# Sphericity Minimum Description Length: Asymptotic Performance under Unknown Noise Variance

Josep Font-Segura

Universitat Pompeu Fabra (UPF)

Roc Boronat 138, 08018 Barcelona, Spain

josep.font@upf.edu

Jaume Riba, and Gregori Vázquez

Technical University of Catalonia (UPC)

Jordi Girona 1-3, 08034 Barcelona, Spain

{jaume.riba, gregori.vazquez}@upc.edu

*Abstract*—**This paper revisits the model order selection problem in the context of second-order spectrum sensing in cognitive radio. Taking advantage of the recent interest on the generalized likelihood ratio (GLR), the asymptotic performance of the minimum description length (MDL) rule under unknown noise variance is addressed. In particular, by exploiting the asymptotically Chi-squared distribution of the GLR, a complete characterization of the error probability is reported, instead of approximating only the missed-detection probability as done in the literature.**

*Index Terms*—**Model order selection, minimum description length, generalized likelihood ratio, noise uncertainty.**

## I. INTRODUCTION

**M**ODEL ORDER SELECTION is the signal processing problem that consists of determining the dimension of the parameter vector of the data model [1]–[3], and has many applications from radar to biomedicine.

Model order selection is cast a composite hypothesis testing problem among all the possible dimension of the parameter vector. As a composite hypothesis testing problem, the estimation of the dimension is driven by the *detection* of a signal satisfying a probability density function (PDF), known up to some unknown parameters. In order to circumvent the unknown parameters, one approach consists of maximum likelihood estimate (MLE) the unknowns and formulate the corresponding generalized likelihood ratio (GLR) [4]. Unluckily, this leads to a systematic overestimation of the dimension of the parameter vector, for which the introduction of dimension-dependent penalty terms is required. Among others, the minimum description length (MDL) [5], also known as the Bayesian information criterion (BIC) [6], is the most commonly adopted model order selection statistic due to its motivation in the information theory community and its tradition in the coding theory. Even though the choice of MDL for the present work has information theoretical fundamentals,

the results obtained in this paper can be extrapolated to other model order selection statistics [7], as the likelihood ratio is a common denominator.

It is worth noting that both MDL and BIC have received recent attention surrounding the discussion of consistency [8], [9] within the signal-to-noise ratio (SNR), or the asymptotic equivalences within the recently proposed exponentially embedded family (EEF) rule [10]. The problem of spectral occupancy estimation has also been recently pointed out in [11], where a maximum a posteriori (MAP) estimation of the signal rank with known noise variance is addressed. Historically, MDL and other information criterion rules had their major role in the problem of estimating the number of sources in an array of sensors [12]. In [12], the detection problem is first presented as a minimization over the MDL rule, instead of the at that time conventional composite hypotheses testing approach. Even though the MDL has been one of the core rules in classical signal processing problems such as model order selection in autoregressive models [13] and sinusoids [14], the asymptotic performance analysis has not been exhaustively addressed.

By considering that the difference between two MDL statistics is normally distributed, the performance of the MDL has been characterized in [15], [16], and more recently in [17] and [18]. As a common denominator in the aforementioned works, the error probability is characterized by the missed-detection probability, i.e., the probability of underestimating the model order, because the MDL tends to underestimate. More specifically, the missed-detection probability is further approximated by the probability of the single event $\text{MDL}_{n-1} < \text{MDL}_n$, i.e., the probability that the MDL rule decides for the model order immediately inferior to the correct order. In this work, on the other hand, the MDL is formulated from a likelihood ratio instead as from a likelihood function. As a result, since the GLR part of the MDL rule is known to be asymptotically Chi-squared distributed, this contributed to a more accurate characterization of the MDL for Gaussian sources. Furthermore, this allows to characterize the complete error probability, i.e., the analysis takes into consideration all the missed-detection and false-alarm as error events.

## II. PROBLEM FORMULATION

This paper formulates the model order selection problem of a zero-mean Gaussian signal with unknown low-rank correlation matrix immersed in zero-mean white Gaussian noise with unknown noise variance. This problem naturally arises in multi-measurement passive radar, rank estimation or cognitive radio problems [19], [20]. As in signal detection, the phenomenon of noise uncertainty is an important problem in model order selection when the noise variance is inaccurately known a priori [21]. Hence, both the signal correlation matrix and the noise variance are incorporated as nuisance parameters. This naturally leads to the formulation of the generalized likelihood ratio GLR statistic [22], which has been recently reported as the sphericity test [23] for the valuable scenario of low-rank Gaussian signals. Because the GLR asymptotically follows a non-central Chi-squared distribution, the statistical characterization of the proposed sphericity MDL is further addressed in this paper together with the MDL statistic with *priorly* known noise variance. Explicit expressions of both the PDFs of the MDLs and the associated error probabilities are obtained as a function of the main detection parameters, i.e., the signal-to-noise ratio (SNR), the observation size, and the dimensions of the signal. A simple approach based on the Jensen's approximation of the expectation of the MDLs is proposed to obtain the non-centrality parameters of the Chi-squared distributions. Numerical results show the accuracy of the proposed statistical characterization in predicting both the PDFs and the error probabilities, and show that the proposed sphericity MDL is robust to noise uncertainty.

In a cognitive radio scenario, the problem of determining the dimension of the primary signal subspace spanned by the eigenvectors of the correlation matrix is cast as follows. Each of the secondary users acquire a data set of $LM$ samples in the silent periods in which the secondary system is inactive. These observations are vectorized in $L$ vector measurements $\mathbf{x}_l$ of dimension $M$ which follow the discrete-time model

$$\mathbf{x}_l = \mathbf{s}_l + \mathbf{w}_l, \tag{1}$$

for $1 \leq l \leq L$, where each realization $\mathbf{s}_l$ and $\mathbf{w}_l$ represent the primary signal and the noise components at secondary user with power $P$ and variance $\sigma^2$, respectively. The average signal-to-noise ratio (SNR) in this problem is defined as

$$\text{SNR} \doteq \frac{NP}{\sigma^2}, \tag{2}$$

where $N$ is the dimension of the primary signal, and $P$ is the power per degree of freedom, i.e., the power associated to each of the $N$ dimensions of the primary signal. For Gaussian signals[1], the primary signal and noise components are distributed as

$$\mathbf{s}_l \sim \mathcal{CN}(\mathbf{0}, \mathbf{R}_s) \tag{3}$$
$$\mathbf{w}_l \sim \mathcal{CN}(\mathbf{0}, \sigma^2\mathbf{I}), \tag{4}$$

---

[1]The Gaussian assumption on both the primary signal and the noise is adopted in this paper as a worst-case scenario in the model order selection problem.

for all $l$, where $\mathbf{R}_s$ represents the correlation matrix of the primary signal. The objective of the model order selection problem is to determine the degrees of freedom of $\mathbf{R}_s$, defined as

$$N \doteq \text{rank}(\mathbf{R}_s), \tag{5}$$

with $N \leq M$, from the dataset $\mathbf{x}_1, \ldots, \mathbf{x}_L$ with unknown signal autocorrelation matrix $\mathbf{R}_s$. This paper presents a generic formulation of $\mathbf{R}_s$ without any structure or constraint rather than (5). Therefore, the results harvested in this paper are valid for arbitrary correlation domains, e.g., time, or space[2].

## III. SPHERICITY MINIMUM DESCRIPTION LENGTH WITH UNKNOWN NOISE VARIANCE

Typically, the primary signal subspace estimation problem involves low SNR regimes, as the secondary system may be far away from the primary system. In the assumption of white noise with unknown noise variance, the primary signal subspace estimation problem requires estimating the following parameters: the noise variance, the dimension of the primary signal subspace, and the eigenvectors associated to the primary signal subspace.

Because the estimation of the dimension of the primary signal subspace is a model order selection problem, the formulation of the following MDL statistic is required [4]

$$\text{MDL}(n) = -\log \text{GLR}(n) + nM\log(L), \tag{6}$$

where $\text{GLR}(n)$ is the GLR which jointly estimates the unknown parameters, i.e., the noise variance and the eigenvectors associated to the primary signal space. That is,

$$\text{GLR}(n) = \frac{\max_{\sigma^2, \mathbf{R}_s} p(\mathbf{x}_1, \ldots, \mathbf{x}_L | \mathcal{H}_n)}{\max_{\sigma^2} p(\mathbf{x}_1, \ldots, \mathbf{x}_L | \mathcal{H}_0)}, \tag{7}$$

where $\mathcal{H}_n$ and $\mathcal{H}_0$ denote the hypotheses that the primary signal has dimension $n$ or zero, respectively, i.e.,

$$\mathcal{H}_n : \text{rank}(\mathbf{R}_s) = n, \tag{8}$$

for $n = 0, 1, \ldots, M$. Solving (7) for $\mathbf{R}_s$ and $\sigma^2$ for a given primary signal subspace size $n$ derives to the rank-$n$ sphericity test [23]. As a result, for $n < M-1$, the MDL for the primary signal subspace estimation problem is given by (9) on the top of the next page, where $\lambda_1, \ldots, \lambda_M$ are the eigenvalues of the sample covariance matrix

$$\hat{\mathbf{R}} = \frac{1}{L}\sum_{l=1}^{L} \mathbf{x}_l \mathbf{x}_l^H. \tag{10}$$

As noted in [23], for $n \geq M-1$ the low rank structure of the signal correlation matrix cannot be exploited, and hence (9) particularizes to the AGM detector [24]. It is important to note that the traditional MDL in array processing (see, e.g., [12, Eq. (14)]) omits the first term of (9), as it is independent of the optimization variable $n$. In this work, however, the

---

[2]It is worth noting that in the case of temporal domain, $\mathbf{R}_s$ has a Toeplitz structure which could be exploited in the formulation of the MDL to further improve performance. However, the MLE of a Toeplitz matrix has not known closed-form solution.

$$\mathrm{MDL}(n) = \begin{cases} LM \log \dfrac{\prod_{m=1}^{M} \lambda_m^{1/M}}{\frac{1}{M} \sum_{m=1}^{M} \lambda_m} - L(M-n) \log \dfrac{\prod_{m=n+1}^{M} \lambda_m^{1/(M-n)}}{\frac{1}{M-n} \sum_{m=n+1}^{M} \lambda_m} + nM \log(L) & n < M-1 \\[4mm] LM \log \dfrac{\prod_{m=1}^{M} \lambda_m^{1/M}}{\frac{1}{M} \sum_{m=1,}^{M} \lambda_m} + nM \log(L) & n \geq M-1 \end{cases} \quad (9)$$

---

likelihood ratio structure of the MDL is preserved in order to take advantage of the asymptotic results of GLR detectors.

From (9), the estimation of the primary signal subspace dimension involves the minimization

$$\hat{N} = \arg\min_{n} \mathrm{MDL}(n), \quad (11)$$

where the search is performed in the set $n = 0, \ldots, M$.

### A. Statistical Characterization

The statistical characterization of the primary signal subspace estimation problem resorts to identify the statistical properties of the MDL function in (11). Because the MDL$(n)$ in (6) implements a GLR statistic for $n > 0$, its distribution is asymptotically (as $L \to \infty$) given by a non-central Chi-squared distribution [4]. More specifically, as

$$2 \log \mathrm{GLR}(n) \sim \mathcal{X}_{r_n}^2(\mu_n), \quad (12)$$

it follows from (6) that

$$2 \cdot [nM \log(L) - \mathrm{MDL}(n)] \sim \mathcal{X}_{r_n}^2(\mu_n), \quad (13)$$

where $r_n$ are the degrees of freedom and $\mu_n$ is the non-centrality parameter. In the following, the Chi-squared parameters are obtained for arbitrary $n$, which will be further use to derive the error probability.

Because the noise variance is a nuisance parameter and MDL$(n)$ involves the estimation of an Toeplitz Hermitian complex matrix of rank $n$, in Appendix A it is proved that the degrees of freedom $r_n$ are given as

$$r_n = 2Mn - n^2 + M - n - 1 \quad (14)$$

On the other hand, the non-centrality parameter $\mu_n$ depends on the occupancy of the primary signal, i.e., the true hypotheses $\mathcal{H}_N$. If the primary signal is not present (i.e., if $\mathcal{H}_0$ is true), the distribution is non-central and $\mu_n = 0$ for all $1 \leq n \leq M$. This task under $\mathcal{H}_N$ is a more difficult problem, as it involves the computation of the non-centrality parameter $\mu_n$. Making use of the Jensen's approximation, a good approximation holds for $n \geq N$:

$$\frac{\mu_n + r_n}{2} \approx L \log \frac{\left(1 + \dfrac{\mathrm{SNR}}{M}\right)^M}{\left(1 + \dfrac{\mathrm{SNR}}{N}\right)^N}. \quad (15)$$

It is worth noting that for $n \geq N$, i.e., when the MDL$(n)$ statistic has reached the true primary signal dimension $N$, the non-centrality parameter does not improve with $n$. This is due to the fact that the second part of the MDL$(n)$ in (9)

uncovers no structure in the noise subspace. Even tough, as the $nM \log(L)$ penalty term in (6) continues to increase with $n$, the MDL will "prefer" $n = N$ in front of higher dimensions. Finally, for $n < N$, the non-centrality parameter is affected by a negative term, i.e.,

$$\frac{\mu_n + r_n}{2} \approx L \log \frac{\left(1 + \dfrac{\mathrm{SNR}}{M}\right)^M}{\left(1 + \dfrac{\mathrm{SNR}}{N}\right)^N}$$
$$- L \log \frac{\left(1 + \dfrac{N-n}{M-n}\dfrac{\mathrm{SNR}}{N}\right)^{M-n}}{\left(1 + \dfrac{\mathrm{SNR}}{N}\right)^{N-n}}, \quad (16)$$

because the second part of the MDL$(n)$ in (9) now evaluates a portion of the signal subspace of dimension $N - n$. A last note on MDL is that by construction, MDL$(0) = 0$ in a deterministic fashion. A sketch of the proof of the non-centrality parameters (15) and (16) is provided in Appendix B. Both the expressions involved in The expressions involved in (15) and (16) have the structure of a *Shannon opportunity measure*, i.e., the difference between capacity terms that quantify the amount of information contained at each subspace. As an example, (15) can be rewritten as

$$\frac{\mu_n + r_n}{2} \approx ML \log\left(1 + \frac{\mathrm{SNR}}{M}\right) - NL \log\left(1 + \frac{\mathrm{SNR}}{N}\right) \quad (17)$$

with $N < M$.

### B. Error Probability

An important application of the statistical characterization of the primary signal subspace estimation problem is the computation of the error probability. In particular, the statistical characterization of the sphericity MDL allows to determine the probability of selecting $n$ as the dimension of the primary signal subspace is employed, i.e.,

$$P_n \doteq \mathbb{P}\left[\arg\min_{m} \mathrm{MDL}(m) = n\right], \quad (18)$$

for $n = 0, \ldots, M$. In Appendix C it is shown that this probability is given by

$$P_n = \prod_{m=0}^{n-1} \mathrm{CDF}_{-\mathcal{X}_{r_n - r_m}^2(\mu_n - \mu_m)}\left(\frac{(m-n)M \log(L)}{2}\right)$$
$$\times \prod_{m=n+1}^{M} \mathrm{CDF}_{\mathcal{X}_{r_m - r_n}^2(\mu_m - \mu_n)}\left(\frac{(m-n)M \log(L)}{2}\right), \quad (19)$$
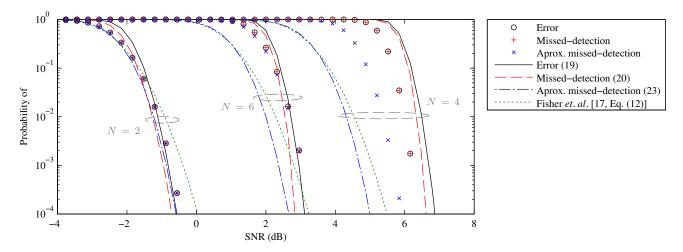
Fig. 1.  Probabilities of error, missed-detection and approximated missed-detection versus SNR for $M = 8$, $N = 2, 4, 6$ and $L = 1,000$.

where $r_n$ and $\mu_n$ are given in (14) and (15)–(16), respectively, under either $\mathcal{H}_0$ or $\mathcal{H}_N$. The error probability is defined as the complementary probability to the probability of detection, i.e,

$$P_e \doteq 1 - P_N, \qquad (20)$$

being $P_N$ the probability of detection given by (19) for $n = N$ under $\mathcal{H}_N$. As it is appreciated from (19), the pairs $(n, m)$ that exhibit smaller $\mu_n - \mu_m$ will contribute to the error probability.

It is worth noting that in cognitive radio an incorrect primary signal subspace detection will produce a different effect, depending if either the estimated dimension is smaller or larger than the true dimension. If $\hat{N} < N$, the secondary user will underestimate the primary signal subspace, hence will cause interference, whereas if $\hat{N} > N$, the secondary user will overestimate the primary signal subspace, hence losing opportunity. From the formulation above, these underestimation and overestimation probabilities, also denoted as missed-detection and false-alarm probabilities respectively, are given by

$$P_{\mathrm{MD}} = \sum_{n=1}^{N-1} P_n, \qquad (21)$$

and

$$P_{\mathrm{FA}} = \sum_{n=N+1}^{M} P_n. \qquad (22)$$

Furthermore, it is also worth mentioning that in the model order selection literature (see, e.g., [15]–[17]) the following approximation for the error probability is commonly adopted

$$P_{\mathrm{MD}} \approx \mathbb{P}\left[\mathrm{MDL}(N-1) - \mathrm{MDL}(N) < 0\right], \qquad (23)$$

due to the fact that the MDL tends to underestimate the model order, and hence the error probability (19) becomes dominated by the term $m = N - 1$ in the product. Using this result, (20) can be approximated as

$$P_{\mathrm{MD}} \approx 1 - \mathrm{CDF}_{-\mathcal{X}^2_{2K}(\mu_N - \mu_{N-1})}\left(-\frac{M \log L}{2}\right), \qquad (24)$$

where the difference of means equals

$$\mu_N - \mu_{N-1} = 2L \log \frac{\left(1 + \dfrac{1}{1+K}\dfrac{\mathrm{SNR}}{N}\right)^{1+K}}{\left(1 + \dfrac{\mathrm{SNR}}{N}\right)} + 2K, \qquad (25)$$

where the noise subspace dimension has been defined in the former expression as $K \doteq M - N$.

## IV. NUMERICAL RESULTS

This Section provides numerical results to show the behavior of the sphericity MDL statistic, and to assess the theoretical characterizations proposed in this work. The theoretical characterization addressed by [17, Eq. (12)] has been included as benchmarking, i.e.,

$$P_{\mathrm{MD}} \approx 1 - Q\left(\frac{\mu}{\sigma}\right), \qquad (26)$$

where $\mu$ and $\sigma$ are the asymptotically mean and variance of the asymptotically (as $L \to \infty$) normally distributed random variable $\mathrm{MDL}(N-1) - \mathrm{MDL}(N)$.

The probability of error, the missed-detection probability and the approximated missed-detection probability are depicted in Figure 1 versus SNR for several values of the signal order $N$ and a fixed observation size of $L = 1,000$. Lines illustrate the theoretical evaluations, and markers correspond to the average of 100,000 simulations. On the other hand, Figure 2 shows the same probabilities versus the number of observations $L$ for a fixed SNR of 2.5 dB.

From both figures, it can be concluded that

1) As known in the literature, the missed detection is the main source of error. The simulation results (markers) in both figures support this claim. On the other hand, the theoretical curves corresponding to the error probability and the missed-detection probability are slightly apart.

2) The behavior of the theoretical characterizations behave distinctly regarding the signal order $N$. For small $N$,
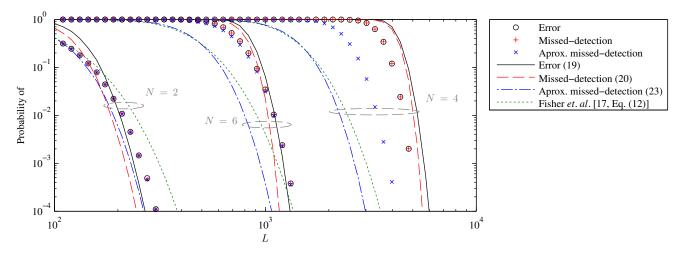
4

Fig. 2. Probabilities of error, missed-detection and approximated missed-detection versus $L$ for $M = 8$, $N = 2, 4, 6$ and SNR= 2.5 dB.

the approximated missed-detection and the true missed-detection probabilities are equivalent. However, as $N$ becomes large, the cross-terms $\text{MDL}_m - \text{MDL}_n$, for $m < n$, cannot be neglected. Therefore, the approximation (22) is not valid.

3) The present work improves the existing work by (i) considering the complete error probability, (ii) providing a more accurate characterization, and (iii) allowing to obtain insights of the problem by inspecting the means of the MDL statsitics, e.g., (14)–(16).

4) The slopes (error exponents) of the proposed theoretical characterization exhibit a more accurate matching with that of the empirical simulations, when compared to the Gaussian assumption adopted by [15]–[17].

## V. CONCLUSIONS

In this paper, we have provide the statistical characterization of the sphericity minimum description length (MDL) for the model order selection problem under unknown noise variance. By exploiting the asymptotically Chi-squared distribution of the generalized likelihood ratio (GLR) statistic, a complete characterization of the error probability for Gaussian sources is addressed.

## REFERENCES

[1] H. Akaike, "A new look at the statistical model identification," *IEEE Trans. Autom. Control*, vol. 19, no. 6, pp. 716–723, Dec. 1974.

[2] P. Stoica, P. Eykhoff, P. Janssen, and T. Söderström, "Model structure selection by cross-validation," *Int. J. Control*, vol. 43, pp. 1841–1878, 1986.

[3] H. Linhart and W. Zucchini, *Model Selection*. New York: Wiley, 1986.

[4] S. M. Kay, *Fundamentals of statistical signal processing*. Upper Saddle River, NJ: Prentice Hall, 1998, vol. 2 (detection theory).

[5] J. Rissanen, "Modeling by shortest data description," *Automatica*, vol. 14, no. 5, pp. 465–471, 1978.

[6] G. Schwarz, "Estimating the dimension of a model," *The Annals of Statistics*, vol. 6, no. 2, pp. 461–464, 2008.

[7] P. Stoica and Y. Selen, "Model-order selection: a review of information criterion rules," *IEEE Signal Process. Mag.*, vol. 21, no. 4, pp. 36–47, 2004.

[8] Q. Ding and S. Kay, "Inconsistency of the MDL: On the performance of model order selection criteria with increasing signal-to-noise ratio," *IEEE Trans. Signal Process.*, vol. 59, no. 5, pp. 1959–1969, May 2011.

[9] D. Schmidt and E. Makalic, "The consistency of MDL for linear regression models with increasing signal-to-noise ratio," *IEEE Trans. Signal Process.*, vol. 60, no. 3, pp. 1508–1510, Mar. 2012.

[10] P. Stoica and P. Babu, "On the exponentially embedded family (EEF) rule for model order selection," *IEEE Signal Process. Lett.*, vol. 19, no. 9, pp. 551–554, Sep. 2012.

[11] K. Beaudet and D. Cochran, "Estimation of subspace occupancy," in *Asilomar Conf. on Signals, Systems, and Computers*, Nov. 2014, pp. 1–5.

[12] M. Wax and T. Kailath, "Detection of signals by information theoretic criteria," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 33, no. 2, pp. 387–392, Apr. 1985.

[13] P. Djuric and S. Kay, "Order selection of autoregressive models," *IEEE Trans. Signal Process.*, vol. 40, no. 11, pp. 2829–2833, Nov. 1992.

[14] P. Djuric, "A model selection rule for sinusoids in white gaussian noise," *IEEE Trans. Signal Process.*, vol. 44, no. 7, pp. 1744–1751, Jul. 1996.

[15] Q.-T. Zhang, K. M. Wong, P. Yip, and J. Reilly, "Statistical analysis of the performance of information theoretic criteria in the detection of the number of signals in array processing," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 37, no. 10, pp. 1557–1567, Oct 1989.

[16] W. Xu and M. Kaveh, "Analysis of the performance and sensitivity of eigendecomposition-based detectors," *IEEE Trans. Signal Process.*, vol. 43, no. 6, pp. 1413–1426, Jun. 1995.

[17] E. Fishler, M. Grosmann, and H. Messer, "Detection of signals by information theoretic criteria: general asymptotic performance analysis," *IEEE Trans. Signal Process.*, vol. 50, no. 5, pp. 1027–1036, May 2002.

[18] Z. Lu and A. Zoubir, "Generalized bayesian information criterion for source enumeration in array processing," *IEEE Trans. Signal Process.*, vol. 61, no. 6, pp. 1470–1480, Mar. 2013.

[19] S. Haykin, D. J. Thomson, and J. H. Reed, "Spectrum sensing for cognitive radio," *Proc. IEEE*, vol. 97, no. 5, pp. 849–877, May 2009.

[20] A. Goldsmith, S. A. Jafar, I. Maric, and S. Srinivasa, "Breaking spectrum gridlock with cognitive radios: An information theoretic perspective," *Proc. IEEE*, vol. 97, no. 5, pp. 894–914, May 2009.

[21] R. Tandra and A. Sahai, "SNR walls for signal detection," *IEEE J. Sel. Topics Signal Process.*, vol. 2, no. 1, pp. 4–17, Feb. 2008.

[22] P. Stoica, Y. Selen, and J. Li, "On information criteria and the generalized likelihood ratio test of model order selection," *IEEE Signal Process. Lett.*, vol. 11, no. 10, pp. 794–797, Oct. 2004.

[23] D. Ramírez, G. Vázquez-Vilar, R. López-Valcarce, J. Vía, and I. Santamaría, "Detection of rank-$P$ signals in cognitive radio networks with uncalibrated multiple antennas," *IEEE Trans. Signal Process.*, vol. 59, no. 8, pp. 3764–3774, Aug. 2011.

[24] J. Mauchly, "Significance test for sphericity of a normal $n$-variate distribution," *Ann. Math. Stat.*, vol. 11, pp. 204–209, 1940.

5