of the area of green and blue zones, it is assigned to yellow zone. This zone displays potentially "critical" samples.

Producing many samples belonging to the yellow and the red zones is an indication of a major flaw in the system. Physical complications, such as sensor loss in the e-nose, or sudden changes in the chemical pattern of the environment, account for all undesirable measurements. Zone assignment, therefore, require some outlier detection algorithms. To define the green and the blue zones, the new samples

Table 1: Description of each zones in validity assessment procedure.

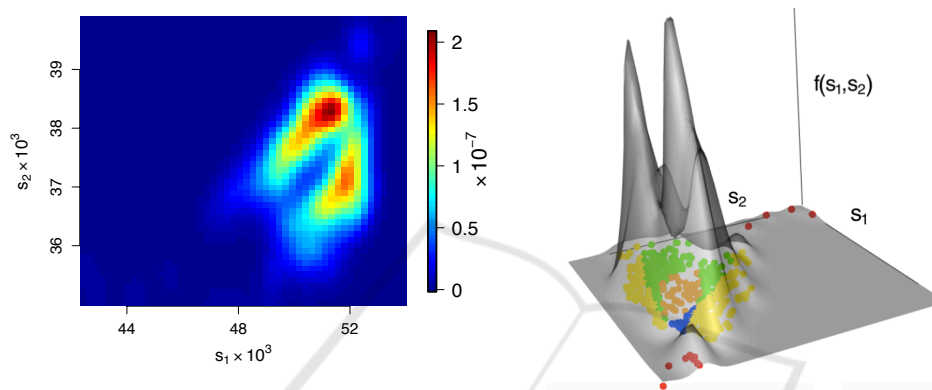| Zone | Description |
|------|-------------|
| Red | Observations that are outliers in terms of AO measure. |
| Green | Observations that are non-outliers in terms of AO measure. Moreover, they fall into the polytope of the Reference 1. |
| Blue | Observations that are non-outliers in terms of AO measure. Moreover, they fall into the polytope of the Reference 2. |
| Orange | Observations that are non-outliers in terms of AO measure. Moreover, they fall into the polytopes of both the Reference 1 and the Reference 2. |
| Yellow | Observations that are non-outliers in terms of AO measure. Moreover, they do not fall into the polytope of neither the Reference 1 nor the Reference 2. |



Figure 4: Validity assessment for about 700 samples based on 2 sensor values. Left panel, the plot illustrates the contour map of estimated density function for the 2 sensors. Right panel, the density function of the samples demonstrated in $3D$ with zones identified for each of the samples in the sensor 1 ($s_1$) versus sensor 2 ($s_2$) plane. Higher density is assigned to green, blue, and orange zones compared to yellow and red zones.

are projected onto a lower dimension subspace. Dimension reduction methods such as principal component analysis (PCA) can serve this purpose (Jolliffe, 2002). PCA transforms a collection of possibly correlated attributes into a set of linearly uncorrelated axes through orthogonal linear transformations. The first $k$ ($k < p$) principal components are the eigenvectors of the covariance matrix $\Sigma$ associated with the $k$ largest eigenvalues. PCA exploits empirical covariance matrix, $\hat{\Sigma}$, which is extremely sensitive to outliers (Prendergast, 2008). Since the data contain many outliers, robust covariance estimation must be applied to avoid misleading results. Robust principal component analysis (Hubert et al., 2005) is employed for dimension reduction purpose throughout this paper. This robust PCA computes the covariance matrix through projection pursuit (Li and Chen, 1985) and minimum covariance determinant (Croux and Haesbroeck, 2000) methods. The robust PCA procedure can be summarized as follows:

1. The matrix of data is pre-processed such that the data spread in the subspace of at most $\min(n - 1, p)$.

2. In the spanned subspace, the most obvious outliers are diagnosed and removed from data. The covariance matrix is calculated for the remaining data, $\hat{\Sigma}_0$.

3. $\hat{\Sigma}_0$ is used to decide about the number of principal components to be retained in the analysis, say $k_0$ ($k_0 < p$).

4. The data are projected onto the subspace spanned by the first $k_0$ eigenvectors of $\hat{\Sigma}_0$.

5. The covariance matrix of the projected points is estimated robustly using minimum covariance determinant method and its $k$ leading eigenvalues are computed. The corresponding eigenvectors are the robust principal components.

To define the red zone, it is required to find the outliers of data as it is being measured by the e-nose through time. As the data fail to follow a Gaussian distribution, outlier detection methods that rely on the assumption of elliptical contoured distribution should be avoided. Here, outliers are flagged by means of *adjusted outlyingness* (AO) criterion (Brys et al., 2006). If a sample is detected as an outlier by AO measure, it belongs to the red zone. For the specification of the

remaining zones, we need to define the polytopes of the samples in Reference 1 and Reference 2. These polytopes are built using the convex hull of the robust principal component *scores*. More specifically, the boundary of the green zone is defined by computing the convex hull of the robust principal component scores of the Reference 1. A short description of each zone is provided in Table 1. Before determining the color tag for each new data, the samples are checked for missing values and are imputed in case needed by *multivariate imputation* methods such as Josse et al. (2011). The idea behind the validity assessment is visualized in Figure 4. For simplicity, only 2 sensors are used for all computations in Figure 4 and a *2D* presentation of zones is plotted using the sensors' coordinates. Suppose that $\mathbf{X}_{N \times 11}$ represents the matrix of sensor values for $N$ samples, $\mathbf{y}_N$ the vector of corresponding odor concentration values and $\mathbf{x}_l^\top$ is the $l$th row of $\mathbf{X}_{N \times 11}$, $l = 1, 2, \ldots, N$. Furthermore, suppose that $n_1$ refers to the number of samples in the proposed set of the sampling and $n_2$ refers to the number of samples in the calibration set. The samples of the proposed set are always available, but not necessary the calibration set. Two different scenarios occur based on the availability of the calibration set. If the calibration set is accessible, then Scenario 1 happens. Otherwise, we only deal with Scenario 2. Scenario 1 is a general case which is explained more in details. The data undergo a preprocessing stage, including imputation and outlier detection, before any further analyses. Having done the pre-processing stage, data are stored as Reference 1, $\mathbf{X}_{n_1 \times 11}$, and Reference 2, $\mathbf{X}_{n_2 \times 11}$. The first $k$, e.g. $k = 2, 3$, robust principal components of $\mathbf{X}_{n_1 \times 11}$ are calculated and the corresponding *loading* matrix is

---

**Sub-Algorithm:** (Scenario 1).

1: **if** the point $\mathbf{x}_l^\top$, $l = 1, 2, \ldots, N$ is identified as an outlier by *AO* measure **then**
2:     $\mathbf{x}_l^\top$ is in red zone,
3: **else if** $\mathbf{x}_l^\top \mathbf{L}_1 \in \text{ConvexHull}^{(1)}$ AND $\mathbf{x}_l^\top \mathbf{L}_1 \notin \text{ConvexHull}^{(2)}$ **then**
4:     $\mathbf{x}_l^\top$ is in green zone,
5: **else if** $\mathbf{x}_l^\top \mathbf{L}_1 \notin \text{ConvexHull}^{(1)}$ AND $\mathbf{x}_l^\top \mathbf{L}_1 \in \text{ConvexHull}^{(2)}$ **then**
6:     $\mathbf{x}_l^\top$ is in blue zone,
7: **else if** $\mathbf{x}_l^\top \mathbf{L}_1 \in \text{ConvexHull}^{(1)}$ AND $\mathbf{x}_l^\top \mathbf{L}_1 \in \text{ConvexHull}^{(2)}$ **then**
8:     $\mathbf{x}_l^\top$ is in orange zone,
9: **else**
10:     $\mathbf{x}_l^\top$ is in yellow zone.
11: **end if**

---

denoted by $\mathbf{L}_1$. The pseudo code of two algorithms for Scenario 1 is provided below. Scenario 2 is a special case of Scenario 1 in which Sub-Algorithm (Scenario 1) is used with $\text{ConvexHull}^{(2)} = \varnothing$ that eliminates the blue and the orange zones. Consequently, there is no model for odor concentration prediction in the Main Algorithm.

---

**Main Algorithm:** (Scenario 1).

**Require:** $\mathbf{X}_{n_1 \times 11}$, $\mathbf{X}_{n_2 \times 11}$, and the loading matrix $\mathbf{L}_1$ using robust PCA over Reference 1, $\mathbf{X}_{n_1 \times 11}$.
1: $\text{ConvexHull}^{(1)} \leftarrow$ the convex hull of the projected values of the Reference 1, $\mathbf{X}_{n_1 \times 11}\mathbf{L}_1$.
2: Train a supervised learning model on Reference 2, $\mathbf{X}_{n_2 \times 11}$, and its odor concentration vector, $\mathbf{y}_{n_2}$.
3: $\text{ConvexHull}^{(2)} \leftarrow$ the convex hull of the projected values of the Reference 2, $\mathbf{X}_{n_2 \times 11}\mathbf{L}_1$.
4: Do **Sub-Algorithm** for new data $\mathbf{x}^*$.
5: Predict the odor concentration for new data $\mathbf{x}^*$ using the trained supervised learning model.

---

The above steps are implemented over 8 months of data collected by the e-nose in Section 6. In order to justify our choice of statistical techniques, the proposed methodology is run over a set of simulated data in a following section.

## 5 SIMULATION

To emphasize on the importance of the assumptions such as non-eliptical contoured distribution and robust estimation considered in our methodology, we examine the methodology on a set of simulated data. Assume the matrix of data $\mathbf{X}_{N \times 2}$, where $\mathbf{x}_l^\top = (x_{l1}, x_{l2})$; $l = 1, 2, \ldots, N$, are generated according to the mixture of Gaussian and the Student's t-distributions, Figure 5 (top left panel). Ignoring the distribution of data and seeking for any classical approach toward outlier detection, renders some observations as outliers mistakenly, Figure 5 (top right panel). The parameters of interest, the mean vector and the covariance matrix, need to be estimated robustly, otherwise the confidence region misrepresents the underlying distribution. In Figure 5 (bottom left panel), the classical confidence region is pulled toward the outlier observations. On the contrary, the robust confidence region perfectly unveil the distribution of the majority of observations because of the robust and efficient estimation of the mean and the covariance matrix. Consequently, the classical principal components are affected by the inefficient esti-
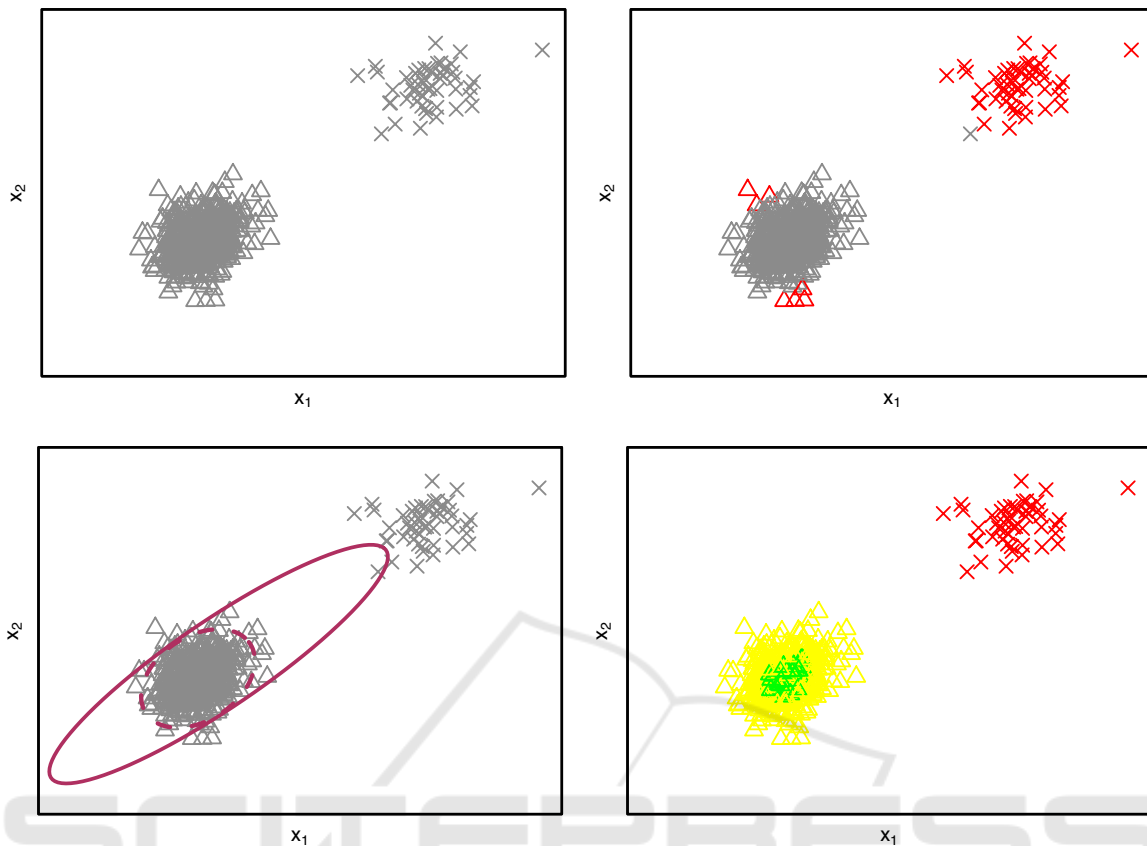
Figure 5: Top left panel, the simulated data from the mixture distribution $f(x) = (1 - \varepsilon)f_1(x) + \varepsilon f_2(x)$ with contamination proportion of $\varepsilon = \frac{1}{10}$, and $f_1$ and $f_2$ being the Gaussian and Student's t-distribution respectively. The data from $f_1$, and $f_2$ are plotted in triangles and crosses correspondingly. Top right panel, the outliers of data are identified and highlighted with red using the classical Mahalonobis distance and 95th percentile of the Chi-square distribution with two degrees of freedom. Bottom left panel, the 95% confidence region for the data is computed using the classical estimates of parameters (solid line) and the robust estimates (dashed line). Bottom right panel, the Main Algorithm is implemented and the zones are graphed by colors described in Table 1.

mation of the covariance matrix. We proposed using methods which deal with contaminated data appropriately. Adjusted outlyingness (AO) measure identifies the outliers of the data correctly. In the Main Algorithm, suppose we take the Gaussian sub-sample as the Reference 1. Figure 5 (bottom right panel) shows the result of our algorithm on the simulated data.

# 6 APPLICATION

For the easy visualization, the first 3 robust principle components of the data are used, $PC1$, $PC2$, $PC3$. These components correspond to the 3 largest eigenvalues of the covariance matrix. In case of sensor failures, the data contain missing values that need to be imputed. First, data are imputed to replace all the missing values, and then the validity of the measure-

ments are identified over the 8 months sampling. Only a subset of 500 samples out of 200 thousands of observations are plotted to make the graphs more readable. In Figure 6, the sample points are drawn in gray and each zone is highlighted using its corresponding color of Table 1. The circles in Figure 6 are also illustrated on $PC1$ and $PC2$ plane for a better demonstration of the zones.

The zones' definition is helpful in interpreting the results. As an example, the green or the blue zone reveals the fact that the sampling points are very close to the samples that have already been observed in either Reference 1 or Reference 2. The observations in reference sets were entirely under control, therefore, the blue and green zones justify the validity of samples. Consequently, the prediction obtained over these samples is expected to be more accurate. On the contrary, the prediction values for the points in the yellow zone are less accurate compared with the green
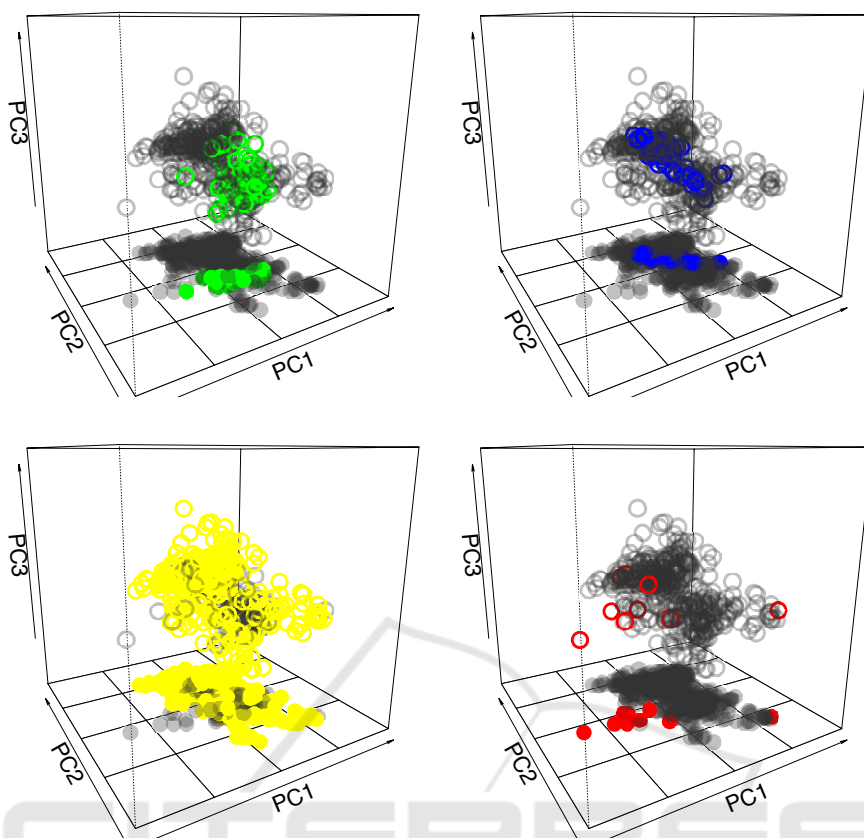
Figure 6: A random sample of size $N = 500$ is plotted over the first three robust principal components coordinates. From top left panel to bottom right panel, the colored blobs represent green, blue, yellow, and red zones respectively.

and the blue zones. In other words, the data that are dissimilar to the already observed data deserve further attention. These points are the potential outliers and are reported in the red zone. Additionally, this also reveals that the predictions values associated with such data can be misleading. Producing a noticeable percentage of samples belonging to the yellow and the red zones referring to the possible failure of the e-nose equipment.

# 7 CONCLUSION

Electronic nose devices have received continuous attention in the field of sensor technology recently. Applications of the e-nose appear in industrial production, processing, and manufacturing including quality control, grading, processing controls, gas leak detection, and monitoring odors. The measurement quality of the e-nose depends on its sensor's performance. Due to the high variability of the gases in the air and the sensitivity of the sensor values, e-nose measurements can fluctuate very often and fail to main-

tain a certain level of precision. An automatic procedure that detects the samples' validity in an online fashion has been a technical shortage and was addressed in this work. This allows administrators to take the subsequent steps like sampling new observations from the field or re-calibrating the system if necessary. Equipping the e-nose device with a computing server that performs the measurement validation and odor concentration prediction in real time, initiates a new era to automatic odor detection. Developing a suitable model for predicting odor concentration might be the next challenge of this emerging technology. We follow this direction in our future work.

# REFERENCES

Artursson, T., Eklov, T., Lundstrom, I., Martensson, P., Sjostrom, M., and Holmberg, M. (2000). Drift correction methods for gas sensors using multivariate methods. *Journal of chemometrics*, 14:711–723.

Bermak, A., Belhouari, S. B., Shi, M., and Martinez, D. (2006). Pattern recognition techniques for odor dis-

crimination in gas sensor array. *Encyclopedia of sensors*, X:1–17.

Brys, G., Hubert, M., and Rousseeuw, P. J. (2006). A robustification of independent component analysis. *Chemometrics*, 19:364–375.

Carlo, S. D. and Falasconi, M. (2012). Drift correction methods for gas chemical sensors in artificial olfaction systems: techniques and challenges. *Advances in chemical sensors*, pages 305–326.

Croux, C. and Haesbroeck, G. (2000). Principal components analysis based on robust estimators of the covariance or correlation matrix: Infulence functions and efficiencies. *Biometrika*, 87:603–618.

Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9:432–441.

Gardner, J. and Bartlett, P. (1994). A brief history of electronic noses. *Sens. Actuat. b: chem.*, 18:211–220.

Gutierrez-Osuna, R. (2002). Pattern analysis for machine olfaction : a review. *IEEE Sensors journal*, 2:189–202.

Hubert, M., Rousseeuw, P. J., and Branden, K. V. (2005). Robpca: A new approach to robust principal component analysis. *Thechnometrics*, 47:64–79.

Hyvarinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10:626 – 634.

Jolliffe, I. (2002). *Principal Component Analysis*. Springer.

Josse, J., Pagès, J., and Husson, F. (2011). Multiple imputation for principal component analysis. *Advances in data analysis and classifications*, 5:231–246.

Kermiti, M. and Tomic, O. (2003). Independent component analysis applied on gas sensor array measurement data. *IEEE, Sensors Journal, IEEE*, 3:218–228.

Li, G. and Chen, Z. (1985). Projection-pursuit approach to robust dispersion matrices and principal components: primary theory and monte carlo. *Journal of the american statistical association*, 80:759–766.

McGinley, P. C. and Inc, S. (2002). Standardized odor measurement practices for air quality testing. *Air and Waste Management Association Symposium on Air Quality Measurement Methods and Technology-San Francisco, CA*.

Padilla, M., Perera, A., Montoliu, I., Chaudry, A., Persaud, K., and Marco, S. (2010). Drift compensation of gas sensor array data by orthogonal signal correction. *Journal of chemometrics and Intelligent labrotory system*, 100:28–35.

Persaud, K. and Dodd, G. (1982). Analysis of discrimination mechanisms in the mammalian olfactory system using a model nose. *Nature*, 299:352–355.

Prendergast, L. (2008). A note on sensitivity of principal component subspaces and the efficient detection of influential observations in high dimensions. *Electronic Journal of Statistics*, 2:454–467.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288.

Zuppa, M., Distante, C., Persaud, K. C., and Siciliano, P. (2007). Recovery of drifting sensor responses by means of DWT analysis. *Journal of Sensors and Actuators*, 120:411–416.