

A Low-cost Approach for Detecting Activities of Daily Living using Audio Information: A Use Case on Bathroom Activity Monitoring

Georgios Siantikos, Theodoros Giannakopoulos and Stasinios Konstantopoulos

Institute of Informatics and Telecommunications, NCSR Demokritos, Aghia Paraskevi, Athens, Greece

Keywords: Audio Analysis, Activities of Daily Living, Health Monitoring, Remote Monitoring, Audio Sensors, RaspberryPI, Audio Event Recognition.

Abstract: In this paper, we present an architecture for recognizing events related to activities of daily living in the context of a health monitoring environment. The proposed approach explores the integration of a Raspberry PI single-board PC both as an audio acquisition and analysis unit. A set of real-time feature extraction and classification procedures has been implemented and integrated on the Raspberry PI device, in order to provide continuous and online audio event recognition. In addition, a tuning and calibration workflow is presented, according to which the technicians installing the device in a fast and user-friendly manner, without any requirements for machine learning expertise. The proposed approach has been evaluated against a particular scenario that is rather important in the context of any healthcare monitoring system for the elder, namely the "bathroom scenario" according to which a single microphone installed on a Raspberry PI device is used to monitor bathroom activity in a 24/7 basis. Experimental results indicate a satisfactory performance rate on the classification process (around 70% for five bathroom-related audio classes) even when less than two minutes of annotated data are used for training in the installation procedure. This makes the whole procedure non demanding in terms of time and effort needed to be calibrated by the technician.

1 INTRODUCTION

Although fully autonomous artificial intelligence is actively researched and advanced, the current state of the art (and at the level of maturity required for commodity electronics) has machine learning methods rely on delicate training and configuration sessions in order to adapt to different environments. When embedding machine learning methods in commodity electronics this is typically worked around by uploading the signal and receiving analysis results from remote centralized services. Recent examples include voice-operated personal assistant applications and companion, toy, and 'pet' applications.

This model, however, suffers from its obvious privacy implications. These implications are further exacerbated in the *telemedicine* domain for two reasons: the data collected by the remote service is not only more sensitive, but the users might also not be able to make informed decisions or might not be offered reasonable alternatives. Home monitoring for the elderly is a prime example: the increase in life expectancy and in the need for long-term care creates a pressure to seek alternatives to institutional healthcare for the

aged population. Advancements in robotics and automation and in artificial intelligence and intelligent monitoring are explored as a way to prolong independent living at home while providing guarantees of safety and adequate medical monitoring ((Barger et al., 2005), (Hagler et al., 2010), (MOTS et al., 2002)). The users of such solutions, however, might be suffering from mild cognitive impairment or be unable to afford conventional monitoring, which makes ethically questionable any consent they provide to upload and analyse raw content of their *activities of daily living (ADL)* in order to extract medical monitoring information. Several methods have been used to detect activities of daily living in real home environments, focusing on elderly population ((Vacher et al., 2013), (Vacher et al., 2010), (Costa et al., 2009), (Botia et al., 2012)) and a wide range of modalities.

In this paper, we present an audio analysis system (Section 2) that explores the integration of the audio sensor and the processing unit as Raspberry PI¹ device. Such a unit is able to execute signal processing and machine learning algorithms in order to eliminate

¹Please cf. <https://www.raspberrypi.org>

the need to provide raw content: the only information that leaves the confines of the integrated unit is an abstract ADL log. Although such information still needs to be managed in full accordance to guidelines pertaining private data, the level of obtrusiveness is greatly reduced by the assurance that no unwarranted analysis or recording can conceivably be done.

Our system is designed to satisfy two key requirements: that the analysis algorithms are computationally efficient so that they can be implemented for the Raspberry PI device; and that they can be tuned and configured for different acoustic environments by technicians without machine learning expertise. In order to evaluate the proposed approach on these requirements, we motivate and present a use case based on bathroom usage (Section 3) and draw conclusions (Section 3).

2 PROPOSED METHOD

2.1 Overall Architecture

The main part of the whole system is a microphone-equipped Raspberry PI single-board PC that is used for all data acquisition and processing. Its small-form factor, low energy consumption and low overall cost make it ideal for installing it in any room/area we want to monitor and its processing power is enough for running our algorithms in real time. In our experiments we used a Raspberry PI model B with a Wolfson audio card.

Communication to/from the PC is made using the MQTT machine-to-machine communication protocol. MQTT is a lightweight messaging protocol that implements the brokered publish/subscribe pattern, created widely used in IoT applications. Without going into technical details, the main idea is that when connected to a specified MQTT ‘broker’, various machines/applications can send messages under a certain topic and others can listen to these when ‘‘subscribed’’ to these topics. In our use case, it is used both for sending commands to the Raspberry PI (for example to start/stop recording) and for remotely receiving the processing results.

For this purpose, two MQTT clients were implemented: The first is installed in the Raspberry PI and is subscribed to a ‘‘command’’ topic in order to receive requests for collecting training data, building audio classes models and finally use them for real-time classification. The second one is bundled in an Android application and is used for sending remotely the corresponding commands and listening to the classification results. The system is designed with ease of use

in mind and the only set-up needed is connecting the two clients to the same broker. By having a dedicated broker this step can be performed automatically, making the whole system ‘‘plug-and-play’’.

2.2 System Calibration

Once setup, the system has to go through a training phase in order to be used for real-life scenarios. This includes recording, feature extraction, manual annotation of the recorded events and classifier tuning / training. Figure 1 shows the proposed calibration procedure. During this phase, the various events are recorded using the Android application as a remote controller of the Raspberry PI device that makes the actual recording and further processing. An audio file is created on user’s demand and the user/technician is informed about the categories and durations of already recorded data. He then provides the current recording’s label (e.g. ‘‘door bell’’). When a reasonably large amount of data is gathered (typically about 1-2 minutes of recordings for each category), the technician uses the mobile application to trigger the training process (that is also executed on the Raspberry PI device). After this process, the Raspberry PI is ready to monitor and recognize sound in the ‘‘learned’’ environment.

2.3 Audio Event Recognition

2.3.1 Audio Features

In total, 34 audio features are extracted on a short-term basis. This process results in a sequence of 34-dimensional short-term feature vectors. In addition, the processing of the feature sequence on a mid-term basis is adopted. According to that the audio signal is first divided into mid-term windows (segments). For each segment, the short-term processing stage is carried out and the feature sequence from each mid-term segment, is used for computing feature statistics (e.g. the average value of the ZCR). Therefore, each mid-term segment is represented by a set of statistics. In this Section we provide a brief description of the adopted audio features. For detailed description the reader can refer to the related bibliography (Giannakopoulos and Pikrakis, 2014) (Theodoridis and Koutroumbas, 2008), (Hyoung-Gook et al., 2005). The time-domain features (features 1 - 3) are directly extracted from the raw signal samples, while the frequency-domain features (features 4-34, apart from the MFCCs) are based on the magnitude of the Discrete Fourier Transform (DFT). The cepstral domain (e.g. used by the MFCCs) results after applying

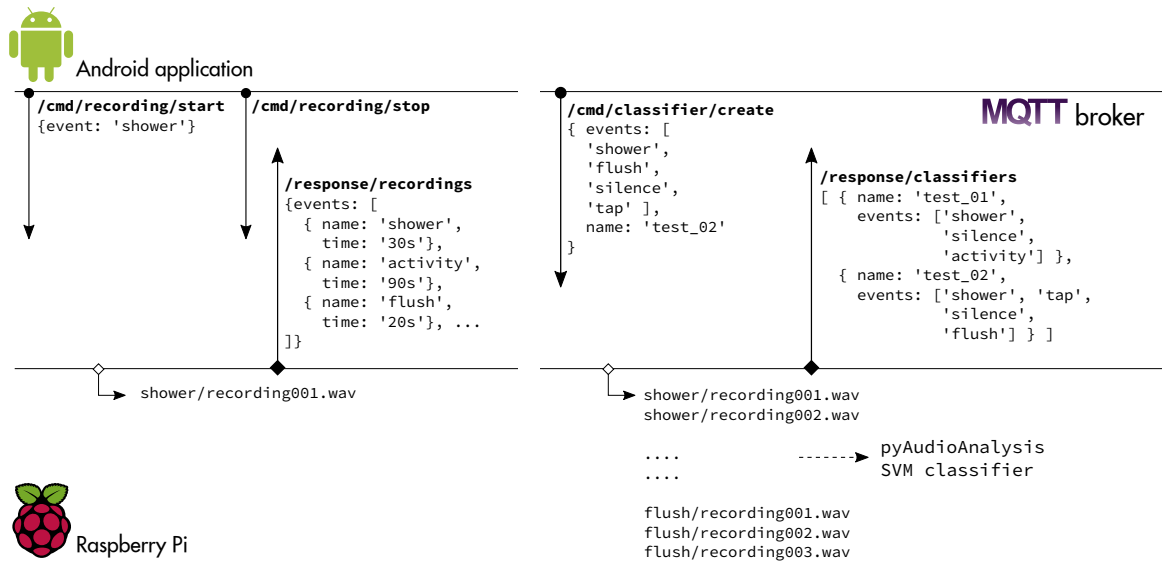


Figure 1: From left to right: User initiates an event recording and the corresponding file is created. When user stops, a response is returned with information about the events recorded so far. When a reasonable amount of data is gathered, an SVM classifier for the desired events can be created using the pyAudioAnalysis library. In this case, the response contains information about the classifiers available for future use.

the Inverse DFT on the logarithmic spectrum. The complete list of features is presented in Table 1.

2.3.2 Classification

As described in Section 2.3.1, the feature extraction process leads to a 68-dimensional feature vector for each 1-second audio segment, i.e. 2 statistics x 34 short-term features. Each unknown audio segment of fixed size (1 second in our case) is therefore represented by a 68-D feature vector. Each of these samples is classified using a Support Vector Machine with probabilistic output. We have selected to use probabilistic SVMs (Platt, 1999) due to their ability to generalize well especially in high dimensional classification problems (Chapelle et al., 1999). The model is trained using a cross-validation procedure to select the optimal SVM parameter, namely the soft margin parameter C .

2.3.3 Audio Analysis Implementation

Audio feature extraction and classification has been implemented using the pyAudioAnalysis library (Giannakopoulos, 15). This is an open-source Python library that implements a wide range of audio analysis functionalities and can be used in several applications. Using pyAudioAnalysis one can classify an unknown audio segment to a set of predefined classes, segment an audio recording and classify homogeneous segments, remove silence areas from a speech recording, estimate the emotion of a speech segment, ex-

tract audio thumbnails from a music track, etc. In this work, pyAudioAnalysis has been used to extract audio features, to train the classification models and to perform cross validation experimentation in order to extract the respective performance measures. pyAudioAnalysis achieves $2\times$ realtime performance on the Raspberry devices, which validates its usage in the context of the particular setup.

3 BATHROOM USE CASE AND EVALUATION

3.1 Use Case and Motivation

As discussed in the introduction, the motivating use case for our approach is medical monitoring. Specifically, we base our evaluation setup on allowing elderly people with mild cognitive impairment to maintain an independent life, at their own home, for longer than what is safely possible today.

In order to have a guideline about what information is used by medical doctors to assess such conditions, we use the *interRAI Long-Term Care Facilities Assessment System (interRAI LTCF)*. interRAI LTCF enables comprehensive, standardized evaluation of the needs, strengths, and preferences of persons receiving care. interRAI has been analysed previously in order to identify assessment items, such as mood and ADL logs, that can be automatically rec-

Table 1: Adopted short-term audio features.

Index	Name	Description
1	Zero Crossing Rate	Rate of sign-changes of the frame
2	Energy	Sum of squares of the signal values, normalized by frame length
3	Entropy of Energy	Entropy of sub-frames' normalized energies. A measure of abrupt changes
4	Spectral Centroid	Spectrum's center of gravity
5	Spectral Spread	Spectrum's second central moment of the spectrum
6	Spectral Entropy	Entropy of the normalized spectral energies for a set of sub-frames
7	Spectral Flux	Squared difference between the normalized magnitudes of the spectra of the two successive frames
8	Spectral Rolloff	The frequency below which 90% of the magnitude distribution of the spectrum is concentrated.
9-21	MFCCs	Mel Frequency Cepstral Coefficients: a cepstral representation with mel-scaled frequency bands
22-33	Chroma Vector	A 12-element representation of the spectral energy in 12 equal-tempered pitch classes of western-type music
34	Chroma Deviation	Standard deviation of the 12 chroma coefficients.

ognized and are useful to medical personnel (RADIO Project, 2015). Among the assessment items listed, we identified those regarding bathroom use as being closest to our concept: these items can be extracted by processing very sensitive content, and being able to provide guarantees about the management and processing of this content would have significant impact on the acceptance of any relevant solution by users.

In this context, we have recorded and manually annotated sounds from bathroom usage. In particular, the following audio classes are trained and evaluated by the proposed methodology:

- Silence - no sound
- Flushing water
- Shower
- Tap water
- Other activities

Note that the selected audio events are location-specific and therefore the adopted calibration workflow can be used in order during the installation phase, as described in Section 2.2.

3.2 Dataset

In order to train and evaluate the proposed event recognition methodology, we have recorded and manually annotated (using the mobile app described earlier in the paper) 4 recordings. The total duration of the dataset is 7 minutes. The extracted feature vector sequences and the respective ground truth of each recording is openly available at <https://iit.demokritos.gr/~tyianak/bathroomScenarioEvents.zip>

3.3 Experimental Evaluation

3.3.1 Performance Measures

Let CM be the confusion matrix, i.e. a $N_c \times N_c$ matrix (N_c is the total number of audio classes), whose rows and columns refer to the true (ground truth) and predicted class labels of the dataset, respectively. In other words, each element, $CM(i, j)$, stands for the number of samples of class i that were assigned to class j by the adopted classification method. The diagonal of the confusion matrix captures the correct classification decisions ($i = j$). CM is normalized row-wise, in order to discard the information that is related to the

size of each class:

$$CM_n(i, j) = \frac{CM(i, j)}{\sum_{n=1}^{N_c} CM(i, n)} \quad (1)$$

Obviously, after the normalization process, the elements of each row sum to unity.

Three useful performance measures are then extracted from the confusion matrix. The first is the overall accuracy, Acc , of the classifier, which is defined as the fraction of samples of the dataset that have been correctly classified:

$$Acc = \frac{\sum_{m=1}^{N_c} CM(m, m)}{\sum_{m=1}^{N_c} \sum_{n=1}^{N_c} CM(m, n)} \quad (2)$$

Apart from the overall accuracy, we have adopted two class-specific measures that describe how well the classification algorithm performs on each class. The first of these measures is the class recall, $Re(i)$, which is defined as the proportion of data with true class label i that were correctly assigned to class i :

$$Re(i) = \frac{CM(i, i)}{\sum_{m=1}^{N_c} CM(i, m)} \quad (3)$$

where $\sum_{m=1}^{N_c} CM(i, m)$ is the total number of samples that are known to belong to class i . In addition, we use the class precision ($Pr(i)$), i.e. the fraction of samples that were correctly classified to class i if we take into account the total number of samples that were classified to that class:

$$Pr(i) = \frac{CM(i, i)}{\sum_{m=1}^{N_c} CM(m, i)} \quad (4)$$

Finally, the F_1 -measure is also computed, which is the harmonic mean of the precision and recall values:

$$F_1(i) = \frac{2Re(i)Pr(i)}{Pr(i) + Re(i)} \quad (5)$$

3.3.2 Results

Table 2 shows the (row-wise normalized) confusion matrix and the respective precision, recall and F_1 measures for each audio class. This is the result of the evaluation process when only one recording is used during the training phase.

These results correspond to the most realistic and less demanding (in terms of calibration-training time). In addition, Table 3 demonstrates the ability of the classifiers to adopt to more data, if they can be available. In particular, 3 shows the same performance measures if a "leave one out" process is used in the evaluation process, using the described dataset. That is, if three whole recordings are used for each training phase. Results indicate an almost 5% performance

boosting. However, using three recordings instead of one means a 300% increase in the calibration time to be carried out by the technicians.

4 CONCLUSIONS

We have presented an architectural scheme that employs a Raspberry PI device both as an audio acquisition and analysis unit, in the context of a health monitoring system that detects ADLs in the living environments of elderly people. Real-time audio feature extraction and classification methods have been implemented and integrated on the device. Also, we propose a procedure for a fast and easy-to-use calibration procedure that actually trains the implemented classifiers in the particular home's sound conditions. Experimental evaluation has demonstrated a 70% classification performance even if a single recording (1 to 1.5 minute long) is used in the training process. The complete software that was used for the experiments can be found in the project's Git repository² under an open-source license.

The proposed system architecture satisfies three vital requirements.

- First the implemented algorithms are computationally efficient so that they can be implemented for the Raspberry PI device, as they demonstrate a $2 \times$ realtime performance on the Raspberry devices. This validates the system's usage in the context of a *low-cost* health monitoring setup as it does not require a workstation or a PC (e.g. (Chen et al., 2005)), but a single Raspberry PI that serves both as an acquisition and an analysis module. In particular, the total cost of both the acquisition and analysis modules is less than 100\$.
- In addition, despite the low-cost characteristics of the proposed approach, the system achieves a satisfactory classification performance, given (a) the low cost and (b) the lack of demand for big training data. Compared to other similar methods for ADL recognition in the context of a smart home environment, our method does not outperform (in terms of overall classification accuracy), however, given the significant differences in terms of overall cost, the proposed approach is preferable for real house applications. For instance, the approach in (Vuegen et al., 2013) achieves a 85% classification accuracy in a ADL recognition task, however the acquisition scenario requires multiple microphone sensors and therefore much higher cost.

²<https://bitbucket.org/radioprojectanalysis/ict4awe2016>

Table 2: Single-recording training: Row-wise normalized confusion matrix, recall precision and F1 measures. Overall F1 measure: 68.1%.

Confusion Matrix (%)					
True ↓	Predicted				
	Shower	Flush	Tap	Silence	Activity
Shower	89.4	1.8	3.0	0.1	5.8
Flush	7.2	70.7	0.4	2.1	19.6
Tap	5.7	4.0	85.8	0.8	3.6
Silence	1.4	4.0	0.0	58.0	36.5
Activity	13.0	11.6	2.6	31.2	41.6

Performance Measurements (% , per class)					
Recall:	89.4	70.7	85.8	58.0	41.6
Precision:	76.6	76.8	93.5	62.9	38.8
F1:	82.5	73.6	89.5	60.3	40.2

Table 3: Three-recording training: Row-wise normalized confusion matrix, recall precision and F1 measures. Overall F1 measure: 73.8%.

Confusion Matrix (%)					
True ↓	Predicted				
	Shower	Flush	Tap	Silence	Activity
Shower	85.6	1.6	2.6	0.2	10.0
Flush	5.7	86.5	0.0	1.5	6.3
Tap	5.2	3.7	85.5	0.7	4.9
Silence	0.1	3.4	0.0	72.5	24.0
Activity	5.8	9.7	1.5	29.0	53.9

Performance Measurements (% , per class)					
Recall:	85.6	86.5	85.5	72.5	53.9
Precision:	83.6	82.5	95.4	69.8	54.4
F1:	84.6	84.5	90.2	71.1	54.2

- Finally, the whole system can be configured and calibrated for different acoustic environments by technicians without machine learning expertise.

(more than 10%) (Lane et al., 2015), (Wang et al., 2015).

Our ongoing and future research work focuses on the following directions:

- Extend the calibration procedure so that it also takes into account a "base dataset", i.e. an initial classification scheme that is tuned in the context of the annotation process and not re-trained from scratch.
- Use long-term temporal knowledge to smooth the results of the classifier, based on prior knowledge regarding the events.
- The most important ongoing research direction focuses on adopting more robust audio classification approaches. We are currently working towards implementing deep learning methods for audio event recognition, which have been proved to achieve a significant performance boosting

ACKNOWLEDGEMENTS

This project has received funding from the European Unions Horizon 2020 research and innovation programme under grant agreement No 643892. Please see <http://www.radio-project.eu> for more details.

REFERENCES

- Barger, T. S., Brown, D. E., and Alwan, M. (2005). Health-status monitoring through analysis of behavioral patterns. *Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, 35(1):22–27.
- Botia, J. A., Villa, A., and Palma, J. (2012). Ambient assisted living system for in-home monitoring of healthy independent elders. *Expert Systems with Applications*, 39(9):8136–8148.

- Chapelle, O., Haffner, P., and Vapnik, V. N. (1999). Support vector machines for histogram-based image classification. *Neural Networks, IEEE Transactions on*, 10(5):1055–1064.
- Chen, J., Kam, A. H., Zhang, J., Liu, N., and Shue, L. (2005). Bathroom activity monitoring based on sound. In *Pervasive Computing*, pages 47–61. Springer.
- Costa, R., Carneiro, D., Novais, P., Lima, L., Machado, J., Marques, A., and Neves, J. (2009). Ambient assisted living. In *3rd Symposium of Ubiquitous Computing and Ambient Intelligence 2008*, pages 86–94. Springer.
- Giannakopoulos, T. (2015–). pyAudioAnalysis: Python audio analysis library: Feature extraction, classification, segmentation and applications. [Online; accessed 2015-04-27].
- Giannakopoulos, T. and Pikrakis, A. (2014). *Introduction to Audio Analysis: A MATLAB® Approach*. Academic Press.
- Hagler, S., Austin, D., Hayes, T. L., Kaye, J., and Pavel, M. (2010). Unobtrusive and ubiquitous in-home monitoring: a methodology for continuous assessment of gait velocity in elders. *Biomedical Engineering, IEEE Transactions on*, 57(4):813–820.
- Hyoung-Gook, K., Nicolas, M., and Sikora, T. (2005). *MPEG-7 Audio and Beyond: Audio Content Indexing and Retrieval*. John Wiley & Sons.
- Lane, N. D., Georgiev, P., and Qendro, L. (2015). Deep-ear: robust smartphone audio sensing in unconstrained acoustic environments using deep learning. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pages 283–294. ACM.
- MOTS, T. M., Linda Fraas OTR, M., and Kathleen Stanton MS, R. (2002). Elder acceptance of health monitoring devices in the home. *Care Management Journals*, 3(2):91.
- Platt, J. C. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in large margin classifiers*. Citeseer.
- RADIO Project (2015). D2.2: Early detection methods and relevant system requirements. Available at <http://radio-project.eu/deliverables>.
- Theodoridis, S. and Koutroumbas, K. (2008). *Pattern Recognition, Fourth Edition*. Academic Press, Inc.
- Vacher, M., Portet, F., Fleury, A., and Noury, N. (2010). Challenges in the processing of audio channels for ambient assisted living. In *e-Health Networking Applications and Services (Healthcom), 2010 12th IEEE International Conference on*, pages 330–337. IEEE.
- Vacher, M., Portet, F., Fleury, A., and Noury, N. (2013). Development of audio sensing technology for ambient assisted living: Applications and challenges. *Digital Advances in Medicine, E-Health, and Communication Technologies*, page 148.
- Vuegen, L., Van Den Broeck, B., Karsmakers, P., Vanrumste, B., et al. (2013). Automatic monitoring of activities of daily living based on real-life acoustic sensor data: a preliminary study. In *Fourth workshop on speech and language processing for assistive technologies (SLPAT): Proceedings*, pages 113–118. Association for Computational Linguistics (ACL).
- Wang, H.-H., Liu, J.-M., You, M., and Li, G.-Z. (2015). Audio signals encoding for cough classification using convolutional neural networks: A comparative study. In *Bioinformatics and Biomedicine (BIBM), 2015 IEEE International Conference on*, pages 442–445. IEEE.