

DETERMINATION OF THE NUMBER OF WIDEBAND ACOUSTICAL SOURCES IN A REVERBERANT ENVIRONMENT

Angela Quinlan⁽¹⁾, Frank Boland⁽¹⁾, Jean-Pierre Barbot⁽²⁾, Pascal Larzabal⁽²⁾

⁽¹⁾ School of Engineering, Trinity College Dublin, Ireland

⁽²⁾ Ecole Normale Supérieure de Cachan, 61, Av. du Président Wilson 94235 Cachan Cedex, France

ABSTRACT

This paper addresses the problem of determining the number of wideband sources in a reverberant environment. In [1] an Exponential Fitting Test (EFT) is proposed based on the exponential profile of the noise only eigenvalues. We consider the performance of this test for the problem in question, and compare it with the results achieved by the well known Akaike Information Criterion (AIC) and Minimum Description Length (MDL). Once reverberation is present in the received signals the EFT is seen to perform much better than the AIC and MDL.

1. INTRODUCTION

For any acoustical source localization scheme the initial step is to determine the number of signals received by the microphone array, based on a finite number of observed data samples. This process, which can be called detection or model order selection, is crucial to allow for good performance of high-resolution direction finding techniques.

Traditionally the number of sources is determined by firstly estimating the covariance matrix of the observed data samples, and then evaluating the multiplicity of the smallest eigenvalue of this matrix. One of the most widely-used approaches is that of Information Theoretic Criteria (ITC) [2]. The best known of this test family are the Akaike Information Criterion (AIC) [3] and the Minimum Description Length (MDL) [4].

However, these approaches are known to exhibit poor performance when there is only a small number of data samples available. In order to obtain accurate estimates of the number of sources in this difficult situation an Exponential Fitting Test (EFT) was proposed in [1], which exploits the ordered noise eigenvalue profile first introduced in [5]. In this paper we examine the performance of the EFT in the situation of determining the number of acoustical sources detected by an array of microphones, and compare the results obtained to those of the AIC and MDL tests. The EFT has previously been shown to outperform the AIC and MDL for the simulated case of a narrowband signal received in the presence of zero mean Complex White Gaussian Noise (CWGN) [5, 6]. In this paper we are concerned with determining the number of speakers in a moderately reverberant meeting room environment.

2. PROBLEM FORMULATION

We consider the model of an array of M microphones located in the sound field generated by d sources. Let $\mathbf{a}(\theta)$ be the steering vector representing the complex gains from a signal impinging on the M microphones with Direction of Arrival (DOA) θ . Then, if $\mathbf{x}(t)$ is the observation vector of size $M \times 1$, $\mathbf{s}(t)$ the emitted vector signal of size $d \times 1$ and $\mathbf{n}(t)$ the additive noise vector of size $M \times 1$, we obtain the following conventional model:

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) + \mathbf{n}(t) = \mathbf{y}(t) + \mathbf{n}(t), \quad (1)$$

where \mathbf{A} is the matrix of the d steering vectors, and $M > d$. The source covariance matrix $\mathbf{R}_s = E\{\mathbf{s}(t)\mathbf{s}^H(t)\}$ is assumed to be non-singular, which is equivalent to assuming that the sources are non-coherent.

The vector $\mathbf{n}(t)$ denotes zero-mean, spatially and temporally uncorrelated circular Gaussian complex noise, i.e. $E\{\mathbf{n}(t)\mathbf{n}^H(t)\} = \sigma^2\mathbf{I}$, $E\{\mathbf{n}(t)\mathbf{n}^T(t)\} = 0$. This corresponds to the noise being zero-mean and having a common variance σ^2 at all the sensors, and also being uncorrelated among all sensors. Such noise is termed spatially white, and is commonly assumed in DOA estimation schemes [7].

Thus, from equation (1), the observation covariance matrix \mathbf{R}_x can be expressed as:

$$\mathbf{R}_x = E[\mathbf{x}(t)\mathbf{x}^H(t)] = \mathbf{A}\mathbf{R}_s\mathbf{A}^H + \sigma^2\mathbf{I} \quad (2)$$

3. MODEL ORDER SELECTION

In the field of high resolution array signal processing a lot of work has been published concerning the model order selection problem.

3.1. Principle of statistical tests based on eigenvalue profile

According to (1), the noiseless observations $\mathbf{y}(t)$ are a linear combination of $\mathbf{a}(\theta_1), \dots, \mathbf{a}(\theta_d)$. Assuming independent source amplitudes $\mathbf{s}(t)$, the random vector $\mathbf{y}(t)$ spans the whole subspace generated by the steering vectors. This is the "signal subspace". Assuming $d < M$ and no array ambiguity, the signal subspace dimension is d . As a consequence, the number of non-zero eigenvalues of \mathbf{R}_y is equal to the number of sources d , with $(M-d)$ eigenvalues being zero.

Now, assuming $\mathbf{n}(t)$ meets the assumptions stated in the previous section, then according to (2), \mathbf{R}_x has the same eigenvectors as \mathbf{R}_y , with eigenvalues $\lambda_x = \lambda_y + \sigma^2$; for \mathbf{R}_x in the white noise case, σ^2 is a degenerate order $(M-d)$ eigenvalue. Then, arranging the eigenvalues in decreasing order, it becomes easy to discriminate between signal and noise eigenvalues.

In practice, \mathbf{R}_x is unknown and an estimate is made using $\hat{\mathbf{R}}_x = \frac{1}{N} \sum_{t=1}^N \mathbf{x}(t)\mathbf{x}^H(t)$. The "signal eigenvalues" are still identified as the d largest ones. But, with the statistical fluctuations in $\hat{\mathbf{R}}_x$, the noise eigenvalues are no longer equal to each other, and the separation between them is only clear in the case of high Signal to Noise Ratio (SNR), when a gap can be observed between signal and noise eigenvalues.

3.2. Classical Tests

For many years Information Theoretic Criteria (ITC) have been the most commonly used approach for the problem of detecting multiple sources. The best known of this test family are the Akaike

Information Criterion (AIC) [3] and the Minimum Description Length (MDL) [4].

The aim of the Akaike Information Criterion (AIC) is to determine the order of a model using information theory. The number of sources is the integer \hat{d} which, for $m \in \{0, 1, \dots, M-1\}$, minimizes the following quantity:

$$AIC(m) = -N(M-m) \log \left(\frac{g(m)}{a(m)} \right) + m(2M-m) \quad (3)$$

where $g(m)$ and $a(m)$ are respectively the geometric and arithmetic means of the $(M-m)$ smallest eigenvalues of the covariance matrix of the observation. The first term in equation 3 stands for the log-likelihood residual error, while the second is a penalty for overfitting. The Minimum Description Length (MDL) criterion [4] differs only in the penalty term (the last term in equation 4).

$$MDL(m) = -N(M-m) \log \left(\frac{g(m)}{a(m)} \right) + \frac{1}{2} m(2M-m) \log N \quad (4)$$

It appears that at low SNR and/or small number of snapshots, the results of ITC based tests degrade rapidly [8] and also perform very poorly with non-stationary signals in even moderately reverberant environments [9]. These approaches are therefore not applicable in the difficult scenario of determining the number of non-stationary sources in a reverberant environment.

4. EXPONENTIAL FITTING TEST

Few works have been concerned with model order selection under the constraint of a limited number of samples. Recently, however an Exponential Fitting Test (EFT) was proposed in [1], which continues to accurately determine the model order when the number of samples is limited, as is the case when dealing with non-stationary sources. This method is based on the profile of the ordered noise eigenvalues originally demonstrated in [5], and recalled here.

4.1. Eigenvalue Profile Under Noise Only Assumption

The sample covariance matrix of noise is $\hat{\mathbf{R}}_n = \frac{1}{N} \sum_{t=1}^N \mathbf{n}(t) \cdot \mathbf{n}(t)^H$. The distribution of the matrix $\hat{\mathbf{R}}_n$ is a Wishart distribution [10]. Finding the decreasing eigenvalue profile of $\hat{\mathbf{R}}_n$ is extremely difficult, if not impossible, and so instead the first and second order moments of the eigenvalues can be used to approximate the decreasing eigenvalue profile.

The error of the covariance matrix is denoted by Ψ , where $\Psi = \hat{\mathbf{R}}_n - \mathbf{R}_n = \hat{\mathbf{R}}_n - E\{\hat{\mathbf{R}}_n\} = \hat{\mathbf{R}}_n - \sigma^2 \mathbf{I}$. The noise eigenvalue profile can then be found by considering the first and second moments of $tr[\Psi]$. From $E\{tr[\Psi]\} = \mathbf{0}$, it can be seen that

$$M\sigma^2 = \sum_{i=1}^M \lambda_i \quad (5)$$

Note that

$$\Psi_{ij} = \frac{1}{N} \sum_{t=1}^N n_i(t) \cdot n_j^*(t) - \sigma^2 \delta_{ij},$$

and therefore $E\{\|\Psi_{ij}\|^2\} = \frac{\sigma^4}{N}$ in the case of white Gaussian complex circular noise. Since the trace of a matrix remains unchanged when the base changes, it follows that

$$E\{tr(\hat{\mathbf{R}}_n - \mathbf{R}_n)^2\} = \sum_{i,j} E\{\|\Psi_{ij}\|^2\} \quad (6)$$

$$= M^2 \frac{\sigma^4}{N} = \sum_{i=1}^M (\lambda_i - \sigma^2)^2 \quad (7)$$

The simple exponential decay model will be used to approximate the eigenvalue distribution:

$$\lambda_i = \lambda_1 r_{M,N}^{i-1}, \quad (8)$$

with $0 < r_{M,N} < 1$. In order to simplify the notation $r_{M,N}$ is denoted by r from now on. From equation (5) we get

$$\lambda_1 = M \frac{1-r}{1-r^M} \sigma^2 = MJ_M \sigma^2,$$

where:

$$J_M = \frac{1-r}{1-r^M}. \quad (9)$$

Considering that $(\lambda_i - \sigma^2) = (MJ_M r^{i-1} - 1) \sigma^2$, the relation (6) gives:

$$\frac{M+N}{MN} = \frac{(1-r)(1+r^M)}{(1-r^M)(1+r)}.$$

Setting $r = e^{-2a}$ ($a > 0$), this becomes:

$$\frac{M \cdot \tanh(a) - \tanh(Ma)}{M \cdot \tanh(Ma)} = \frac{1}{N},$$

where \tanh is the hyperbolic tangent function. An order-4 expansion gives the following biquadratic equation in a .

$$a^4 - \frac{15}{M^2+2} a^2 + \frac{45M}{N(M^2+1)(M^2+2)} = 0 \quad (10)$$

for which the positive solution is given by

$$a(M, N) = \sqrt{\frac{1}{2} \left\{ \frac{15}{M^2+2} - \sqrt{\frac{225}{(M^2+2)^2} - \frac{180M}{N(M^2-1)(M^2+2)}} \right\}}.$$

4.2. Test principle

The relations in the previous section can be extended to the case where the observations consist of d sources corrupted by additive noise. Under these conditions, eigenvalues associated with the covariance matrix can be broken down into two complementary subspaces: the source subspace \mathcal{E}_s (of dimension d) and the noise subspace \mathcal{E}_n (of dimension $Q = M - d$ for non coherent sources). Consequently the eigenvalue profile established in the previous section still holds if we replace M with Q . This will now be used to select the model order.

The basis of the EFT is to detect the eigenvalue index at which a break point appears between the observed eigenvalue profile and the theoretical one provided by the exponential model, valid only for noise eigenvalues.

For this purpose, a recursive test is performed on the noise subspace dimension, P , beginning with $P = 1$. Assuming that the last P eigenvalues are noise eigenvalues, we test to see if the previous eigenvalue (λ_{M-P}) is that of noise (assumption H_{P+1}) or of a signal (assumption \bar{H}_{P+1}). This is done by extending the previous relations to a noise subspace of dimension $P+1$:

$$\hat{\lambda}_{M-P} = (P+1) J_{P+1} \hat{\sigma}^2, \quad (11)$$

$$\text{with : } J_{P+1} = \frac{1 - r_{P+1,N}}{1 - (r_{P+1,N})^{P+1}} \quad (12)$$

$$\text{and : } \hat{\sigma}^2 = \frac{1}{P+1} \sum_{i=0}^P \lambda_{M-i} \quad (13)$$

Let us define the following two hypotheses:

$$H_{P+1} : \lambda_{M-P} \text{ is a noise eigenvalue.}$$

$$\bar{H}_{P+1} : \lambda_{M-P} \text{ is a signal eigenvalue.}$$

In order to choose between the above hypotheses, we predict the value of λ_{M-P} as explained before and compare the real value to the predicted one. The decision is made by comparing the relative error between the predicted and observed eigenvalues to a threshold:

$$H_{P+1} : \left| \frac{\lambda_{M-P} - \hat{\lambda}_{M-P}}{\hat{\lambda}_{M-P}} \right| \leq \eta_P \quad (14)$$

$$\bar{H}_{P+1} : \left| \frac{\lambda_{M-P} - \hat{\lambda}_{M-P}}{\hat{\lambda}_{M-P}} \right| > \eta_P \quad (15)$$

If this error is less than the corresponding threshold, the eigenvalue is determined to be a noise eigenvalue. Otherwise, it is a signal eigenvalue. The estimated dimension of the noise subspace is the first value of P for which H_{P+1} is chosen against \bar{H}_{P+1} . Consequently, the estimated dimension of the signal subspace is $\hat{d} = M - \hat{P}$.

4.3. Threshold Selection

The threshold is selected by considering the profile of the ordered eigenvalues in the noise only case [6]. Experimental recordings are taken of the background noise when there are no sources present. The threshold is then established as follows:

- The recordings are split into frames $\mathbf{n}(t)$, $t = 1, \dots, N$.
- The sample covariance matrix $\hat{\mathbf{R}}_n = \frac{1}{N} \sum_{t=1}^N \mathbf{n}(t) \mathbf{n}(t)^H$ is estimated for each frame, and the ordered eigenvalues, $(\lambda_1, \dots, \lambda_M)$ are computed.
- Using equation (11) the predicted eigenvalues, $(\hat{\lambda}_1, \dots, \hat{\lambda}_M)$ are then found for every $\hat{\mathbf{R}}_n$.
- The relative difference $\left| \frac{\lambda_m - \hat{\lambda}_m}{\hat{\lambda}_m} \right|$ is then found for each step $m = 1, \dots, M$, and the distribution of the relative differences is considered.

As these are the relative differences in the noise-only case, we would ideally like to set the threshold for each step greater than the maximum value found for this step. However, while reducing the probability of false alarm, such a threshold will increase the probability of non-detection. We therefore select a threshold value at each step that allows for a false alarm probability of 0.25%, with the overall false alarm probability across the four steps equal to 1%.

The use of experimental recordings for determination of the threshold allows the EFT to learn the statistics of the noise and reverberation levels in the room, resulting in the EFT being much greater robustness to the effects of non-gaussian noise than the AIC and MDL methods.

The ability to restrict the false alarm probability for the EFT is also an advantage, as both the AIC and MDL tend to over-estimate the correct model order, particularly in the presence of reverberation, resulting in a false alarm probability that is too high for many operational purposes and will lead to serious degradation of any subsequent localization steps.

5. COMPARISON WITH CLASSICAL TESTS

Denoting the true number of sources as d , and \hat{d} as the estimated number of sources the following criteria are used to evaluate the

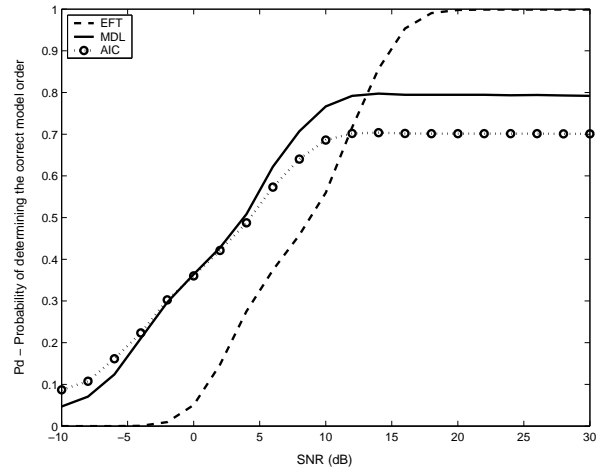


Fig. 1. Probability of detection for the EFT, the MDL and the AIC for simulated case of complex Gaussian White Noise and no reverberation.

performance of the tests:

$$\begin{aligned} \hat{d} = d &: \text{correct detection,} & P_d &= \text{Prob} \left[\hat{d} = d \right] \\ \hat{d} > d &: \text{over-estimation,} & P_{oe} &= \text{Prob} \left[\hat{d} > d \right] \\ \hat{d} < d &: \text{under-estimation,} & P_{ue} &= \text{Prob} \left[\hat{d} < d \right] \end{aligned}$$

5.1. Simulation Results in the Presence of White Noise

We begin by considering the performance of the three tests for simulations of the case of the two male speakers. The speech signals are received by an array of 5 microphones in the presence of Complex White Gaussian Noise, the sampling frequency $f_s = 16\text{kHz}$, and a frame length of 10 samples (0.625ms) with a 50% overlap between frames is used for the covariance estimation. In order to establish the threshold for the EFT, the relative difference between the noise-only eigenvalue distribution and the predicted distribution is considered for 10000 trials. In fig 1 it can be seen that for this situation that the EFT outperforms the AIC and MDL methods, and as the SNR increases detects the correct number of sources with a probability of 1. It should also be noted that the probability of over-estimating the number of sources P_{oe} is much lower than that of the AIC and MDL tests.

5.2. Room Response Simulation

The previous simulation does not take into the consideration the presence of reverberation, and therefore does not give an accurate indication of the performance that can be expected in a practical situation. We therefore consider the performance of tests in the presence of reverberation using room Enhanced Acoustic Simulator for Engineers (EASETM) simulation software. This software models the transmission of the signals from the specified sources to the microphone array, taking into account the dimensions and other acoustical characteristics of the room.

The responses are modelled for two sources located at angles of 70° and 110° from a Uniform Linear Array (ULA) of 5 omni-directional microphones with an inter-microphone spacing of 3.4cm. The horizontal distance between the array and the sources is 1.5m, allowing for the far-field assumption to be made, and the height of the microphone array and the sources

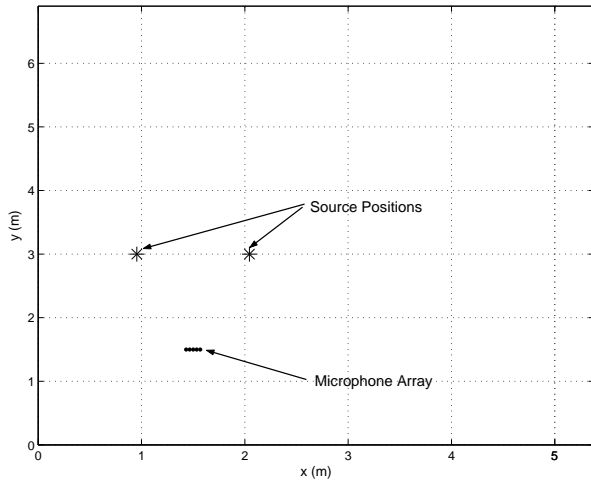


Fig. 2. Microphone and sound source positions.

are equal. The dimensions of the modelled room were $5.38m \times 6.9m \times 2.44m$, and the time required for the level of the mean-square sound pressure to decrease by $60dB$ after the source has been stopped, $T_{60} \approx 0.5s$. For the x - and y -coordinates of the sound source positions see fig. (2).

Once again the sampling frequency $f_s = 16kHz$, and for this case the frame length is 100 samples ($6.25ms$). The impulse responses produced by this software give a good indication of the effect of reverberation on the test, although they usually simulate higher reverberation levels than those encountered experimentally. The source signals used are two male speakers, and the additive background noise is simulated complex White Gaussian Noise, as in the previous case.

For this case the thresholds are established as described in section 4.3, however instead of using simulated white noise recordings taken in the experimental environment being modelled, when there were no sources present. A hilbert transform is then performed on the recorded noise signals, and these recordings are then used to find the eigenvalue distribution in the noise-only case.

We consider the performance of the tests as the Useful-to-detrimental ratio is increased. The Useful-to-detrimental ratio is a measure of the strength of the beneficial early-arriving reflections, to the later arriving sounds and the background noise [11]. The cut-off time between the beneficial and detrimental reflections depends on the room impulse response, and for rooms such as offices or meeting rooms has been shown to lie between $25 - 30ms$ [12]. In this paper, we consider the energy of the first $25ms$ of the impulse response as the early arrivals energy E , and the energy of the impulse response from $25ms$ to $500ms$ (T_{60}) as the energy of the late arrivals L . The power of the background noise present is N , and the Useful-to-detrimental ratio, U_{25} is defined as:

$$U_{25} = 10 \log \left\{ \frac{E}{L+N} \right\} dB.$$

From the results, shown in fig. (3) it can be seen that the introduction of the reverberation causes a deterioration in the accuracy of all three model order selection tests. However, the EFT can be clearly seen to outperform the AIC and MDL tests, both of which perform very poorly. The results of the AIC and MDL vary only slightly with increasing U_{25} , and it is clear that both these tests are highly unsuitable in the presence of reverberation.

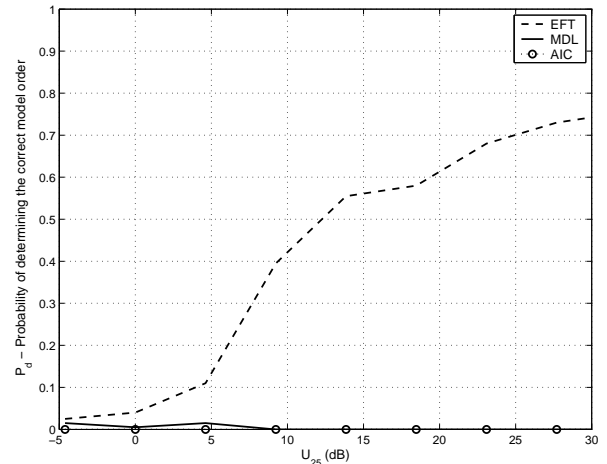


Fig. 3. Probability of detection for the EFT, the MDL and the AIC using room simulator EASE, as the Useful-to-Detrimental Ratio, U_{25} , is increased.

	P_d (%)	P_{oe} (%)	P_{ue} (%)
EFT	74.04	0	25.96
AIC	0.88	99.12	0
MDL	2.36	97.64	0

Table 1. Results found by EFT, AIC and MDL tests using experimental recordings of two different male speakers.

6. EXPERIMENTAL RESULTS

In this section the performance of the EFT, MDL and the AIC methods are compared for the experimental situation described in the previous section. Firstly we consider the case of two male speakers, then the case of the two female speakers and finally the case of one male and one female speaker. The thresholds are once again established using recordings of the noise taken when there are no sources present, as described in section (5.2). Tables 1- 3 respectively, show the results for the cases of: two male speakers, two female speakers, and one male and one female speaker.

The results from the experimental recordings confirm those of the simulated room response. Once again the AIC and MDL methods perform very poorly in the presence of reverberation, and they continuously over-estimate the true number of sources. The EFT greatly outperforms the other two methods, showing its suitability for determining the number of speakers in such an environment.

7. CONCLUSIONS

In this paper we have considered the performance of three model order selection tests, the EFT, the AIC and the MDL, for the problem of determining the number of speakers in a moderately reverberant environment. We have shown that in the absence of reverberation all three tests perform well for SNR levels of higher than $5dB$. Using EASE (Enhanced Acoustic Simulator for Engineers) room simulation software, we have then shown that once reverber-

	P_d (%)	P_{oe} (%)	P_{ue} (%)
EFT	54.74	0	45.26
AIC	0	100	0
MDL	2.46	97.54	0

Table 2. Results found by EFT, AIC and MDL tests using experimental recordings of two different female speakers.

	P_d (%)	P_{oe} (%)	P_{ue} (%)
EFT	62.61	0	37.39
AIC	0.59	99.41	0
MDL	3.56	96.44	0

Table 3. Results found by EFT, AIC and MDL tests using experimental recordings of a male and a female speaker.

ation is introduced the performance of all three tests is reduced, and the EFT greatly outperforms the AIC and MDL methods, both of which consistently overestimate the number of sources present. These results were then confirmed using experimental recordings.

8. ACKNOWLEDGMENTS

Thanks to Dennis Mc Carthy for his helpful discussions on reverberation. This work is supported by a post-graduate award from the Irish Research Council for Science Engineering and Technology (IRCSET).

9. REFERENCES

- [1] A. Quinlan, J. P. Barbot, and P. Larzabal, "Automatic determination of the number of targets present when using the time reversal operator (TRO)," *Journal Acoustical Society of America (JASA)*, Accepted for publication.
- [2] S. Valaee and P. Kabal, "An information theoretic approach to source enumeration in array signal processing," *IEEE Trans. Signal Processing*, vol. 52, pp. 1190–1196, 2004.
- [3] H. Akaike, "A new look at the statistical model identification," *IEEE Trans. Automat. Contr.*, vol. 19, no. 6, pp. 1361–1373, 1974.
- [4] J. Rissanen, "Modeling by shortest data description length," *Automatica*, vol. 14, pp. 465–471, 1978.
- [5] J. Grouffaud, P. Larzabal, and H. Clergeot, "Some properties of ordered eigenvalues of a Wishart matrix: Application in detection test and model order selection," in *Proc. ICASSP*, 1996, pp. 2463–2466.
- [6] A. Quinlan, J. P. Barbot, P. Larzabal, and M. Haardt, "Model order selection for short data: An exponential fitting test (EFT)," *EURASIP JASP (European Journal Applied Signal Processing)*, submitted for publication.
- [7] B. Ottersten, M. Viberg, P. Stoica, and A. Nehorai, "Exact and large sample maximum likelihood techniques for parameter estimation and detection in array processing," in *Radar*

Array Processing, S. Haykin, J. Litva, and T. J. Shepherd, Eds. Berlin: Springer-Verlag, 1993, ch. 4, pp. 99–151.

- [8] A. P. Liavas and P. A. Regalia, "On the behavior of information theoretic criteria for model order selection," *IEEE Trans. Signal Processing*, vol. 49, no. 8, pp. 1689–1695, 2001.
- [9] H. Teutsch and W. Kellerman, "Estimation of the number of wideband sources in an acoustic wave field using eigenbeam processing for circular apertures," in *Proc. Elektronische Sprachsignalverarbeitung (ESSV)*, Karlsruhe, Germany, 2003.
- [10] N. L. Johnson and S. Kotz, in *Distributions in Statistics: Continuous Multivariate Distributions*. New York: John Wiley, 1972, ch. 38–39.
- [11] J. S. Bradley, R. D. Reich, and S. G. Norcross, "On the combined effects of signal-to-noise ratio and room acoustics on speech intelligibility," *Journal Acoustical Society of America (JASA)*, 1999.
- [12] H. Golzer and M. Kleinschmidt, "Importance of early and late reflections for automatic speech recognition in reverberant environments," in *Proc. IEEE Workshop on Application of Signal Processing to Audio and Acoustics*, New Paltz, NY, 2005.