Pattern Classification Techniques for Early Lung Cancer Diagnosis using an Electronic Nose

Rossella Blatt¹ and **Andrea Bonarini**¹ and **Elisa Calabró**³ and **Matteo Matteucci**¹ and **Matteo Della Torre**² and **Ugo Pastorino**³

Abstract. We present a method to diagnose lung cancer by the analysis of breath using an electronic nose. This device can react to a gas substance by providing signals that can be analyzed to classify the input. It is composed of a sensor array (6 MOS sensors, in our case) and a pattern classification process based on machine learning techniques. During the first phase of our research, we have evaluated the possibility and accuracy of lung cancer diagnosis by classifying the olfactory signal associated to exhalations of subjects. The second part of the research, still in progress, is aimed at assessing the possibility of discriminating also the different types and stages of the disease. At the end of the first phase, results have been very satisfactory and promising: we achieved an average accuracy of 92.6%, sensitivity of 95.3% and specificity of 90.5%. In particular we analyzed the breath of 101 individuals, of which 58 control subjects, and 43 suffer from different types of lung cancer (primary and not) at different stages. In order to find the components able to discriminate between the two classes 'healthy' and 'sick' at best, and to reduce the dimensionality of the problem, we have extracted the most significant features and projected them into a lower dimensional space using Non Parametric Linear Discriminant Analysis. Finally, we have used these features as input to several supervised pattern classification algorithms, based on different k-nearest neighbors (k-NN) approaches (classic, modified and Fuzzy k-NN), linear and quadratic discriminant classifiers and on a feed-forward artificial neural network (ANN). The observed results have all been validated using cross-validation. These results pushed us to begin the second phase of the project to investigate the possibility of early lung cancer diagnosis: we are involving a larger number of subjects, partioned in different classes according to the type and stage of the disease. The research demonstrates that the electronic nose is a promising alternative to current lung cancer diagnostic techniques: the obtained predictive errors are lower than those achieved by present diagnostic methods, and the cost of the analysis, both in money, time and resources, is lower. The introduction of this technology will lead to very important social and business effects: its low price and small dimensions allow a large scale distribution, giving the opportunity to perform non invasive, cheap, quick, and massive early diagnosis and screening.

Keywords: Electronic Nose, E-Nose, Olfactory Signal, Pattern Classification, Fuzzy *k*-NN, MOS Sensor Array, Lung Cancer.

1 Introduction

Lung Cancer causes 3000 deaths each day in the world. The surviving rate after 5 years of treatment is less then 10% in Europe and around 15% in the USA; these percentages increase to more than 50% if the cancer is discovered in its earliest stage. The diagnosis in an advanced stadium represents the main cause of the therapeutic failure. Current diagnostic techniques (e.g. TAC, PET) are invasive, very expensive, have a high risk of complications and a not so good accuracy; moreover, results on early detection and treatment, in last decades, have been very poor and unsatisfying. This calls for the necessity of a non invasive, accurate and cheaper diagnostic technique, able to identify the presence of lung cancer in its early stages. We found an indication of the solution in the Fundamental Principle of Clinical Chemistry, which affirms that every pathology changes people chemical composition, modifying the concentration of some chemicals in the human body. This is true also in lung cancer: it has been demonstrated that the presence of lung cancer alters the percentage of some volatile organic compounds (VOCs) in the human breath [7, 12]. This means that these VOCs can be considered as lung cancer markers: the analysis of the olfactory signal of patients breath and the recognition of these VOCs in it, allow to determine the presence of lung cancer.

An instrument that allows to acquire, detect and process the olfactory signal is the electronic nose, which mimics the non separative mechanism of human olfaction [1]. Nowadays the research on olfactory systems has become very lively, most of all because of the multitude of applications in which it has been successfully used (e.g., medical diagnosis, food control quality, mines detection, drug detection, environmental analysis etc.) [10].

We divided our study in two phases: during the first one, already completed, we have evaluated the possibility of distinguishing between healthy persons and lung cancer diseased patients only analyzing their breath; the second part of the project, still in act, invastigates the possibility of performing an early diagnose by the analysis of the olfactory signal to discriminate among different lung cancer types and stages.

2 Methodology

The experiment has been developed within the Italian MILD (Multicentric Italian Lung Detection) project, promoted by the Istituto Nazionale Tumori, Italy. We analyzed the breath of 101 volunteers, of which 58 healthy and 43 suffering from different types of lung cancer. All cases were hospitalized at the Istituto Nazionale Tumori of Milan. Among them, 23 have a primary lung cancer, while 20 of

¹ IIT Unit Artificial Intelligence and Robotics Laboratory AIRLab, Politecnico di Milano, Italy, email: blatt@elet.polimi.it, bonarini@elet.polimi.it, matteucci@elet.polimi.it

² Automation & Inspection Systems, SACMI Imola S.C., Italy, email: matteo.della.torre@sacmi.it

³ Toracic Surgery Department, Istituto Nazionale dei Tumori, Milano, Italy, email: elisa.calabro@istitutotumori.mi.it, ugo.pastorino@istitutotumori.mi.it



Figure 1. Basic functioning of an electronic nose. The olfactory signal is acquired by the sensor array and preprocessed to reduce the impact of any form of noise and to reduce the dimensionality of the problem. Then the best features are used to perform the classification of the signal.

them have different kinds of pulmonary metastasis. Control people have no pulmonary disease and have negative chest CT scan. The study has been approved from the Ethical Committee of the Institute and we asked everybody to sign an agreement for the participation to the study.

The breath acquisition has been made by inviting all volunteers to blow into a nalophan bag of approximately $400cm^3$. As the breath exhaled directly from lung is contained only in the last part of exhalation, we have decided to collect only the final portion of it. From each bag we took two measures, obtaining a total of 202 measurements, of which 116 correspond to the breath of healthy people and 86 to diseased ones.

In the second phase of the research we are involving a larger number of volunteers and we are partitioning diseased patients according to the type and stage of lung cancer. The experimental procedure is the same of the previous analysis, except that now the bag is directly connected to the electronic nose while the subject blows into the bag and we introduced a new system to better control humidity effects.

3 Acquisition and analysis of the olfactory signal

An electronic nose is an instrument able to detect and recognize odours, namely the volatile organic compounds present in an analyzed substance [1]. It is composed of an array of electronic chemical sensors with partial specificity able to convert a physical or chemical information into an electrical signal and of a pattern analysis system, able to recognize or classify odours. Each sensor reacts in a different way to the analyzed substance providing multidimensional data that can be considered as an olfactory blueprint of the substance itself. An electronic nose consists in three principal components (Figure 1): the first component regards the Gas Acquisition System, that is done through a sensor array that measures a given physical or chemical quantity; in this research we used an array of six MOS sensors (developed by Group Sacmi). This choice is due to the fact that MOS sensors are characterized by high sensitivity (in the order of parts per billion ppb), low cost, high speed response and a relatively simple electronics. This aspects take on great importance if we consider that most of the VOCs markers of lung cancer are present in the diseased people's breath in very small quantities, varying from parts per million to parts per billion. The MOS sensors react to gases with a variation of resistance [11]; in Figure 2 it is possible to see a typical response of a MOS sensor.

The second component concerns the *pre-processing and dimensionality reduction* phase: after the electronic nose has acquired the olfactory signal it is necessary to reduce the effect of humidity, to



Figure 2. Example of a MOS sensor response. Each measure consists of three main phases: before each measure the instrument inhales the reference air, showing in its graph a relatively constant curve; after this short period it

inhales the analyzed gas, producing a change of the sensors' resistance; finally the instrument returns to the reference line, ready for a new measure.

normalize the acquired signal and to manipulate the baseline. After pre-processing we performed dimensionality reduction to extract the most relevant information from the signal. We defined ten descriptors from the sensors' responses able to represent data characteristics in the most efficient way. In particular, these features have been based on the variation of resistance, the course of the curve, its derivative, its integral and its Fast Fourier Transform (FFT). Some of these features returned more than one value (like the FFT), for a total of 39 descriptors for each measurement. Considering that we used 6 sensors, each measure would be described by 234 descriptors. Among all features it has been necessary to find those able to maximize the informative components and, thus, to contribute to improve the accuracy of the classifier. For this reason we applied the Mann-Wilcoxon non-parametric test with a significance level equal to $\alpha = 0.0001$ to select only discriminat descriptors. The choice of using a nonparametric test instead of a parametric one, is due to a previous analysis of the features distribution and a Lilliefors test. In order to evaluate the discriminative ability of the combination of more features, we performed an Analysis of Variance (ANOVA) and several scatter plots. We found that the most discriminative features between the two classes 'healthy' and 'sick' were the following descriptors (R(t)) is the curve representing the resistance variation during the measurement and R_0 the value of the resistance at the beginning of the measurement - as indicated in Figure 2 -):

• *Delta*: resistance change of sensors during measurement:

$$\delta = R_0 - \min(R(t)) \tag{1}$$



Figure 3. The results of dimensionality reduction through PCA on the left and NPLDA on the right.

• *Classic*: the ratio between the reference line and the minimum value of resistance reached during the measurement:

$$C = R_0 / \min(R(t)) \tag{2}$$

• Relative Integral: calculated as:

$$I = \int R(t)/(t \cdot R_0) \tag{3}$$

• *Phase Integral*: the closed area determined by the plot of the state graph of the measurement [9]:

$$x = R; \quad y = dR/dt \tag{4}$$

• *Single Point*: the minimum value of resistance reached during the measurement.

$$S = \min(R(t)) \tag{5}$$

After feature selection we performed data projection: we considered Principal Component Analysis (PCA) [5] and Nonparametric Linear Discriminant Analysis (NPLDA) [6]: PCA transforms data in a linear way projecting features into the directions with maximum variance, the latter is based on nonparametric extensions of the commonly used Fisher's linear discriminant analysis [5]. It is important to notice that PCA does not consider category labels; this means that the discarded directions could be exactly the most suitable for classification purpose. This limit can be overcome by NPLDA, which looks for the projection able to maximize differences between different classes and minimize those intra-class. In particular, NPLDA removes the unimodal gaussian assumption by computing the between scatter-matrix S_b using local information and the k nearest neighbors rule; as a result of this, the matrix S_b is full-rank, allowing to extract more that c-1 features (where c is equal to the number of considered classes) and the projections are able to preserve the structure of the data more closely [6]. As evident from Figure 3, NPLDA is able to separate the projected features more clearly than PCA, which plot shows a more evident overlap of samples. This means that NPLDA is more suitable, for the considered problem, in terms of classification performance. Moreover, the plot and the obtained eigenvalues clearly indicated that only one principal component is needed.

After the most significative dimension has been obtained, it has been possible to perform *classification*, that represents the third main component of an electronic nose. We considered three families of classifiers: Nearest Neighbors Classifiers (*k*-NN), Linear and Quadratic Discriminant Function based Classifiers (LD and QD) and an Artificial Neural Network (ANN).

3.1 Nearest Neighbors

The basic idea of this simple and powerful algorithm is to assign a sample to the class of the k closest samples in the training set. This method is able to do a non linear classification starting from a small number of samples. The algorithm is based on a measure of the distance (in this case, the Euclidean one) between the normalized features, and it has been demonstrated [5], that the k-NN is formally a non parametric approximation of the Maximum A Posteriori MAP criterion. The asymptotic performance of this algorithm, is almost optimum: with an infinite number of samples and setting k=1, the minimum error is never higher than the double of the Bayesian error (that is the theoretical lower bound reachable) [4].

One of the most critical aspects of this method regards the choice of parameter k when having a limited number of samples: if k is too large, then the problem is too much simplified and the local information loses its relevance. On the other hand, a too small k leads to a density estimation too sensitive to outliers.

For this reason, in addition to the classic k-NN, we implemented two other versions of this technique: the Modified k-NN and the Fuzzy k-NN. In the former, k means the number of closest neighbors to look for (as in the classic k-NN), but all belonging to the same class. This dynamically modify the neighborhood according to the noise in the dataset. Fuzzy k-NN, a variation of the classic k-NN based on a Fuzzy logic approach [14], assigns a fuzzy class membership to each sample and provides an output in a fuzzy form. In particular, the membership value of unlabeled sample x to i^{th} class is influenced by the inverse of the distances from neighbors and their class memberships:

$$\mu_i(x) = \frac{\sum_{j=1}^k \mu_{ij} (\|x - x_j\|)^{\frac{-2}{m-1}}}{\sum_{j=1}^k (\|x - x_j\|)^{\frac{-2}{m-1}}}$$
(6)

where μ_{ij} represents the membership of labeled sample x_j to the i^{th} class. This value can be crisp or it can be calculated according to a particular fuzzy rule: in this work we defined a fuzzy triangular

Classifier	NER	TPR	TNR	PREC _{POS}	$PREC_{NEG}$
Classic 9-NN	90.1%	89.5%	90.5%	87.5%	92.1%
Confidence Interval	[85.7-94.5]	[85.3-93.8]	[86.0-95.0]	[81.6-93.4]	[86.8-97.4]
Modified 9-NN	91.1%	91.9%	90.5%	87.8%	93.7%
Confidence Interval	[86.8-95.4]	[87.9-95.9]	[86.0-95.0]	[81.9-93.7]	[89.1-98.4]
Fuzzy k-NN	92.6 %	95.3%	90.5%	88.2%	96.3%
Confidence Interval	[88.5-96.7]	[91.8-98.9]	[86.0-95.0]	[82.3-94.1]	[93.2-99.4]
LD	89.6%	96.5 %	84.5%	82.2%	97.0 %
Confidence Interval	[85.0-94.2]	[93.7-99.3]	[79.1-89.9]	[75.2-89.1]	[93.9-100]
QD	92.6 %	95.3%	90.5%	88.2%	96.3%
Confidence Interval	[88.5-96.7]	[91.8-98.9]	[86.0-95.0]	[82.3-94.1]	[93.2-99.4]
ANN	91.6%	91.9%	91.3793 %	88.8 %	93.8%
Confidence Interval	[87.4-95.8]	[87.9-95.9]	[87.0-95.8]	[84.1-93.4]	[88.2-99.4]

Table 1. Performance indexes and corresponding confidence intervals (CI=95%) obtained from considered algorithms. Features have been previously projected by NPLDA and only the first principal component has been kept for classification. For k-NN techniques, we considered k=1,3,5,9,101; for classic and modified k-NN we show the best achieved results (when k=9). On the contrary, Fuzzy k-NN led to the same results independently from k's values.

membership function with maximum value at the average of the class and null outside the minimum and maximum values of it. In this way, the closer the sample j is to the average point of class i, the closer its membership value μ_{ij} will be to 1 and vice versa. The parameter m determines how heavily the distance is weighted when calculating each neighbor's contribution to the membership value [8]; we chose m = 2, but almost the same error rates have been obtained on these data over a wide range of values of m.

3.2 Discriminant Functions Classifier

Classification based on discriminant functions represents a geometric approach where the features space is divided in c decision regions each one corresponding to a particular class. The classifier is represented as a family of discriminant functions $g_i(x)$ with only one output that minimizes a given cost function. We considered two types of discriminant functions: the linear (LD) and the quadratic one (QD). A classifier based on a linear discriminant function divides the features space by planes and it is therefore optimum when the problem is linearly separable. In any case, this technique is able to lead to good performances also when the problem is not linearly separable. We implemented the Minimum Distance to Means (MDM) approach, in which the representatives of each class have been calculated as the mean value of samples belonging to that class. This approach is very simple and lead to good generalization; the drawback is that it compresses all information in only one representative value. If the problem is not linearly separable, a quadratic discrimination function could be more suitable, as it has been verified also in this work.

3.3 Artificial Neural Network

Artificial Neural Networks (ANN) are non-linear statistical modeling tools that can be used to model complex relationships between inputs and outputs or to find patterns in data. It can be demonstrated that an ANN, given a sufficient number of sigmoidal neurons in the hidden levels, is able to approximate any non linear function on a compact set. Moreover ANNs asymptotically (with an infinite number of examples) approximate the a-posteriori probability as with the Bayesian classifiers [2].

One of the main drawbacks of this method regards the impossibility to decide a priori the best topology to use. This choice has therefore been made through an empirical approach. In particular, we chose to use a feedforward neural network with one hidden layer, in which inputs are the first principal component obtained by NPLDA and the output is a single neuron assuming the value 1 if the presence of the disease is detected and 0 otherwise. All neurons have a sigmoidal function as activation function. The net has been trained using the Resilient Backpropagation algorithm, based on the gradient descent approach, in which only the sign of the derivative is used to determine the direction of the weights update. This choice is due to the fact that this algorithm was able to offer the best compromise between the error on the validation and convergence. Finally, we set the number of neurons in the hidden layer equal to 3; this value has been obtained by training a set of networks with increasing number of hidden neurons and picking the smallest one with a good validation error. Since ANN's results depend on the values of the initialization, we trained the net 20 times and we choose the best configuration (according to the early stopping error) to evaluate the test set.

4 Results

The performance of the classifiers has been evaluated through the obtained confusion matrices and performance indexes, defined as:

- accuracy (*Non Error Rate NER*), that represents the probability of doing a generic correct classification;
- sensitivity (*True Positive Rate TPR*): the probability to classify a person as sick when this is true;
- specificity (*True Negative Rate TNR*): the probability of classifying a person as healthy when this is true;
- precision w.r.t. diseased people (*PREC_{POS}*): the probability that, having assigned a sample to the class of diseased people, it actually belongs to that class;
- precision w.r.t. healthy people (*PREC_{NEG}*): the probability that, having assigned a sample to the class of healthy people, it actually belongs to that class.

To obtain indexes able to describe in a reliable way the performances of the algorithms, it is necessary to evaluate these parameters on new and unknown data, validating the obtained results. Considering the not so big dimension of population and that for every person we had two samples, we opted for a modified Leave-One-Out approach where each test set is composed by the pair of measurements corresponding to the same person, instead of a single measure as would be in the normal Leave-One-Out method. Doing this way, we avoided that one of these two measures could belong to the training set, while using the other in the test set.

 Table 2.
 Confusion matrix obtained from Fuzzy k-NN and Quadratic

 Discriminant Functions algorithms. Positive samples correspond to diseased subjects, while negative samples represent healthy volunteers.

CONFUS	ION	TRUE LABELS		
MATR	IX	Positive	Negative	
ESTIMATED	Positive	82	11	
LABELS	Negative	4	105	

All implemented algorithms have demonstrated a good ability to discriminate the two classes 'healthy' and 'sick'. Performance indexes are reported in Table 1, where we considered the first principal component obtained from NPLDA. The first consideration regards the similarity of Modified and Classic k-NN: results are strongly comparable, but a slight improvement is shown by Modified k-NN. Moreover Modified k-NN is able to achieve the same performance as Classic k-NN with a lower k value. Another relevant consideration regards the robustness of Fuzzy k-NN to k changes: we considered different values of k (k=1,3,5,9,101), but the algorithm demonstrated to be robust to these changes, keeping its results invariant.

In diagnostic field, sensitivity is more important than specificity because it is more relevant to recognize correctly a sick person instead of a healthy one; in the same way, precision on negative samples is more important than precision on positive ones, because it is worse to classify a person as healthy when he or she is actually sick, than the opposite. Considering larger importance of sensitivity and precision w.r.t. healthy samples, we can affirm that the fuzzy k-NN and the quadratic classifier are the algorithms able to achieve best results for the considered problem. The confusion matrix obtained by these algorithms is shown in Table 2, where elements along the principal diagonal represent respectively the TruePositive (TP) and the TrueNegative (TN) values, while those off-diagonal are respectively the FalsePositive (FP) and the FalseNegative (FN) values.

Performing a Student's t-test between all pair of classifiers, no relevant differences emerged; this means that implemented classifiers' results are comparable for the considered problem.

5 Conclusion and Further Direction of Research

The use of an electronic nose as lung cancer diagnostic tool is reasonable if it gives some advantage compared to current lung cancer diagnostic techniques, namely Computed Axial Tomography (CAT) and Positron Emission Tomography (PET). Not only this is verified in terms of performances, as illustrated in Table 3, but also because the electronic nose, unlike the classical approaches, is a low cost, robust, small (and thus eventually portable), very fast and, above all, non invasive instrument. This means that this instrument allows a massive and quick lung cancer diagnosis.

In order to improve the sensors technology, the necessity to develop longer-lyfe and more stable sensors emerged. Moreover, the development of hybrid systems is desirable, in order to obtain both selective and sensitive sensors.

According to the classification techniques, our work could be improved evaluating other classification algorithms (as support vector machines, Bayesian approaches or other topologies of ANN), as well as improving the feature selection algorithm. It could be also very interesting to train the ANN in presence of noise, since it has been demonstrated that ANNs can compensate humidity, drift and temperature variation phenomenons [3] that affect olfactory signals.

According to the scientific literature, there are no studies on the variation of VOCs in the breath before and after the surgery: it may be interesting to evaluate the resolution of the disease due to surgery.

 Table 3. Comparison of lung cancer diagnosis performance reached with the electronic nose presented in this work and current diagnostic techniques (data from [13]).

Indexes	CAT	PET	E-Nose
Accuracy (NER)	Nd	Nd	92.6%
Confidence Interval			[88.5-96.7]
Sensitivity (TPR)	75%	91%	95.3%
Confidence Interval	[60-90]	[81-100]	[91.8-98.9]
Specificity (TNR)	66%	86%	90.5%
Confidence Interval	[55-77]	[78- 94]	[86.0-95.0]
PRECPOS	Nd	Nd	88.2%
Confidence Interval			[82.3-94.1]
PREC _{NEG}	Nd	Nd	96.3%
Confidence Interval			[93.2-99.4]

An ambitious research prospective regards the individuation of risk factors connected to lung cancer (as smoke or food).

After the promising results that we obtained in the first phase of our research, we believed that the most important prospective of research that we should had followed was the evaluation of performing early dioagnosis. We want to understand to what extent the electronic nose is able to detect lung cancer even when it is at its earliest stage. If the results will confirm our assumption, the social and economics effects will be of strong impact: the low price and small dimensions of the electronic nose, allow a large scale distribution, giving the opportunity to perform non invasive, cheap, quick, and massive early diagnosis and screening.

REFERENCES

- P.N. Bartlett and J.W. Gardner, *Electronic Noses: Principles and Appli*cations, Oxford Univ Press: Oxford, 1999.
- [2] C.M. Bishop, *Neural Networks for Pattern Recognition*, Clarendon Pr, 1995.
- [3] J. Brezmes, N. Canyellas, E. Llobet, X. Vilanova, and X. Correig. Application of artificial neural networks to the design and implementation of electronic olfactory systems, 2000.
- [4] R.O. Duda, P.E. Hart, and D.G. Stork, *Pattern Classification*, Wiley-Interscience, 2000.
- [5] K. Fukunaga, Introduction to Statistical Pattern Recognition, Academic Pr, 1990.
- [6] K. Fukunaga and JM Mantock, 'Nonparametric discriminant analysis(for pattern feature extraction)', *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5, 671–678, (1983).
- [7] SM Gordon, JP Szidon, BK Krotoszynski, RD Gibbons, and HJ O'Neill, 'Volatile organic compounds in exhaled air from patients with lung cancer', *Clinical Chemistry*, **31**(8), 1278–1282, (1985).
- [8] JM Keller, 'A fuzzy k-nearest neighbor algorithm', *IEEE Transactions m Systems, Man, and Cybernetics*, 15(4), 580–585, (1985).
- [9] E. Martinelli, C. Falconi, A. D'Amico, and C. Di Natale, 'Feature Extraction of chemical sensors in phase space', *Sensors and Actuators B: Chemical*, 95(1), 132–139, (2003).
- [10] HT Nagle, R. Gutierrez-Osuna, and SS Schiffman, 'The how and why of electronic noses', *Spectrum, IEEE*, 35(9), 22–31, (1998).
- [11] M. Pardo and G. Sberveglieri, 'E lectronic Olfactory Systems Based on Metal Oxide Semiconductor Sensor Arrays', *MRS BULLETIN*, 29(10), 703, (2004).
- [12] M. Phillips, R.N. Cataneo, A.R.C. Cummin, A.J. Gagliardi, K. Gleeson, J. Greenberg, R.A. Maxfield, and W.N. Rom, 'Detection of Lung Cancer With Volatile Markers in the Breath', *Chest*, **123**, 2115–2123, (2003).
- [13] R.M. Pieterman, J.W.G. van Putten, J.J. Meuzelaar, E.L. Mooyaart, W. Vaalburg, G.H. Koeter, V. Fidler, J. Pruim, and H.J.M. Groen, 'Preoperative Staging of Non-Small-Cell Lung Cancer with Positron-Emission Tomography', *The New England journal of medicine*, **343**(4), 254–261, (2000).
- [14] LA Zadeh, 'Fuzzy sets [J]', Information and Control, 8(3), 338–353, (1965).