

ASSESSMENT OF AUDIO FEATURES FOR AUTOMATIC COUGH DETECTION

Thomas Drugman, Jerome Urbain, Thierry Dutoit

TCTS Lab, Faculty of Engineering, University of Mons
 31, Boulevard Dolez - 7000 Mons - Belgium
 phone: + (32) 65 37 47 49, fax: + (32) 65 37 47 29, email: thomas.drugman@umons.ac.be
 web: <http://tcts.fpms.ac.be/~drugman/>

ABSTRACT

This paper addresses the issue of cough detection using only audio recordings, with the ultimate goal of quantifying and qualifying the degree of pathology for patients suffering from respiratory diseases, notably mucoviscidosis. A large set of audio features describing various aspects of the audio signal is proposed. These features are assessed in two steps. First, their intrinsic potential and redundancy are evaluated using mutual information-based measures. Secondly, their efficiency is confirmed relying on three classifiers: Artificial Neural Network, Gaussian Mixture Model and Support Vector Machine. The influence of both the feature dimension and the classifier complexity are also investigated.

1. INTRODUCTION

For children as well as for adults, cough is in pneumology the commonest syndrom. It is a daily and very frequent reason of seeking advices to the general practitioner (around 20% of consultations for children below 4 years old), the paediatrician and the pneumologist (for whom chronic cough represents one third of consultations). The impact of cough, notably chronic coughing, on life quality can be important.

The severity of cough can be evaluated by asking patients to fill in forms about their perception of the syndrom. However such a subjective assessment of cough has been shown [1] to be only slightly correlated to its objective characterization (using audio or video recordings for example). Medical literature on this topic therefore underlines the lack of a tool allowing the automatic, objective and reliable quantification of this symptom. This latter step is notably anterior to any correct evaluation of possible treatments.

Some approaches have been proposed to address the automatic detection of cough [2]. These systems generally couple various sensors to the audio signal (see [2] and references in it): air coupled microphones, accelerometer, lapel microphone, free field microphone, throat microphone or contact sensor. Although reported results are encouraging [2], there is currently no standardization and very few of these approaches led to a commercialization. In addition, following the patient in ambulatory and 24h-long conditions (while preserving his daily habits) remains an open problem.

As a result, cough quantification in the majority of hospitals is still nowadays performed by a tedious task of manual counting from audio recordings, or for validation by comparison using simultaneous video recordings.

This paper focuses on the automatic detection of cough using only the audio signal, as a preliminary and necessary study for its further integration within a multimodal system. On an acoustic point of view, cough is described as a forced expulsive manoeuvre against a closed glottis that is associ-

ated with a characteristic sound [3]. The main difficulty in detecting cough from audio recordings lies in its efficient discrimination with other audio non-cough events such as speech, laugh, or ambient noise.

The goal of this paper is to study which audio features and classifier are the most suited for automatic cough detection. For this, it is structured as follows. Section 2 proposes a large set of possible audio features for this purpose. The experimental protocol used for assessing these features is described in Section 3. Section 4 first evaluates their significance using mutual information-based measures. These features are then integrated within three classifiers in Section 5: Artificial Neural Network (ANN), Gaussian Mixture Model (GMM) and Support Vector Machine (SVM). Finally Section 6 concludes the paper.

2. AUDIO FEATURES FOR COUGH DETECTION

The various audio features that are used throughout this study are briefly presented in the following. These features can be divided into three categories: features describing the spectral contents, measures of noise, and prosody-related features. In our experiments, we also added the first and second derivatives for each of these features in order to integrate the sound dynamics. This leads to a total set of 105 descriptors whose relevance will be assessed in Sections 4 and 5.

2.1 Features Describing the Spectral Contents

Several features characterizing the spectral shape have been proposed in [4]. For a comprehensive description of the magnitude spectrum, the well-known *Mel Frequency Cepstral Coefficients* (MFCCs, [5]) are extracted. 13 MFCCs (including the 0th coefficient) are used to represent the spectral distribution within 13 perceptual sub-bands. Besides, several parameters describing the spectral shape are also employed. The *Spectral Centroid* is defined as the barycenter of the amplitude spectrum. Similarly, the *Spectral Spread* is the dispersion of the spectrum around its mean value. The *Spectral Decrease* is a perceptual measure quantifying the amount of decreasing of the spectral amplitude [4]. Finally, the *Spectral Variation* and *Spectral Flux* characterize the amount of variations of spectrum along time and are based on the normalized cross-correlation between two successive amplitude spectra [4].

2.2 Measures of Noise

Quantifying the level of noise in the audio signal is of interest for describing the cough sound. For this purpose, several measures are here suggested. First, the *Harmonic to Noise Ratio* (HNR) is calculated for the fre-

quency ranges [0-0.5kHz], [0-1.5kHz], [0-2.5kHz] and [0-3.5kHz] using the Voice Sauce toolkit freely available at: <http://www.ee.ucla.edu/~spapl/voicesauce/index.html>. These latter parameters are respectively denoted *HNR05*, *HNR15*, *HNR25* and *HNR35* in the remainder of this paper. The *Cepstral Peak Prominence* (CPP) is used as it has been shown to be correlated with the degree of breathiness in voice [6]. The *Spectral Flatness* measures the noisiness/sinusoidality of a spectrum (or a part of it). As suggested in [4], we here calculate the spectral flatness in the four following frequency bands: [0.25-0.5kHz], [0.5-1kHz], [1-2kHz] and [2-4kHz]. The *Zero-Crossing Rate* quantifies the number of times the signal crosses the zero axis. It is expected that the greater the amount of noise, the higher the amount of zero-crossing. As a last parameter quantifying the amount of noise in the audio signal, the *Chirp Group Delay* (chirp GD) is a phase-based measure proposed in [7] for highlighting turbulences during glottal production.

2.3 Prosody-related Features

In speech processing, prosody refers to the rhythm, stress and intonation of speech. It is generally reflected by clues such as volume, pitch and duration. We therefore use measures of energy and loudness which basically are informative mainly about the presence of audio activity. As it is known [3] that for a three-phase cough sound, the last phase presents voicing, the fundamental frequency is estimated using the STRAIGHT technique [8].

3. EXPERIMENTAL PROTOCOL

The database consists of audio signals captured by a cheap standard MP3 recorder in an hospital context. They were kindly provided by the belgian mucoviscidosis center at the Cliniques Universitaires Saint-Luc. Subjects are patients suffering from mucoviscidosis who had to spend a night at the hospital. The recorder was placed on their bedside table during the evening. Recordings then contain parasitical signals such as talking, laughing and TV, music or other types of noise, which can be confusing for detecting cough. The database is made of 5 minute-long recordings from 9 different patients, manually labeled in cough and non-cough segments.

Audio signals were downsampled from 44.1 kHz to 16 kHz. Features introduced in Section 2 were extracted every 10 ms on Hanning windows whose length is 25 ms. The relevance of these features is assessed in Sections 4 and 5. First, an evaluation based on the Mutual Information (MI) is led in Section 4. This approach is advantageous as it is independent of any classifier. A method of feature selection based on MI is employed to reduce dimensionality. In a second step, these features are assessed in Section 5 by being integrated within three classifiers: ANN, GMM and SVM.

4. MUTUAL INFORMATION-BASED ASSESSMENT AND FEATURE SELECTION

4.1 Background on Mutual Information

The problem of automatic classification consists in finding a set of features X_i such that the uncertainty on the determination of classes C is reduced as much as possible [9]. For this, Information Theory [10] allows to assess the relevance of features for a given classification problem, by making use

of the following measures (where $p(\cdot)$ denotes a probability density function):

- The entropy of classes C is expressed as:

$$H(C) = -\sum_c p(c) \log_2 p(c) \quad (1)$$

and can be interpreted as the amount of uncertainty on their determination.

- The mutual information between one feature X_i and classes C :

$$I(X_i; C) = \sum_{x_i} \sum_c p(x_i, c) \log_2 \frac{p(x_i, c)}{p(x_i)p(c)} \quad (2)$$

can be viewed as the information the feature X_i conveys about the considered classification problem, i.e. the discrimination power of one individual feature.

- The joint mutual information between two features X_i , X_j , and classes C can be expressed as:

$$I(X_i, X_j; C) = I(X_i; C) + I(X_j; C) - I(X_i, X_j; C) \quad (3)$$

and corresponds to the information that features X_i and X_j , when *used together*, bring to the classification problem. The last term can be written as:

$$I(X_i, X_j; C) = \sum_{x_i} \sum_{x_j} \sum_c p(x_i, x_j, c) \cdot \log_2 \frac{p(x_i, x_j) p(x_i, c) p(x_j, c)}{p(x_i, x_j, c) p(x_i) p(x_j) p(c)} \quad (4)$$

An important remark has to be underlined about the sign of this term. It can be noticed from Equation 3 that a positive value of $I(X_i, X_j; C)$ implies some **redundancy** between the features, while a negative value means that features present some **synergy** (depending on whether their association brings respectively less or more than the addition of their own individual information).

4.2 Mutual Information-based Assessment

To evaluate the significance of the audio features proposed in Section 2, the following measures are computed:

- the *relative intrinsic information* of one individual feature $\frac{I(X_i; C)}{H(C)}$, i.e. the percentage of relevant information conveyed by the feature X_i ,
- the *relative redundancy* between two features $\frac{I(X_i, X_j; C)}{H(C)}$, i.e. the percentage of their common relevant information,
- the **relative joint information** of two features $\frac{I(X_i, X_j; C)}{H(C)}$, i.e. the percentage of relevant information they convey together.

For this, Equations 1 to 4 are calculated. Probability density functions are estimated by a histogram approach. The number of bins is set to 50 for each feature dimension, which results in a trade-off between an adequately high number for an accurate estimation, while keeping sufficient samples per bin. Class labels correspond to the presence or not of a cough event.

Since 105 audio features were extracted, an exhaustive presentation of results cannot be detailed here. For the sake of clarity, Table 1 displays the MI-based values for the 14

	MFCC 0	Zero-Crossing	HNR05	MFCC 1	Chirp GD	MFCC 4	MFCC 3	MFCC 8	F0	HNR15	MFCC 5	MFCC 2	MFCC 6	Flatness 2-4
MFCC 0	47.2	59.0	56.3	57.5	55.4	54.7	54.2	53.6	54.2	54.0	53.0	53.0	52.9	52.9
Zero-Crossing	15.1	26.8	41.8	35.0	38.0	41.8	40.4	38.9	36.1	50.0	40.2	37.0	39.1	38.0
HNR05	7.6	1.8	15.9	45.0	37.6	35.9	28.1	29.8	32.7	37.9	32.1	30.4	30.2	28.0
MFCC 1	20.5	22.6	2.6	30.8	39.2	43.4	43.3	39.9	37.4	51.0	41.0	36.5	40.4	39.0
Chirp GD	18.0	15.0	5.4	17.8	26.2	43.7	42.9	35.2	33.0	46.1	38.7	32.2	37.1	31.8
MFCC 4	16.8	9.2	5.2	11.7	6.7	24.2	33.2	38.5	32.6	42.8	35.7	35.0	38.7	37.3
MFCC 3	4.1	-2.5	0.1	-1.4	-5.6	2.1	11.1	26.5	20.5	37.1	33.4	24.2	27.2	26.7
MFCC 8	7.4	1.7	0.9	4.7	4.8	-0.5	-1.6	13.8	22.9	38.9	34.9	26.3	29.3	25.2
F0	-1.6	-3.9	-11.1	-1.3	-1.4	-3.0	-3.9	-3.6	5.0	45.4	27.6	21.8	24.0	21.1
HNR15	22.9	6.7	6.9	9.7	10.0	11.4	4.2	4.9	-10.7	29.0	41.5	37.7	40.3	38.1
MFCC 5	15.3	7.7	5.9	10.9	8.6	9.9	-1.3	0.0	-1.1	9.7	21.1	30.1	29.6	30.4
MFCC 2	7.7	3.3	0.1	7.8	7.5	2.7	0.4	0.9	-2.9	5.9	4.4	13.5	28.2	24.9
MFCC 6	10.6	3.9	3.1	6.6	5.3	1.7	0.1	0.7	-2.3	6.1	7.8	1.5	16.2	27.1
Flatness 2-4	7.3	1.8	2.0	4.8	7.4	-0.1	-2.6	1.6	-2.7	4.9	3.6	1.5	2.1	13.0

Table 1: Mutual information-based measures for the 14 first selected features (respecting the ranking). *On the diagonal*: the relative intrinsic information. *In the bottom-left part*: the relative redundancy between the two considered features. *In the top-right part*: the relative joint information of the two considered features.

first features (respecting the ranking) selected by the algorithm that will be described in Section 4.3. The diagonal indicates the percentage of relevant information conveyed by each feature. It is worth noting that the selection technique accounts for the redundancy and synergy between features. Selected features are therefore not necessarily the ones presenting the highest individual discrimination power. In our results, we observed that features conveying the greatest relative intrinsic information are: MFCC 0 (47.25%), the loudness (46.56%), a measure of energy (39.88%), HNR35 (38.74%) and HNR25 (35.41%). The three first features are related to the signal energy and are particularly informative about the presence of an audio event. Although individually interesting, these features are strongly redundant, with e.g a value of 41.63% of relative redundancy between MFCC 0 and the loudness. A strong redundancy (30.84%) is also observed between HNR35 and HNR25. The algorithm of feature selection presented in Section 4.3 therefore tends to give priority to slightly redundant (or even synergic) features.

The top-right part of Table 1 contains the values of relative joint information of two features, while the bottom-left part shows the relative redundancy between two features. The best combination of two features is MFCC 0 with the zero-crossing rate, bringing together 59% of relative joint information. Inspecting the values of redundancy, it is worth observing that F0 extracted with STRAIGHT is synergic with all 13 other features. The set of 14 features is relatively weakly redundant, with a maximum relative redundancy of 22.9% between MFCC 0 and HNR15, and a maximum synergy value of -11.1% between F0 and HNR05. Note the absence of first or second derivative features in the selected subset.

4.3 Mutual Information-based Feature Selection

Several techniques of feature selection have been proposed in the literature [9]. An important category of such methods is the approach relying on mutual information [11]. Computing

MI from data requires the estimation of probability densities, which cannot be accurately done in high dimensions. This is why a majority of feature selection algorithms use measures based on up to three variables (two features plus the class label). Therefore, various MI-based strategies for feature selection have been proposed, all trying to deal with the issue of redundancy management. In this paper, we use the following algorithm which is known [11] to provide among the best results. Let us denote $F=\{X_1, X_2, \dots, X_N\}$ the initial set of N features, and S_k the selected subset (with $S_k \subseteq F$) of k features at step k . The method is a greedy algorithm which starts from an empty set and which selects at each step k the feature Y_k maximizing:

$$Y_k = \arg \max_{X_i \in F \setminus S_{k-1}} [I(X_i; C) - \max_{Y_j \in S_{k-1}} I(X_i; Y_j; C)] \quad (5)$$

considering that the redundancy between X_i and the selected subset S_{k-1} is dominated by the most redundant feature in it.

It is confirmed in Table 1 that selected features exhibit weak redundancy values, as it is penalized via the term in $I(X_i; Y_j; C)$ in Equation 5. It is also interesting to note that selected features arise from the three categories: prosody-related characteristics (MFCC 0 and F0), noise measures (zero-crossing rate, HNR05, chirp GD, HNR15 and flatness 2-4), as well as spectral-based parameters (MFCC 1 to 8). Since these features arise from complementary sources of information, it can be expected that redundancy has been appropriately taken into account.

5. CLASSIFIER-BASED ASSESSMENT

The use of three types of classifiers is here investigated: ANN, GMM and SVM. We rely on Matlab implementations for ANN and GMM, and on the Torch toolbox [12] for SVM. Evaluation is achieved using a 10-fold cross validation framework. This means that training is led on 90% of the

database (randomly chosen), and the 10% remaining are used for the test. This operation is repeated 10 times (with exclusive subsets for testing), so as to cover the whole database for the evaluation. The system is then generally assessed through its averaged error rate. However, given that the database is strongly unbalanced, i.e the proportion of cough events (compared to non-cough) is highly under-represented, we preferred to rely on Receiver Operating Characteristic (ROC) curves. A ROC curve shows the True Positive Rate (TPR, or sensitivity) as a function of the False Positive Rate (FPR, or 1-specificity) as a discrimination threshold θ is varied. As a single measure of performance of the ROC curve, we defined the Revised Error Rate (RER) as:

$$RER = \min_{\theta} \sqrt{(1 - TPR(\theta))^2 + FPR(\theta)^2} \quad (6)$$

Indeed, an ideal classifier being characterized by a $TPR = 100\%$ and a $FPR = 0\%$, a single measure of performance is the Euclidian distance from the top-left corner to the ROC curve. As a consequence, the lower RER, the better the system. This criterion implies that an equal importance is given to both TPR and FPR. Based on a medical advice, TPR or FPR could be emphasized by weighting its importance in Equation 6.

5.1 ANN-based Classification

An Artificial Neural Network (ANN) is a method of classification using an interconnected group of artificial neurons, and which allows a non-linear statistical modeling of the class posterior. It is here used for its ability to model complex relationships between inputs (audio features) and outputs (posterior probability of belonging to a given class).

Our ANN implementation relies on the Matlab Neural Network toolbox. The ANN is made of a single hidden layer with a variable number of neurons whose activation function is an hyperbolic tangent sigmoid transfer function. Figure 1 displays the evolution of the ROC curves as a function of the number of neurons using the 20 first selected features. It is observed that the performance increases with the number of neurons. For information, a RER of 13.5% is achieved with 2 neurons, 9.26% with 32 and 8.78% with 64 neurons.

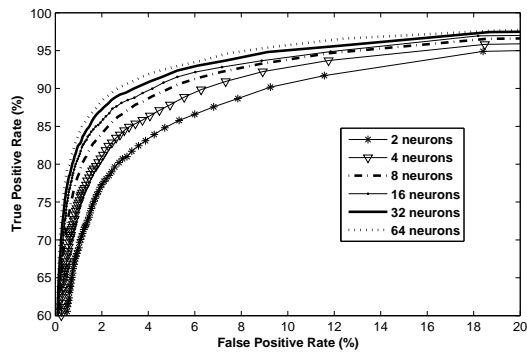


Figure 1: ROC curves obtained with the ANN classifier using 20 features and various numbers of neurons in the hidden layer.

The impact of the number of features on the classifier performance is illustrated in Figure 2, using 64 neurons in the

hidden layer. Performance with 5 or 10 features is largely under what is obtained with more than 20 features. However, ROC curves carried out with 20, 50 and 105 features are very close, with respective RERs of 8.78%, 8.13% and 7.94%. In other words, thanks to the efficient feature selection algorithm described in Section 4.3, using only 20 features gives similar results to what is reached with 105 features, allowing an important dimensionality reduction.

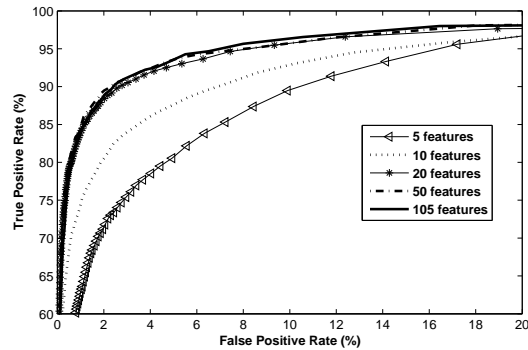


Figure 2: ROC curves obtained with the ANN classifier using 64 neurons in the hidden layer and various numbers of features.

For the best ANN configuration (64 neurons with 105 features), the following performance measures are obtained: **TPR=94.27%, FPR=5.50% and RER=7.94%**.

5.2 GMM-based Classification

A Gaussian Mixture Model (GMM) is a technique of classification in which the conditional probability for each class is approximated by a mixture of Gaussian distributions. In our Matlab implementation, GMMs are first initialized by a K-Means clustering step. The same number of Gaussians is used to model each class. Figure 3 plots the ROC curves using 20 features and various numbers of Gaussians in the mixture. It is observed that cough detection gets better with an increasing number of Gaussians.

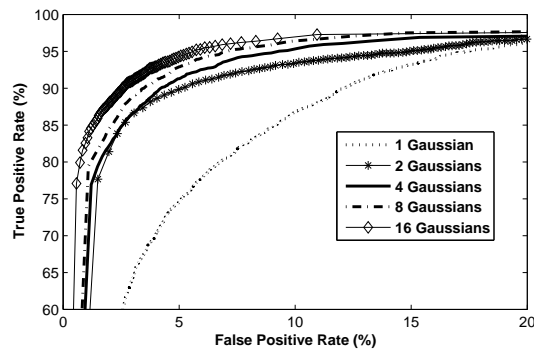


Figure 3: ROC curves obtained with the GMM classifier using 20 features and various numbers of Gaussians in the mixture.

In order to illustrate the influence of the feature dimension on the system, Figure 4 displays the evolution of RER as a function of the number of features using 8 Gaussians.

As it was the case for the ANN classifier, it turns out that using 20 features gives among the best results, and that the contribution when considering more features is minor.

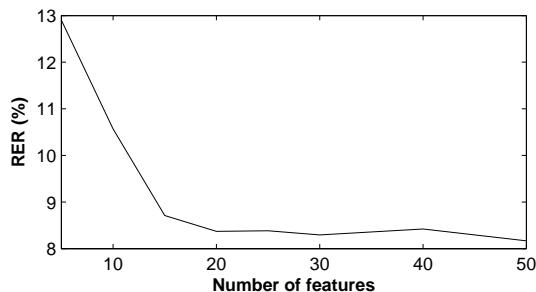


Figure 4: Evolution of RER as the number of features increases for a GMM classifier with 8 Gaussians.

For the best GMM configuration (16 Gaussians with 20 features), the following performance measures are obtained: **TPR=95.20%**, **FPR=5.73%** and **RER=7.48%**. First, it is interesting to note that for audio-based cough detection GMM outperforms ANN, with a reduction of 0.46% of RER. Secondly, it is worth emphasizing that only 20 features were used to reach that performance. For the same feature dimension, ANN obtains a RER of 8.78% with 64 neurons in the hidden layer.

5.3 SVM-based Classification

A Support Vector Machine (SVM) is a method of supervised learning able to analyze data and recognize patterns. It is here used as a non-probabilistic binary linear classifier. The initial feature space is mapped using a Gaussian kernel so as to maximize the final linear separability between classes. The criterion of good separability is that the hyperplane of decision should have the largest distance to the closest training data points of any class.

Experiments are here performed based on the SVM implementation available in the Torch toolbox. Using 20 features, as it was shown with ANN and GMM to convey almost all the information contained in the large feature set, we obtained the following performance measures: **TPR=81.87%**, **FPR=0.32%** and **RER=18.13%**. SVM is then clearly outperformed by the 2 other classifiers, the GMM approach providing the best identification rates.

6. CONCLUSION

This paper focused on the problem of cough detection relying only on audio recordings, as a preliminary and necessary study before integrating other sensors. A large set of features characterizing various aspects of the audio signal was proposed. These features were first assessed based on information theoretical measures, evaluating not only their intrinsic discrimination power, but also their redundancy and complementarity. Secondly, cough detection on recordings from patients suffering from mucoviscidosis was performed with three types of classifier: SVM, ANN and GMM. Impacts of feature dimension (reduced using a mutual information-based feature selection algorithm) as well as of the classifier complexity were analyzed. The best results were obtained with the GMM approach using only 20 features, reporting a sensitivity of 95.2% and a specificity of 94.3%.

7. ACKNOWLEDGMENTS

Thomas Drugman is supported by the “Fonds National de la Recherche Scientifique” (FNRS). The authors would like to thank the belgian mucoviscidosis center at the Cliniques Universitaires Saint-Luc for providing the audio recordings, as well as Pascaline Delchambre for her preliminary experiments. The authors would like to thank the Walloon Region, Belgium, for its support (grant WIST 3 COMPTOUX # 1017071).

REFERENCES

- [1] S. Decalmer, D. Webster, A. Kelsall, K. McGuinness, A. Woodcock, and J. Smith. Chronic cough : how do cough reflex sensitivity and subjective assessments correlate with objective cough counts during ambulatory monitoring? In *Thorax*, volume 62, pages 329–334, 2007.
- [2] J. Smith. Cough: Assessment and equipment. In *The Buyers Guide to Respiratory Care Products*, pages 96–101, 2008.
- [3] A. Morice, G. Fontana, M. Belvisi, S. Birring, and K. Chung et al. ERS guidelines on the assessment of cough. In *European Respiratory Journal*, volume 29, pages 1256–1276, 2007.
- [4] G. Peeters. A large set of audio features for sound description (similarity and classification) in the cuidado project. 2003.
- [5] F. Zheng, G. Zhang, and Z. Song. Comparison of different implementations of mfcc. In *J. Computer Science and Technology*, volume 16(6), 2001.
- [6] Y. Shue, G. Chen, and A. Alwan. On the interdependencies between voice quality, glottal gaps, and voice-source related acoustic measures. In *Interspeech Conference*, volume 34:37, 2010.
- [7] T. Drugman, T. Dubuisson, and T. Dutoit. Phase-based information for voice pathology detection. In *Int. Conf. on Acoustics, Speech and Signal Processing*, 2011.
- [8] H. Kawahara, H. Katayose, A. de Cheveigne, and R. Patterson. Fixed point analysis of frequency to instantaneous frequency mapping for accurate estimation of f0 and periodicity. In *Proc. Eurospeech*, volume 6, pages 2781–2784, 1999.
- [9] L. Huan and H. Motoda. Feature selection for knowledge discovery and data mining. In *The Springer International Series in Engineering and Computer Science*, volume 454, 1998.
- [10] T. Cover and J. Thomas. Elements of information theory. In *Wiley Series in Telecommunications, New York*, 1991.
- [11] T. Drugman, M. Gurban, and J-P Thiran. Relevant feature selection for audio-visual speech recognition. In *IEEE International Workshop on Multimedia Signal Processing*, 2007.
- [12] R. Collobert, S. Bengio, and J. Marithoz. Torch: a modular machine learning software library. In *Technical Report IDIAP-RR 02-46*, 2002.