

Consistency of Model Order Selection Rules

Niklas Wahlström

May 20, 2011

These notes present techniques to be used for analyzing the consistency of model order selection rules. The discussion is focused on the array model and the Bayesian Information Criteria, as it is presented in [4]. However, the analysis can be extended to more general models as well as other model order selection rules, which will be discussed in the lecture.

Furthermore, in the Appendix, the asymptotic distribution of a likelihood ratio is presented and proved, inspired by the treatment in [3].

1 Problem formulation

In array processing, observations are described by the model

$$\mathbf{x}(t) = \sum_{i=1}^q \mathbf{a}(\boldsymbol{\theta}_i) s_i(t) + \mathbf{n}(t) \quad (1)$$

where $s_i(t) \in \mathbb{C}$ is j th emitters signal, $\mathbf{a}(\boldsymbol{\theta}_i) \in \mathbb{C}^p$ the sensor response and $\mathbf{n}(t) \in \mathbb{C}^p$ additive noise. In matrix form the following familiar equation is obtained

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) + \mathbf{n}(t) \quad (2)$$

where

$$\mathbf{A} = [\mathbf{a}(\boldsymbol{\theta}_1), \dots, \mathbf{a}(\boldsymbol{\theta}_q)]$$

. We assume that

- $q < p$, i.e. that there are more sensors than emitted signals
- $\mathbf{s}(t)$ and $\mathbf{n}(t)$ are stationary, white, zero mean, circular symmetric Gaussian random processes. Especially we assume that $\mathbf{S} = \mathbb{E}[\mathbf{s}(t)\mathbf{s}^*(t)]$ is positive definite and that $\mathbb{E}[\mathbf{n}(t)\mathbf{n}^*(t)] = \sigma^2\mathbf{I}$.
- \mathbf{A} has full rank.

The problem is to detect the number of signals q from the observations $\mathbf{x}(t_1), \dots, \mathbf{x}(t_N)$.

From the model description we can compute the covariance matrix of $\mathbf{x}(t)$ as

$$\mathbf{R} = \boldsymbol{\Psi} + \sigma^2\mathbf{I} \quad (3)$$

where

$$\mathbf{\Psi} = \mathbf{A}\mathbf{S}\mathbf{A}^* \quad (4)$$

Since \mathbf{S} is a positive definite symmetric matrix we have

$$\text{rank}(\mathbf{\Psi}) = \text{rank}(\mathbf{A}\mathbf{S}\mathbf{A}^*) = \text{rank}(\mathbf{A}\sqrt{\mathbf{S}}(\mathbf{A}\sqrt{\mathbf{S}})^*) \quad (5)$$

$$= \text{rank}(\mathbf{A}\sqrt{\mathbf{S}}) = \text{rank}(\mathbf{A}) = \min(p, q) = q \quad (6)$$

Thus, $\mathbf{\Psi}$ has exactly q non-zero eigenvalues. This gives the following eigenvalues of \mathbf{R}

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_q > \underbrace{\lambda_{q+1} = \dots = \lambda_p}_{=\sigma^2} \quad (7)$$

The number of signals q can hence be determined from the multiplicity of the smallest eigenvalue of \mathbf{R} .

2 Estimation

From previous section we know that the observations $\mathbf{x}(t_1), \dots, \mathbf{x}(t_N)$ are independent Gaussian random vectors with the covariance matrix (3). For estimation, we therefore consider

$$\mathbf{R}^{(k)} = \mathbf{\Psi}^{(k)} + \sigma^2 \mathbf{I} \quad (8)$$

where $\mathbf{\Psi}^{(k)}$ is a positive semi-definite matrix of rank k . With a spectral representation of $\mathbf{R}^{(k)}$ we get

$$\mathbf{R}^{(k)} = \sum_{i=1}^k (\lambda_i - \sigma^2) \mathbf{V}_i \mathbf{V}_i^* + \sigma^2 \mathbf{I} \quad (9)$$

and $\mathbf{R}^{(k)}$ can thus be parametrized with

$$\mathbf{\Theta}^{(k)} = (\lambda_1, \dots, \lambda_k, \sigma^2, \mathbf{V}_1^T, \dots, \mathbf{V}_k^T)^T \quad (10)$$

Now, we have

$$f(\mathbf{x}(t_1), \dots, \mathbf{x}(t_N) | \mathbf{\Theta}^{(k)}) = \prod_{i=1}^N \frac{1}{\pi^p |\mathbf{R}^{(k)}|} \exp(-\mathbf{x}(t_i)^* (\mathbf{R}^{(k)})^{-1} \mathbf{x}(t_i)) \quad (11)$$

and the log-likelihood function then becomes

$$\begin{aligned} \log l(\mathbf{\Theta}^{(k)}) &= \sum_{i=1}^N (-\log(|\mathbf{R}^{(k)}|) - \mathbf{x}(t_i)^* (\mathbf{R}^{(k)})^{-1} \mathbf{x}(t_i)) + \text{const.} \\ &= -N \log(|\mathbf{R}^{(k)}|) - \sum_{i=1}^N \text{tr} \left(\mathbf{x}(t_i)^* (\mathbf{R}^{(k)})^{-1} \mathbf{x}(t_i) \right) + \text{const.} \\ &= -N \log(|\mathbf{R}^{(k)}|) - \text{tr} \left((\mathbf{R}^{(k)})^{-1} \sum_{i=1}^N \mathbf{x}(t_i) \mathbf{x}(t_i)^* \right) + \text{const.} \\ &= -N \log |\mathbf{R}^{(k)}| - N \text{tr}((\mathbf{R}^{(k)})^{-1} \hat{\mathbf{R}}) + \text{const.} \end{aligned} \quad (12)$$

where $\hat{\mathbf{R}}$ is the sample-covariance matrix

$$\hat{\mathbf{R}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}(t_i) \mathbf{x}(t_i)^* \quad (13)$$

It can be shown that the MLE of $\Theta^{(k)}$ is

$$\hat{\lambda}_i = l_i \quad i = 1, \dots, k \quad (14a)$$

$$\hat{\sigma}^2 = \frac{1}{p-k} \sum_{i=k+1}^p l_i \quad (14b)$$

$$\hat{\mathbf{V}}_i = \mathbf{C}_i \quad i = 1, \dots, k \quad (14c)$$

where $l_1 > l_2 \dots > l_p$ and $\mathbf{C}_1 \dots \mathbf{C}_p$ are eigenvalues and eigenvectors of $\hat{\mathbf{R}}$. Now, from the MLE in (14) we can construct

$$\hat{\mathbf{R}}^{(k)} = \sum_{i=1}^k (\hat{\lambda}_i - \hat{\sigma}^2) \hat{\mathbf{V}}_i \hat{\mathbf{V}}_i^* + \hat{\sigma}^2 \mathbf{I} \quad (15)$$

By substituting (14) and (15) into (12) then gives

$$\log l(\hat{\Theta}^{(k)}) = -N \log |\hat{\mathbf{R}}^{(k)}| - N \text{tr}((\hat{\mathbf{R}}^{(k)})^{-1} \hat{\mathbf{R}}) + \text{const.} \quad (16)$$

Since the eigenvalues of $\hat{\mathbf{R}}^{(k)}$ are $l_1, \dots, l_k, \hat{\sigma}^2, \dots, \hat{\sigma}^2$, the first term in (16) will be

$$\begin{aligned} \log |\hat{\mathbf{R}}^{(k)}| &= \log |\hat{\mathbf{R}}^{(k)}| \\ &= \log \left(\prod_{i=1}^k l_i \prod_{i=k+1}^p \hat{\sigma}^2 \right) \\ &= \log \left(\prod_{i=1}^k l_i \right) + (p-k) \log(\hat{\sigma}^2) \\ &= -\log \left(\prod_{i=k+1}^p l_i \right) + \log \left(\prod_{i=1}^p l_i \right) + (p-k) \log(\hat{\sigma}^2) \\ &= -\log \left(\frac{\sqrt[p-k]{\prod_{i=k+1}^p l_i}}{\frac{1}{p-k} \sum_{i=k+1}^p l_i} \right)^{(p-k)} + \log \left(\prod_{i=1}^p l_i \right) \end{aligned} \quad (17)$$

For the second term in (16) we have to compute the eigenvalues of $(\hat{\mathbf{R}}^{(k)})^{-1} \hat{\mathbf{R}}$. By construction, the eigenvectors of $\hat{\mathbf{R}}$ are also eigenvectors to $\hat{\mathbf{R}}^{(k)}$ as well as $(\hat{\mathbf{R}}^{(k)})^{-1}$ and the eigenvalues of $(\hat{\mathbf{R}}^{(k)})^{-1} \hat{\mathbf{R}}$ can therefore be computed by element-wise multiplication of the eigenvalues of the two matrices. This gives

$$\begin{aligned} \text{tr}((\hat{\mathbf{R}}^{(k)})^{-1} \hat{\mathbf{R}}) &= \left(\sum_{i=1}^k \frac{l_i}{l_i} + \sum_{k+1}^p \frac{l_i}{\hat{\sigma}^2} \right) \\ &= \left(k + \frac{1}{\frac{1}{p-k} \sum_{i=k+1}^p l_i} \sum_{i=k+1}^p l_i \right) = p \end{aligned} \quad (18)$$

which gives

$$\begin{aligned} \log l(\hat{\Theta}^{(k)}) &= N \log \left(\frac{\sqrt[p-k]{\prod_{i=k+1}^p l_i}}{\frac{1}{p-k} \sum_{i=k+1}^p l_i} \right)^{(p-k)} - N \log \left(\prod_{i=1}^p l_i \right) - Np + \text{const.} \\ &= N \log \left(\frac{\sqrt[p-k]{\prod_{i=k+1}^p l_i}}{\frac{1}{p-k} \sum_{i=k+1}^p l_i} \right)^{(p-k)} + \text{const.} \end{aligned} \quad (19)$$

Now, the Bayesian Information Criterion Rule is given by minimizing the negative log-likelihood function with an extra penalty coefficient

$$\text{BIC}(k) = -2 \log l(\hat{\Theta}^{(k)}) + n_k \log(N) + \text{const.} \quad (20)$$

where n_k is the dimension of the state (10).

3 Consistency

For consistency we demand when $N \rightarrow \infty$ that

$$P(\arg \min \text{BIC}(k) = q) = 1 \implies \quad (21)$$

$$\begin{cases} P(\text{BIC}(k) > \text{BIC}(q) | k < q) = 1 & \text{no underfitting} \\ P(\text{BIC}(k) > \text{BIC}(q) | k > q) = 1 & \text{no overfitting} \end{cases} \quad (22)$$

Thus, we can prove consistency by proving that the quantity

$$\text{BIC}(q) - \text{BIC}(k) = 2 \log \frac{l(\hat{\Theta}^{(k)})}{l(\hat{\Theta}^{(q)})} + (n_q - n_k) \log N \quad (23)$$

is negative w.p.1 both in case of underfitting $k < q$ as well in case of overfitting $k > q$.

3.1 Underfitting

Consider the case $k < q$ and divide the quantity (23) with $2N$

$$\frac{\text{BIC}(q) - \text{BIC}(k)}{2N} = \frac{1}{N} \log \frac{l(\hat{\Theta}^{(k)})}{l(\hat{\Theta}^{(q)})} + (n_q - n_k) \frac{\log N}{2N} \quad (24)$$

and consider the first term of (24)

$$\begin{aligned}
\frac{1}{N} \log \frac{l(\hat{\Theta}^{(k)})}{l(\hat{\Theta}^{(q)})} &= \log \left(\frac{\prod_{i=k+1}^p l_i}{\left(\frac{1}{p-k} \sum_{i=k+1}^p l_i\right)^{p-k}} \cdot \frac{\left(\frac{1}{p-q} \sum_{i=q+1}^p l_i\right)^{p-q}}{\prod_{i=q+1}^p l_i} \right) \\
&= \log \left(\frac{\prod_{i=k+1}^q l_i}{1} \cdot \frac{\left(\frac{1}{p-q} \sum_{i=q+1}^p l_i\right)^{p-q}}{\left(\frac{1}{p-k} \sum_{i=k+1}^p l_i\right)^{p-k}} \right) \\
&= \log \left(\frac{\prod_{i=k+1}^q l_i}{\left(\frac{1}{q-k} \sum_{i=k+1}^q l_i\right)^{q-k}} \right) \\
&\quad + \log \left(\frac{\left(\frac{1}{p-q} \sum_{i=q+1}^p l_i\right)^{p-q} \left(\frac{1}{q-k} \sum_{i=k+1}^q l_i\right)^{q-k}}{\left(\frac{1}{q-k} \sum_{i=k+1}^q l_i\right)^{q-k}} \right) \\
&= \log \left(\frac{\sqrt[q-k]{\prod_{i=k+1}^q l_i}}{\frac{1}{q-k} \sum_{i=k+1}^q l_i} \right)^{q-k} \\
&\quad + \log \left(\frac{\left(\frac{1}{p-q} \sum_{i=q+1}^p l_i\right)^{\frac{p-q}{p-k}} \left(\frac{1}{q-k} \sum_{i=k+1}^q l_i\right)^{\frac{q-k}{p-k}}}{\frac{1}{p-k} \sum_{i=k+1}^p l_i} \right)^{p-k} \quad (25)
\end{aligned}$$

The first term of (25) is non-positive due to the AM-GM inequality

$$\frac{1}{q-k} \sum_{i=k+1}^q l_i \geq \sqrt[q-k]{\prod_{i=k+1}^q l_i} \quad (26)$$

Note that, in contrast to what is stated in [4], this inequality *can* be fulfilled with equality. This is the case when all $\{l_i\}_{k+1}^q$ are equal, which especially is the case when $k+1 = q$.

For the second term, the generalized AM-GM inequality can be applied

$$\sum_{i=1}^n w_i A_i \geq \prod A_i^{w_i} \quad \sum_{i=1}^n w_i = 1 \quad (27)$$

which is fulfilled with equality if and only if all A_i with $w_i > 0$ are equal. By identifying

$$A_1 = \frac{1}{p-q} \sum_{i=q+1}^p l_i, \quad A_2 = \frac{1}{q-k} \sum_{i=k+1}^q l_i, \quad w_1 = \frac{p-q}{p-k}, \quad w_2 = \frac{q-k}{p-k} \quad (28)$$

we get

$$\begin{aligned} & \overbrace{\left(\frac{1}{p-k} \sum_{i=k+1}^p l_i \right)}^{= \frac{1}{p-k} \sum_{i=k+1}^p l_i} \\ & \frac{p-q}{p-k} \left(\frac{1}{p-q} \sum_{i=q+1}^p l_i \right) + \frac{q-k}{p-k} \left(\frac{1}{q-k} \sum_{i=k+1}^q l_i \right) > \\ & \left(\frac{1}{p-q} \sum_{i=q+1}^p l_i \right)^{\frac{p-q}{p-k}} \cdot \left(\frac{1}{q-k} \sum_{i=k+1}^q l_i \right)^{\frac{q-k}{p-k}} \end{aligned} \quad (29)$$

This inequality *can not* be fulfilled with equality since

$$\frac{1}{p-q} \sum_{i=q+1}^p l_i - \frac{1}{q-k} \sum_{i=k+1}^q l_i \rightarrow \frac{1}{p-q} \sum_{i=q+1}^p \lambda_i - \frac{1}{q-k} \sum_{i=k+1}^q \lambda_i > 0$$

when $N \rightarrow \infty$. Thus, we have that the second term in (25) goes to a negative constant as $N \rightarrow \infty$. Since the last term of (24) goes to zero as $N \rightarrow \infty$, the whole difference (24) will be negative w.p.1 in the large-sample limit.

3.2 Overfitting

Now, we consider the case $k > q$. The first term of the right hand side in (23) can be considered as a *generalized likelihood ratio* (GLR). We can find the asymptotic distribution of this ratio by using Lemma A.2. First, we identify

$$\Theta^{(q)} = \boldsymbol{\theta}_s, \quad \Theta^{(k)} = \begin{bmatrix} \boldsymbol{\theta}_r \\ \boldsymbol{\theta}_s \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\theta}_{r_0} = \mathbf{0}. \quad (30)$$

Furthermore, since q is the true model order, the correct hypothesis is

$$\mathcal{H}_0 : \boldsymbol{\theta}_r = \boldsymbol{\theta}_{r_0}.$$

According to Lemma A.2 we then have

$$2 \log \frac{l(\hat{\Theta}^{(k)})}{l(\hat{\Theta}^{(q)})} \sim \chi_{n_k - n_q}^2 \quad (31)$$

This gives

$$P(\text{BIC}(q) - \text{BIC}(k) < 0) = P(\chi_{n_k - n_q}^2 < (n_k - n_q) \log N) \rightarrow 1 \quad \text{as} \quad N \rightarrow \infty \quad (32)$$

By this we have shown that both the probability of underfitting as well as the probability of overfitting goes to zero in the large sample limit and thus, BIC is consistent.

A Appendix

Lemma A.1. *Consider the symmetric matrix*

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \quad (33)$$

Then

$$\min_y \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^T \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{x}^T \Delta_{\mathbf{C}} \mathbf{x} \quad (34)$$

where $\Delta_{\mathbf{C}} = \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T$ is the Schur complement of \mathbf{C} in \mathbf{M}

Proof. By using the block triangular factorization (see, e.g. [1])

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{C}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \Delta_{\mathbf{C}} & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C}^{-1}\mathbf{B}^T & \mathbf{I} \end{bmatrix} \quad (35)$$

we get

$$\min_y \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^T \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \min_y \mathbf{x}^T \Delta_{\mathbf{C}} \mathbf{x} + \underbrace{(\mathbf{C}^{-1}\mathbf{B}^T \mathbf{x} - \mathbf{y})^T \mathbf{C} (\mathbf{C}^{-1}\mathbf{B}^T \mathbf{x} - \mathbf{y})}_{=0 \text{ when } \mathbf{y}=\mathbf{C}^{-1}\mathbf{B}^T \mathbf{x}} = \mathbf{x}^T \Delta_{\mathbf{C}} \mathbf{x} \quad (36)$$

□

Lemma A.2. Consider a PDF $p(\mathbf{x}; \boldsymbol{\theta})$ where $\boldsymbol{\theta}$ is a vector of unknown parameters. Assume that the MLE of $\boldsymbol{\theta}$ is both consistent and efficient. Partition $\boldsymbol{\theta}$ as

$$\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_r \\ \boldsymbol{\theta}_s \end{bmatrix} = \begin{bmatrix} r \times 1 \\ s \times 1 \end{bmatrix} \quad (37)$$

and consider the parameter test

$$\begin{aligned} \mathcal{H}_0 &: \boldsymbol{\theta}_r = \boldsymbol{\theta}_{r_0} \\ \mathcal{H}_1 &: \boldsymbol{\theta}_r \neq \boldsymbol{\theta}_{r_0} \end{aligned}$$

as well as the test statistic

$$L_G(\mathbf{x}) = \frac{p(\mathbf{x}; \hat{\boldsymbol{\theta}}_{r_1}, \hat{\boldsymbol{\theta}}_{s_1})}{p(\mathbf{x}; \boldsymbol{\theta}_{r_0}, \hat{\boldsymbol{\theta}}_{s_0})} = \frac{\max_{\boldsymbol{\theta}_r, \boldsymbol{\theta}_s} p(\mathbf{x}; \boldsymbol{\theta}_r, \boldsymbol{\theta}_s)}{\max_{\boldsymbol{\theta}_s} p(\mathbf{x}; \boldsymbol{\theta}_r = \boldsymbol{\theta}_{r_0}, \boldsymbol{\theta}_s)}. \quad (38)$$

Then, as $N \rightarrow \infty$ the statistic $\log L_G(\mathbf{x})$ has the the PDF

$$2 \log L_G(\mathbf{x}) = \begin{cases} \chi_r^2 & \text{under } \mathcal{H}_0 \\ \chi_r^2(\lambda) & \text{under } \mathcal{H}_1 \end{cases} \quad (39)$$

The noncentrality parameter is

$$\lambda = (\boldsymbol{\theta}_{r_1} - \boldsymbol{\theta}_{r_0})^T \Delta_{\mathcal{I}_{ss}} (\boldsymbol{\theta}_{r_1} - \boldsymbol{\theta}_{r_0}). \quad (40)$$

where $\Delta_{\mathcal{I}_{ss}} = \mathcal{I}_{rr} - \mathcal{I}_{rs} \mathcal{I}_{ss}^{-1} \mathcal{I}_{rs}^T$ is the Schur complement of \mathcal{I}_{ss} in the information matrix

$$\begin{bmatrix} \mathcal{I}_{rr} & \mathcal{I}_{rs} \\ \mathcal{I}_{rs}^T & \mathcal{I}_{ss} \end{bmatrix} = \mathcal{I} = -\mathbb{E} \left(\frac{\partial^2 \log p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right) \Bigg|_{\boldsymbol{\theta}=\boldsymbol{\theta}_1} \quad (41)$$

and where $\boldsymbol{\theta}_1 = [\boldsymbol{\theta}_{r_1}^T \boldsymbol{\theta}_{s_1}^T]^T$ is the true value of $\boldsymbol{\theta}$ under \mathcal{H}_1 .

Proof. Since the MLE of $\boldsymbol{\theta}$ is efficient, we know that [2]

$$\frac{\partial \log p(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathcal{I}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \quad (42)$$

where $\mathcal{I}(\boldsymbol{\theta})$ is the Fisher information matrix. Since MLE $\boldsymbol{\theta}$ is consistent we have

$$\mathcal{I}(\boldsymbol{\theta}) = \mathcal{I}(\hat{\boldsymbol{\theta}}) + \mathcal{O}(\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|) \implies \quad (43)$$

$$\mathcal{I}(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = (\mathcal{I}(\hat{\boldsymbol{\theta}}) + \mathcal{O}(\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|))(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \quad (44)$$

$$= \mathcal{I}(\hat{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) + \mathcal{O}(\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|^2) \rightarrow \mathcal{I}(\hat{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \quad (45)$$

as $N \rightarrow \infty$ since $\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\| \rightarrow 0$ due to consistency. Integrating and taking the exponential gives

$$p(\mathbf{x}; \boldsymbol{\theta}) = K \exp \left[-\frac{1}{2}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T \mathcal{I}(\hat{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \right] =: K \exp(\boldsymbol{\theta}; \hat{\boldsymbol{\theta}}, \mathcal{I}(\hat{\boldsymbol{\theta}})) \quad (46)$$

for some constant K . Note, that this can be seen as a Laplace approximation of the likelihood around its maximum $\hat{\boldsymbol{\theta}}$ and where the curvature is given by the information matrix. Furthermore, since $\mathcal{I}(\hat{\boldsymbol{\theta}}) \rightarrow \mathcal{I}(\boldsymbol{\theta})$ as $N \rightarrow \infty$, we will from here on skip the argument and simply referring to the information matrix as \mathcal{I} .

By using the partitioning (37) and (41), we can write

$$p(\mathbf{x}; \boldsymbol{\theta}_r, \boldsymbol{\theta}_s) = K \exp \left(\begin{bmatrix} \boldsymbol{\theta}_r \\ \boldsymbol{\theta}_s \end{bmatrix}; \begin{bmatrix} \hat{\boldsymbol{\theta}}_{r_1} \\ \hat{\boldsymbol{\theta}}_{s_1} \end{bmatrix}, \begin{bmatrix} \mathcal{I}_{rr} & \mathcal{I}_{rs} \\ \mathcal{I}_{rs}^T & \mathcal{I}_{ss} \end{bmatrix} \right) \quad (47)$$

By this we evaluate the nominator of (39) as

$$p(\mathbf{x}; \hat{\boldsymbol{\theta}}_{r_1}, \hat{\boldsymbol{\theta}}_{s_1}) = K \quad (48)$$

For the denominator, Lemma A.1, which gives us

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\theta}_{r_0}, \hat{\boldsymbol{\theta}}_{s_0}) &= \max_{\boldsymbol{\theta}_s} K \exp \left(\begin{bmatrix} \boldsymbol{\theta}_{r_0} \\ \boldsymbol{\theta}_s \end{bmatrix}; \begin{bmatrix} \hat{\boldsymbol{\theta}}_{r_1} \\ \hat{\boldsymbol{\theta}}_{s_1} \end{bmatrix}, \begin{bmatrix} \mathcal{I}_{rr} & \mathcal{I}_{rs} \\ \mathcal{I}_{rs}^T & \mathcal{I}_{ss} \end{bmatrix} \right) \\ &= K \exp(\boldsymbol{\theta}_{r_0}; \hat{\boldsymbol{\theta}}_{r_1}, \boldsymbol{\Delta}_{\mathcal{I}_{ss}}) \end{aligned} \quad (49)$$

where $\boldsymbol{\Delta}_{\mathcal{I}_{ss}} = \mathcal{I}_{rr} - \mathcal{I}_{rs} \mathcal{I}_{ss}^{-1} \mathcal{I}_{rs}^T$ is the Schur complement of \mathcal{I}_{ss} in \mathcal{I} . This leads to

$$\begin{aligned} 2 \log L_G(\mathbf{x}) &= 2 \log \frac{K}{K \exp(\boldsymbol{\theta}_{r_0}; \hat{\boldsymbol{\theta}}_{r_1}, \boldsymbol{\Delta}_{\mathcal{I}_{ss}})} \\ &= (\hat{\boldsymbol{\theta}}_{r_1} - \boldsymbol{\theta}_{r_0})^T \boldsymbol{\Delta}_{\mathcal{I}_{ss}} (\hat{\boldsymbol{\theta}}_{r_1} - \boldsymbol{\theta}_{r_0}). \end{aligned} \quad (50)$$

Since $\hat{\boldsymbol{\theta}}_{r_1}$ is an efficient and unrestricted MLE of $\boldsymbol{\theta}_r$, it will attain the CRLB which can be computed as the corresponding sub-matrix of the inverse of the information matrix $(\mathcal{I}^{-1})_{rr} = (\boldsymbol{\Delta}_{\mathcal{I}_{ss}})^{-1}$. Furthermore, due to consistency, $\hat{\boldsymbol{\theta}}_{r_1}$ will attain its true value in the large sample limit. Therefore, under \mathcal{H}_0 and as $N \rightarrow \infty$ we have

$$\hat{\boldsymbol{\theta}}_{r_1} \sim \mathcal{N}(\boldsymbol{\theta}_{r_0}, (\boldsymbol{\Delta}_{\mathcal{I}_{ss}})^{-1}) \quad (51)$$

By setting $\mathbf{z}_0 = \sqrt{\Delta_{\mathcal{I}_{ss}}}(\hat{\boldsymbol{\theta}}_{r_1} - \boldsymbol{\theta}_{r_0}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and we get under \mathcal{H}_0

$$2 \log L_G(\mathbf{x}) = \mathbf{z}_0^T \mathbf{z}_0 \sim \chi_r^2 \quad (52)$$

Under \mathcal{H}_1 we have

$$\hat{\boldsymbol{\theta}}_{r_1} \sim \mathcal{N}(\boldsymbol{\theta}_{r_1}, (\Delta_{\mathcal{I}_{ss}})^{-1}) \quad (53)$$

By setting $\mathbf{z}_1 = \sqrt{\Delta_{\mathcal{I}_{ss}}}(\hat{\boldsymbol{\theta}}_{r_1} - \boldsymbol{\theta}_{r_1}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\boldsymbol{\mu}_1 = \sqrt{\Delta_{\mathcal{I}_{ss}}}(\boldsymbol{\theta}_{r_1} - \boldsymbol{\theta}_{r_0})$ we get

$$2 \log L_G(\mathbf{x}) = (\mathbf{z}_1 + \boldsymbol{\mu}_1)^T (\mathbf{z}_1 + \boldsymbol{\mu}_1) \sim \chi_r^2(\lambda) \quad (54)$$

where

$$\lambda = \boldsymbol{\mu}_1^T \boldsymbol{\mu}_1 = (\boldsymbol{\theta}_{r_1} - \boldsymbol{\theta}_{r_0})^T \Delta_{\mathcal{I}_{ss}} (\boldsymbol{\theta}_{r_1} - \boldsymbol{\theta}_{r_0}). \quad (55)$$

□

References

- [1] T. Kailath, A. H. Sayed, and B. Hassibi. *Linear Estimation*. Prentice-Hall, Inc., Upper Saddle River, NJ, 2000.
- [2] S. M. Kay. *Fundamentals of Statistical Signal Processing: Volume 1: Estimation Theory*. Prentice-Hall, PTR, Upper Saddle River, NJ, USA, 2009.
- [3] S. M. Kay. *Fundamentals of Statistical Signal Processing: Volume 2: Detection Theory*. Prentice-Hall, PTR, Upper Saddle River, NJ, USA, 2009.
- [4] M. Wax and T. Kailath. Detection of signals by information theoretic criteria. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 33(2):387 – 392, apr 1985.