



**HAL**  
open science

# Primary investigation of sound recognition for a domotic application using support vector machines

Mohamed El Amine Sehili, Dan Istrate, Jérôme Boudy

► **To cite this version:**

Mohamed El Amine Sehili, Dan Istrate, Jérôme Boudy. Primary investigation of sound recognition for a domotic application using support vector machines. SINTES 2010: 14th International Conference on System Theory and Control, University of Craiova, Oct 2010, Sinaia, Romania. pp.503-506. hal-01308127

**HAL Id: hal-01308127**

**<https://hal.science/hal-01308127v1>**

Submitted on 27 Apr 2016

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

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

# Primary Investigation of Sound Recognition for a domotic application using Support Vector Machines

M. A. Sehili<sup>1,2</sup>, D. Istrate<sup>1</sup>, and J. Boudy<sup>2</sup>

<sup>1</sup>LRIT-ESIGETEL, 1 Rue du Port de Valvins, 77210 Avon, France

<sup>2</sup>Telecom SudParis, 9 Rue Charles Fourier, 91000 Evry, France

**Abstract**—The advent of modern communications and the low cost of some kinds of devices have resulted in a desire to equip elderly peoples' homes with sensors to monitor their activities and be forewarned of abnormal situations. In such an environment, sound may represent a rich source of information that can be exploited and this is considered as one of the most ergonomic and least intrusive solutions. However, this solution is often adversely affected by noise that is to say, mostly sounds of a type not taken into account in the creation of this system. Several methods were used to make it possible to classify sounds. In this work we tested Support Vector Machines to classify sounds in a domotic environment.

## I. INTRODUCTION

Sound classification is a problem of pattern recognition where one aims to distinguish the class of a given sound from other classes. In a domotic environment, there are many kinds of daily sounds which require detection in order to obtain information about the status of elderly people and their activities. There are also some sounds considered as noise that the system should ignore. Speech is considered as one of the most informative sounds, it is by far the most important class for a telemonitoring system. In fact, a speech signal can carry useful information like emotions and may contain a distress expression. This is what has motivated researchers to attempt sound classification in a hierarchical fashion as in [7] where speech was first distinguished from other sounds before being transmitted to a second classification engine.

This research work take place in the framework of the Sweet-Home project which search to provide a domotic HMI based on direct/indirect Speech/Sound recognition. The aim of this project is the safety of the persons and of goods using audio techniques. The interesting sound classes for this project are everyday life sounds (door clap, phone ring, dshes sounds,...) and abnormal sounds (screams, glass breaking, object falls,...).

The problem of sound classification can be compared to that of speaker identification as both are a multiclass pattern recognition task and rely on extracting and modeling the relevant features from the signal in order to differentiate them. In recent years, several statistical methods which have been successfully used for speaker identification for example; Hidden Markov Models (HMMs), Gaussian Mixture Models (GMMs) [4] and Dynamic Time Warping (DWT); were used for sound classification. Previous work of the ANASON team applied GMMs to sound classification following the model

described above. A combination of two or more classification methods was also used like in [12] and [5].

Support Vector Machines (SVMs) is a hyperplane based method that has gained increasing attention in the pattern recognition community over the last few years and has been successfully applied to tasks like speaker identification and verification, and face recognition. From a theoretical point of view, this discrimination method is quite robust. For a linear classification problem, it attempts to choose a hyperplane that best separates data points from two classes. Moreover, it has been shown to perform a non-linear classification with accuracy via the use of appropriate Kernel functions. This makes the SVMs extremely valuable for the task of sound classification.

## II. SUPPORT VECTOR MACHINES

SVMs belong to the family of binary classifiers. That means that an SVM attempts to assign one of two labels to data points from two distinct classes. The goal is to assign the exact label to each point given a set of labeled examples used to train the classifier.

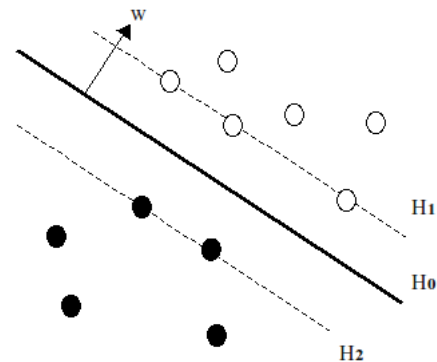


Fig. 1. Example of a linear classifier

The basic idea behind this method is to find a decision surface hyperplane which maximizes the margin between positive and negative examples. This implements the principle of structural risk minimization (SRM) [2] (Figure 1). The hyperplane  $H_0$  and the points which are mapped on it satisfy:

$$w \cdot x + b = 0 \quad (1)$$

The vector  $w$  is the normal to the hyperplane and  $b$  is the bias of the hyperplane from the origin. Given a set of  $N$  training examples  $(x_i, y_i)$ , where  $x_i \in R^p$  are the points and  $y_i \in \{-1, 1\}$  are the associated labels, we need to find the maximum margin subject to the constraints:

$$w \cdot x_i - b \geq 1$$

for  $y_i = 1$ , and

$$w \cdot x_i - b \leq -1$$

for  $y_i = -1$ , which can be written as:

$$y_i(x_i \cdot w + b) - 1 \geq 0, \forall i \quad (2)$$

We find that the distance between the two margins  $H_1$  and  $H_2$  is  $\frac{2}{\|w\|}$ . Thus, the problem can be stated as minimize  $\|w\|$  subject to (2).

The problem can be put as a quadratic programming problem as follows:

$$L_p = \frac{1}{2}\|w\|^2 - \sum_{i=1}^N \alpha_i y_i (x_i \cdot w + b) + \sum_{i=1}^N \alpha_i \quad (3)$$

where the  $\alpha_i$  are the Lagrange multipliers.

In figure 1 it can be seen that few examples can be found on the margins  $H_1$  et  $H_2$ . These are the support vectors and their associated  $\alpha_i$  are greater than 0.

In most cases the data examples are not perfectly separable. In other words, there exists no hyperplane that can separate all points without making any erroneous classification. This has motivated to introduce *slack* variables,  $\xi_i$ , to allow some degree of misclassification for some examples while still maximizing the distance to the nearest cleanly separated examples. The problem becomes:

Minimize:

$$\frac{1}{2}\|w\|^2 + C \sum_{i=1}^N \xi_i$$

subject to:

$$y_i(x_i \cdot w + b) \geq 1 - \xi_i, \forall i \quad (4)$$

where  $C$  is the penalty parameter of the error term.

The above theory works well as long as the data is linearly separable. In many problems, including sound classification, the data is far from being linearly separable. To deal with such problems one solution is to map the data into an extremely high dimensional feature space so that a linear separation becomes possible. However, dealing with data from a high dimensional feature space can easily lead to high computation costs [1]. This can be avoided by using Kernel functions. Typically used Kernel functions are:

$$\text{Linear:} \quad K(x, y) = x \cdot y \quad (5)$$

$$\text{Polynomial:} \quad K(x, y) = (\gamma x \cdot y + c)^p \quad (6)$$

$$\text{RBF:} \quad K(x, y) = \exp(-\Gamma|x - y|^2) \quad (7)$$

The final decision function takes the form:

$$f(x) = \sum_{i=1}^{N_{sv}} \alpha_i y_i K(x, x_i) + b \quad (8)$$

and the sign of the function  $f$  represents the label of the input vector  $x$ .

### III. APPLICATION TO SOUND CLASSIFICATION

In most cases a system has to deal with more than two classes of sounds. However SVM is a binary classification method. Although there exists a variant of SVM which can do multiclass classification, most researchers prefer splitting the problem into multiple binary problems and then using a binary classifier for each problem. There are two schemes most commonly used to do this; the one-to-all scheme and the one-to-one scheme. In the one-to-all scheme,  $C$  classifiers are created to represent  $C$  classes. Each classifier is trained by labeling examples from one class as +1 and examples from all the other classes as -1. An input example is thus evaluated using all the classifiers and is attributed to the class that yields the best distance. In the one-to-one scheme, a classifier is trained for each couple of classes and the final decision is achieved using a tree structure or a Directed Acyclic Graph (DAG) [6].

In most cases, a sound consists of more than one vector (i.e. frame). In [13] where SVMs are applied to speaker identification, the score of an utterance of  $N$  vectors is simply the arithmetic mean of the scores of the vectors it contains:

$$S = \frac{1}{N} \sum_{j=1}^N \left( \sum_i \alpha_i y_i K(x_j, x_i) + b \right) \quad (9)$$

Nevertheless we can also classify a sound using a majority voting on its vectors. This technique allow avoiding the influence of only some vectors missclassified.

Another way to use SVMs is to use an ensemble of classifiers. This may be very fruitful for sound classification especially when data is noisy. The idea is to obtain a set of classifiers for the same classification problem [13]. This can be achieved using bootstrapping or boosting [9].

#### A. Acoustical parameters

The SVM are not applied directly on the time signal but on spectral extracted vectors named acoustical parameters. The acoustical parameters can be the MFCC (Mel Frequency Cepstral Coefficients), LFCC (Linear Frequency Cepstral Coefficients), LPC (Linear Prediction Coefficients), LPCC (Linear Prediction Cepstral Coefficients), etc. In this paper we have used LFCC because are more adapted for for sound with high frequencies components. LFCCs are cepstral coefficients commonly used in speaker/speech recognition systems. They are commonly calculated as shown in figure 2 [10].

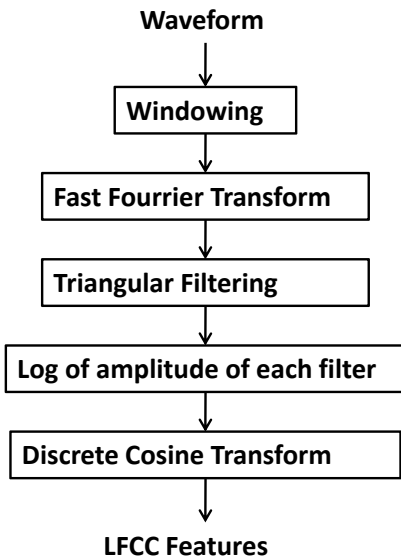


Fig. 2. Steps to derive LFCC

In the first step the signal is divided into frames, usually by using a rectangular windowing function at fixed intervals and overlap. Thus, each frames can be considered as a cepstral feature vector. The discrete Fourier Transform is then applied to each frame and triangular filter of uniformly spaced frequency bins are applied (Figure 3). The logarithm is computed on each output energy of each triangular filter. The components are finally decorrelated using the Discrete Cosine Transform. This has the advantage to reduce the final number of features in each vector.

#### IV. FIRST EXPERIMENTS

In order to experiment with SVMs for sound classification we have used the SVM-light library [8]. We first made a test on a part of the dataset created by the ANASON team. The dataset consists of seven categories of sound related to daily human activities. Table I shows the classes used in these experiments.

These sounds are 16kHz, 16 bits wav files. For this test, 24 order LFCCs (Linear Frequency Cepstral coefficients), energy and Zero Crossing Rate (ZCR) features were used. The frames were 16 ms of length with an overlap of 50%.

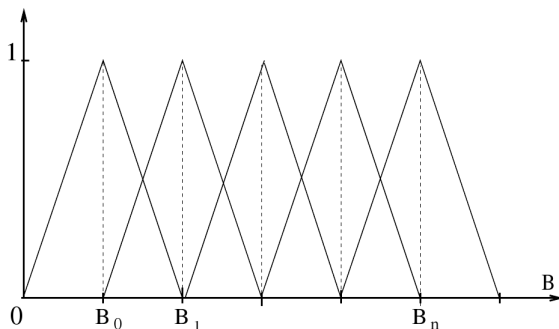


Fig. 3. Uniform frequency scale

TABLE I  
CLASSES OF SOUND FROM THE DATASET

Sound category	# of files
Cough	42
Door bell	14
Laugh	10
Sliding door	19
Sneeze	26
Snore	20
Yawn	21

The multiclass scheme used is the one-to-one, so a classifier is trained for each pair of classes. For each class, 50% of files are used for training and 50% for testing. To attribute a sound to a class, we first used the method consisting of calculating the sum the scores obtained by its vectors and then and then choosing the class according to the sign of the sum. This strategy yielded poor results. We then adopted a majority voting strategy which improved them moderately.

#### V. RESULTS

The proposed algorithm was evaluated on the data base through the good classified rate. The accuracy of the whole database is obtained by dividing the number of correctly classified files by the total number of files. Table II shows the results obtained.

The method used is time consuming because of the non-linear kernel (RBF in our case) where almost all training examples are retained as support vectors. This results in huge models.

In order to better use SVMs and and improve the performances, many methods can be used to train a model like the use of hold-out set or cross-validation [11]. In this work we used the techniques described in [3] which consist in scaling, grid search and cross-validation.

The goal of scaling is to constraint each feature value to be in a specific range, for example  $[-1, +1]$  or  $[0, 1]$ . This has the advantage to avoid features with greater values dominating those smaller values and to avoid numerical difficulties during calculation [3]. A grid search is used to find the couple of  $C$  and  $\Gamma$  which achieve the best accuracy on training data. Many combinations of these two parameters

TABLE II  
THE SCORES OBTAINED WITH THE FIRST TESTS USING TWO STRATEGIES OF CLASSIFICATION

	Classification strategy	
	Score sum	Majority voting
Cough	0.33	0.57
Door bell	0.57	1.00
Laugh	1.00	1.00
Sliding door	0.00	0.00
Sneeze	0.15	0.23
Snore	0.40	0.80
Yawn	0.18	0.18
Whole dataset	0.31	0.48

are thus used to train and test a classifier. One way to do this is to split the training data into two parts, train a classifier using one part and use the rest of data to determine which values of  $C$  and  $\Gamma$  allows for better performance.

A better way to determine the best parameters is to use cross-validation. In  $n$ -fold cross-validation the training dataset is split into  $n$  subsets of equal size. Each subset is then used to test the classifier trained on the other subsets. In our experiments we used 5-fold cross validation.

Tables III shows the results obtained after using the procedures above. It can be seen that these results outperform the previous one. Furthermore, in table III, and contrary to table II, the performances of the two strategies are almost comparable. This is due to scaling the data before training and test. We have also noticed that the models obtained after scaling the data are fairly of smaller size than those obtained with non scaled data. This may be very interesting for real-time systems as the time required to classify one vector is closely related to the size of the model.

## VI. CONCLUSIONS

This paper presents an application of SVMs to classify sound in a domestic environment. The sound classification is a multiclass problem but SVM are binary classifiers; two techniques was used one-against-one and one-against-all. The use of techniques like scaling and detecting the best parameters by using cross-validation allows to improve the performances. Although the first obtained results are encouraging, there are still several methods that can be used to better exploit SVMs and deal with the noise like the use of ensemble of classifiers with bootstrapping or boosting.

Next tests will aim to evaluate the noise influence on the SVM recognition performances and also the possibility to combine GMM with SVM in order to obtain a better system through score fusion.

## ACKNOWLEDGMENTS

We would like to thank the ANR (French National Research Agency) and, especially VERSO program, for funding the Sweet-Home project, the framework of this research activity.

TABLE III  
THE SCORES OBTAINED AFTER SCALING THE DATA AND USING  
CROSS-VALIDATION

	Classification strategy	
	Score sum	Majority voting
Cough	0.90	0.95
Door bell	1.00	1.00
Laugh	1.00	1.00
Sliding door	0.20	0.00
Sneeze	0.38	0.38
Snore	0.70	0.70
Yawn	0.18	0.18
Whole dataset	0.61	0.60

## REFERENCES

- [1] Joseph Picone Aravind Ganapathiraju. Hybrid svm/hmm architectures for speech recognition. 2000.
- [2] Burges Christopher J. C. A tutorial on support vector machines for pattern recognition. *Data Min. Knowl. Discov.*, pages 121–167, 1998.
- [3] Chih-Chung Chang Chih-Wei Hsu and Chih-Jen Lin. A practical guide to support vector classification. 2003.
- [4] Richard C. Rose Douglas A. Reynolds. Robust text-independent speaker identification using gaussian mixture speaker models. *IEEE Transactions on Speech and Audio Processing*, pages 72–80, 1995.
- [5] Zhou Xianzhong Luo Wen He Xin, Guo Ling. Hybrid support vector machine and general model approach for audio classification. In *ISNN '07: Proceedings of the 4th international symposium on Neural Networks*, pages 434–440, 2007.
- [6] Seong-Wan Lee Hyeran Byun. A survey on pattern recognition applications of support vector machines. *International Journal of Pattern Recognition and Artificial Intelligence*, pages 459–486, 2003.
- [7] D. Istrate J.E. Rougui and W. Soudiene. Audio sound event detection for distress situations and context awareness. *31st Annual International Conference of the IEEE EMBS*, pages 3501–3504, 2009.
- [8] Thorsten Joachims. Making large-scale support vector machine learning practical, 1998.
- [9] Hyun-Chul Kim, Shaoning Pang, Hong-Mo Je, Daijin Kim, and Sung Yang Bang. Constructing support vector machine ensemble. *Pattern Recognition*, 36(12):2757 – 2767, 2003.
- [10] Beth Logan. Mel frequency cepstral coefficients for music modeling.
- [11] John C. Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Large Margin Classifiers*, pages 61–74. MIT Press, 1999.
- [12] Hocine Bourouba Rafik Djemili, Mouldi Bedda. A hybrid gmm/svm system for text independent speaker identification. *International Journal of Computer Science and Engineering*, pages 22–28, 2007.
- [13] Zhaohui Wu Zhenchun Lei, Yingchun Yang. Ensemble of support vector machine for text-independent speaker recognition. *International Journal of Computer Science and Network Security*, pages 163–167, 2006.