

Analyzing the parameter bias when an ARMAX model is fitted to noise-corrupted data

Torsten Söderström and Umberto Soverini

October 10, 2022

Abstract

When an ARMAX model is fitted to noise-corrupted data using the prediction error method, biased estimates are obtained. The bias is examined, with emphasis on the situation when the system is almost non-identifiable. In contrast to the case of using an output error model, no general results on the size of the bias seem to apply.

1 Introduction

The aim of these notes is to analyze the asymptotic bias when an ARMAX model is used and the measured input data contains measurement noise. In particular it is assumed that the system is nearly unidentifiable, and the effect of the pole-zero separation on the bias is of particular interest.

Some comparisons are made to the case when instead an output error model is fitted to the data.

The model structure considered is

$$A(q^{-1})y(t) = B(q^{-1})u(t) + C(q^{-1})\varepsilon(t) , \quad (1)$$

and the parameter vector to be estimated is assumed to be

$$\theta = (a_1 \quad \dots \quad a_{n_a} \quad b_1 \quad \dots \quad b_{n_b} \quad c_1 \quad \dots \quad c_{n_c})^T . \quad (2)$$

Let θ_0 denote the true value of the parameter vector.

The input contains a noise-free part, and an additional measurement noise, as

$$u(t) = u_0(t) + \tilde{u}(t) . \quad (3)$$

The underlying data is assumed to fulfil

$$A(q^{-1})y(t) = B(q^{-1})u_0(t) + C(q^{-1})e(t) . \quad (4)$$

It is further assumed:

- A prediction error method is used for fitting the model to the data.
- The asymptotic case when the number of data tends to infinity is considered.
- The input noise $\tilde{u}(t)$ is assumed to be white, with zero mean and variance λ_u^2 .

The general problem was analyzed in [1], where it was shown that the parameter bias $\tilde{\theta} = \hat{\theta} - \theta_0$ can be approximated by

$$\tilde{\theta} = - (V_{\theta\theta})^{-1} V_{\theta} , \quad (5)$$

where $V_{\theta\theta}$ and V_{θ} are the Hessian and the gradient, respectively, of the asymptotic prediction error criterion (for an infinite amount of data)

$$V(\theta) = E \{ \varepsilon^2(t, \theta) \} , \quad (6)$$

and $\varepsilon(t, \theta)$ is the prediction error. The expression (5) is particularly accurate for small noise levels, i.e. when λ_u^2 is small.

In [1], [2] it was shown, with reference to an output error model structure, that the bias vector $\tilde{\theta}$ in (5) can be large when the system is almost non-identifiable due to almost pole-zero cancellation.

2 Expressing the almost non-identifiability

For the ARMAX model (1), the system is almost non-identifiable when *all the three polynomials* A , B , C have almost a common zero. Here we formalize this situation by assuming that

$$A = A_1 \bar{A}, \quad B = B_1 \bar{B}, \quad C = C_1 \bar{C} , \quad (7)$$

$$\bar{A} = 1 + \bar{a}q^{-1}, \quad \bar{B} = 1 + \bar{b}q^{-1}, \quad \bar{C} = 1 + \bar{c}q^{-1} , \quad (8)$$

$$\delta \triangleq \max (|\bar{a} - \bar{b}|, |\bar{a} - \bar{c}|, |\bar{b} - \bar{c}|) \text{ is small.} \quad (9)$$

In what follows we will examine the (approximate) bias when the separation δ is small.

When an output error method is used, only two polynomials are involved. It was shown earlier in [1], [2] that for such a case, the bias will be of order $O(1/\delta)$.

In the ARMAX case we first develop an expression for the gradient $\psi(t) = \frac{\partial \varepsilon}{\partial \theta}(t)$, which will be splitted into two parts. The first part, denoted $\psi_1(t)$ will correspond to the case when $\bar{A} = \bar{B} = \bar{C}$, thus that is when there is a common zero to the three polynomials A , B and C .

In the expressions below let the running indeces in general be $i = 1, \dots, n_a$, $j = 1, \dots, n_b$, $k = 1, \dots, n_c$. The gradient can be written as

$$\begin{aligned}
\psi^T(t) &= \left(\frac{1}{C}y(t-i) \quad -\frac{1}{C}u(t-j) \quad -\frac{1}{C}\varepsilon(t-k) \right) \\
&= \left(\frac{1}{C}\frac{B}{A}u(t-i) + \frac{1}{A}e(t-i) \quad -\frac{1}{C}u(t-j) \quad -\frac{1}{C}\varepsilon(t-k) \right) \\
&= \left(\frac{1}{C}\frac{B_1}{A_1}u(t-i) + \frac{1}{A_1}e(t-i) \quad -\frac{1}{C}u(t-j) \quad -\frac{1}{C_1A}\varepsilon(t-k) \right) \\
&\quad + \left(\frac{1}{C}\left(\frac{B}{A} - \frac{B_1}{A_1}\right)u(t-i) \quad 0 \quad \left(\frac{1}{C_1A} - \frac{1}{C_1C}\right)\varepsilon(t-k) \right) \\
&= \left(\frac{B_1}{A_1C}u(t-i) + \frac{1}{A_1A}e(t-i) \quad -\frac{1}{C}u(t-j) \quad -\frac{1}{AC_1}\varepsilon(t-k) \right) \\
&\quad + \left(\frac{1}{AA_1C}(A_1B - AB_1)u(t-i) \quad 0 \quad \frac{\bar{C}-\bar{A}}{C_1AC}\varepsilon(t-k) \right) \\
&\triangleq \psi_1^T(t) + \bar{\psi}^T(t) . \tag{10}
\end{aligned}$$

Next introduce the parameter vector θ_1 of dimension $n_a + n_b + n_c$, corresponding to the reduced order model $A_1y(t) = B_1u(t) + C_1e(t)$

$$\theta_1 = \left(1 \quad (A_1)_1 \quad \dots \quad (A_1)_{n_a-1} \quad 0 \quad (B_1)_1 \quad \dots \quad (B_1)_{n_b-1} \right. \\
\left. 1 \quad (C_1)_1 \quad \dots \quad (C_1)_{n_c-1} \right)^T . \tag{11}$$

Note in particular that

$$\psi_1^T(t)\theta_1 = 0 . \tag{12}$$

Next we examine the part $\bar{\psi}(t)$:

$$\begin{aligned}
\bar{\psi}^T(t) &= \left(\frac{A_1B_1}{AA_1C}(1 + \bar{b}q^{-1} - 1 - \bar{a}q^{-1})u(t-i) \quad 0 \quad \frac{\bar{c}-\bar{a}}{CA}q^{-1}\varepsilon(t-k) \right) \\
&= \left(\frac{B_1}{AC}(\bar{b} - \bar{a})q^{-1}u(t-i) \quad 0 \quad -\frac{(\bar{a}-\bar{c})q^{-1}}{CA}e(t-k) \right) \\
&= O(\delta) . \tag{13}
\end{aligned}$$

Thus the norm of the part $\bar{\psi}(t)$ is of order $O(\delta)$ for all values of t .

In summary so far, we can write

$$V_{\theta\theta} = (P_{\psi_0} + P_{\bar{\psi}}) , \tag{14}$$

where $P_{\psi_0} = E \{ \psi_1(t) \psi_1^T(t) \}$ is rank deficient and positive semidefinite. The null space of this matrix has dimension 1, and is spanned by the vector θ_1 , cf. (12). Further,

$$P_{\tilde{\psi}} = E \left\{ \psi_1(t) \bar{\psi}^T(t) + \bar{\psi}(t) \psi_1^T(t) + \bar{\psi}(t) \bar{\psi}^T(t) \right\} \quad (15)$$

is of order $O(\delta)$ due to (13).

3 The determinant of the Hessian

In this section the determinant of the Hessian (14) is discussed. After a reordering of columns and rows the matrix can first be written as

$$V_{\theta\theta} = \begin{pmatrix} P_0 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{pmatrix}. \quad (16)$$

Here we assume, for generality, that $V_{\theta\theta}$ is $n \times n$, and that P_0 is $(n-k) \times (n-k)$ and positive definite. The second term in (16) is partitioned similarly, so P_{22} is of dimension $k \times k$.

Referring to the specific case treated in (1), (7) we have here

$$n = n_a + n_b + n_c, \quad k = 1. \quad (17)$$

Note that the submatrices P_{ij} , $i, j = 1, 2$ are all of order $O(\delta)$, and P_{22} is $O(\delta^2)$. Using a standard relation for the determinant of a partitioned matrix, see [3], we have

$$\det(V_{\theta\theta}) = \det(P_{22}) \det(P_0 + P_{11} - P_{12} P_{22}^{-1} P_{21}) = O(\delta^{2k}). \quad (18)$$

In the particular case studied in these notes $k = 1$ holds and thus

$$\det(V_{\theta\theta}) = O(\delta^2). \quad (19)$$

What can be said about the behaviour of the approximate bias

$$\tilde{\theta} = - (V_{\theta\theta})^{-1} V_{\theta} ? \quad (20)$$

A naive approach would be to note that the inverse in (20) includes the determinant in the denominator, and thus that the bias would be $O(1/\delta^2)$. This would be in contrast to the result that when an output error structure is used, where the bias is proved to be of order $O(1/\delta)$.

A more careful inspection of (20) would reveal that all elements of the bias can indeed be written as ratios between two polynomials in δ . Further,

the denominator will always contain a factor δ^2 as a consequence of (19). It is though so far an open question how the numerator polynomials behave. Possibly some of them will also have one or several zeros in $\delta = 0$. Specifically one can write for an arbitrary element of $\tilde{\theta}$:

$$\tilde{\theta} = \frac{\beta_0 + \beta_1\delta + \beta_2\delta^2 + \dots + \beta_m\delta^m}{\delta^2(1 + O(\delta))} . \quad (21)$$

The form of the denominator in (21) follows from (19). At this stage the values of the coefficients β_0, β_1, \dots are not known.

In case $\beta_0 \neq 0$, then $\tilde{\theta} = O(1/\delta^2)$. In case $\beta_0 = 0$ and $\beta_1 \neq 0$, then $\tilde{\theta} = O(1/\delta)$, etc.

In the next section we examine the bias for some simple examples, that illustrate the following:

In contrast to the output error case, there is no simple general rule for how the bias behaves for small values of δ .

4 Some examples for a first order system

4.1 Introduction

In this section the bias expression (20) is evaluated numerically for some different cases. All elements of $V_{\theta\theta}$ and V_{θ} consists of covariance elements, and there are standard ways for how to compute them numerically for a finite order model with given parameters.

In the examples considered below it is assumed that

$$n_a = 1, n_b = 2, n_c = 1 . \quad (22)$$

The residual and its gradient can be evaluated as follows:

$$\begin{aligned} \varepsilon(t) &= \frac{A}{C}y(t) - \frac{B}{C}u(t) = \frac{A}{C} \left(\frac{B}{A}u_0(t) + \frac{C}{A}e(t) \right) - \frac{B}{C} (u_0(t) + \tilde{u}(t)) \\ &= e(t) - \frac{B}{C}\tilde{u}(t) , \end{aligned} \quad (23)$$

$$\begin{aligned} \psi^T(t) &= \left(\frac{q^{-1}}{C}y(t) \quad -\frac{q^{-1}}{C}u(t) \quad -\frac{q^{-2}}{C}u(t) \quad -\frac{q^{-1}}{C}\varepsilon(t) \right) \\ &= \left(\frac{q^{-1}B}{AC}u_0(t) + \frac{q^{-1}}{A}e(t) \quad -\frac{q^{-1}}{C}(u_0(t) + \tilde{u}(t)) \quad -\frac{q^{-2}}{C}(u_0(t) + \tilde{u}(t)) \right. \\ &\quad \left. -\frac{q^{-1}}{C}e(t) + \frac{q^{-1}B}{C^2}\tilde{u}(t) \right) . \end{aligned} \quad (24)$$

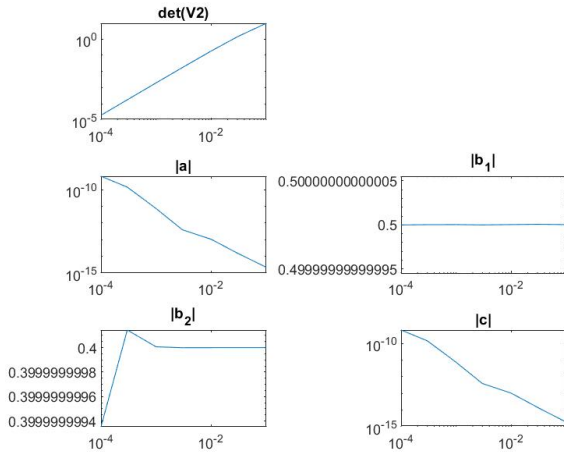


Figure 1: Determinant and parameter biases versus δ for Example 1.

In the examples below, we evaluate the expressions numerically and study the behaviour for small values of δ . By plotting $\det(V_{\theta\theta})$ as well as the absolute values of the elements of the bias vector θ versus δ using log-log scales, one can easily check whether the behaviour is $O(1/\delta^2)$ or something else.

4.2 The examples

Example 1. Let in this example B and C have a joint zero. The parameters are chosen as

$$a = -0.8 + \delta, \quad b_1 = 1, \quad b_2 = -0.8, \quad c = -0.8. \quad (25)$$

Let the noise-free input be white noise of zero mean and variance σ^2 . The variances are chosen as

$$\sigma^2 = 1, \quad \lambda_u^2 = 1, \quad \lambda_e^2 = 1. \quad (26)$$

The obtained results are displayed in Figure 1.

The results are striking. The plots indicate indeed:

- The determinant behaves as $O(1/\delta^2)$ as expected.
- The biases of a and c are zero.
- The bias of b_1 is 0.5 and that of b_2 is -0.4 , independent of the value of δ .

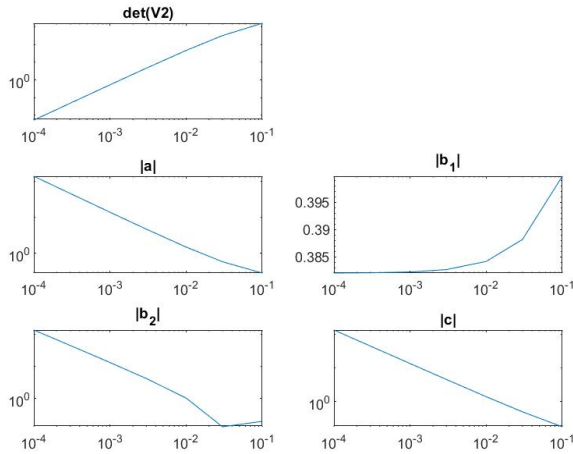


Figure 2: Determinant and parameter biases versus δ for Example 2.

We will return to an explanation of this behaviour a little later, see Section 4.4.

Example 2. We take in this example the same system and parameters as in Example 1, with the modification that the noise-free input is now no longer white noise, but a first order regression

$$u_0(t) + fu_0(t-1) = v(t), \quad Ev^2(t) = \sigma^2, \quad f = -0.9. \quad (27)$$

The obtained results are displayed in Figure 2.

The results differ drastically from Example 1. The plots indicate indeed:

- The determinant behaves as $O(1/\delta^2)$ as expected.
- The bias of a , b_2 and c behave as $O(1/\delta)$.
- The bias of b_1 varies very slowly with δ .

Example 3. Next we modify the original example slightly in another way. In this case we force A and B to have a joint zero. The noise-free input is still assumed to be white noise. The parameters for this example are thus

$$a = -0.8, \quad b_1 = 1, \quad b_2 = -0.8, \quad c = -0.8 + \delta. \quad (28)$$

Let the noise-free input be white of zero mean and variance σ^2 .

The obtained results are displayed in Figure 3.

The plots indicate :

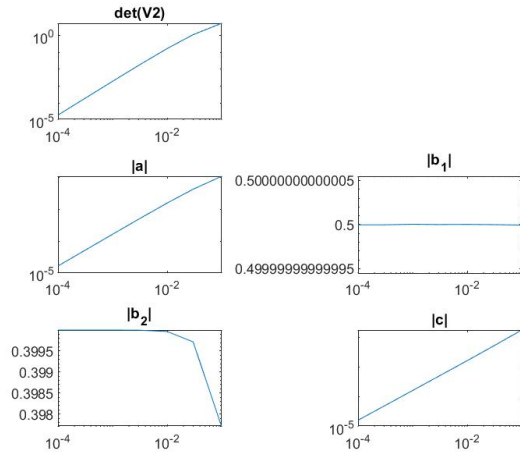


Figure 3: Determinant and parameter biases versus δ for Example 3.

- The determinant behaves as $O(1/\delta^2)$ as expected.
- The biases of a and c behave as $O(\delta)$.
- The biases of b_1 and b_2 vary very slowly with δ .

Example 4. Modify the original example slightly in still another way. In this case force A and C to have a joint zero. The noise-free input is still assumed to be white noise. The parameters for this example are thus

$$a = -0.8, \quad b_1 = 1, \quad b_2 = -0.8 + \delta, \quad c = -0.8. \quad (29)$$

Let the noise-free input be white of zero mean and variance σ^2 .

The obtained results are displayed in Figure 4.

One can now conclude that

- The determinant of $V_{\theta\theta}$ is $O(1/\delta^2)$.
- The bias terms $|a|$ and $|c|$ are both $O(\delta)$.
- The bias term $|b_1|$ does not vary with δ . It is equal to 0.5.
- The bias term $|b_2|$ varies very slowly with δ . For small values of δ it is equal to -0.4 .

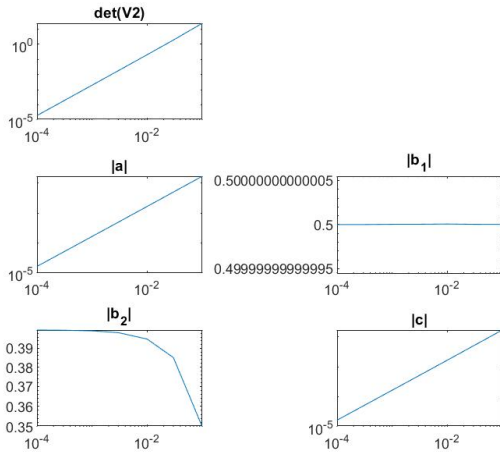


Figure 4: Determinant and parameter biases versus δ for Example 4.

4.3 Conclusions so far

A general conclusion from the four examples is that

The behaviour of the bias as a function of the pole-zero separation δ is quite different in the examples, despite the fact that the examples themselves differ only by slightly modifying the system parameters.

It holds in all the four examples that the determinant of $V_{\theta\theta} = O(\delta^2)$. Further the different examples are characterized as follows:

- Example 1. B is proportional to C , u_0 is white noise. $\tilde{\theta}$ does not depend on δ .
- Example 2. B is proportional to C , u_0 is correlated noise. $|\tilde{\theta}| = O(1/\delta)$.
- Example 3. A is proportional to B . $|\tilde{\theta}| = O(1)$.
- Example 4. A is proportional to C . $|\tilde{\theta}| = O(1)$.

It is not clear from the examples precisely what factors that 'control' the behaviour of $\tilde{\theta}$.

4.4 Analysis of Example 1

The behaviour of the biases in Example 1 can be further analyzed. The obtained numerical results indicate that the biases of the four parameters are constant, and do not depend on δ . Here, we will prove this fact.

Indeed, we claim that in this particular example, the biases satisfies

$$\tilde{\theta} = b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} \begin{pmatrix} 0 \\ 1 \\ c \\ 0 \end{pmatrix}. \quad (30)$$

What is needed to show is that

$$E \{ \psi(t) \psi^T(t) \} \tilde{\theta} = E \{ \psi(t) \varepsilon(t) \},$$

or yet,

$$E \left\{ \psi(t) \left[\psi^T(t) \tilde{\theta} - \varepsilon(t) \right] \right\} = 0. \quad (31)$$

For the specific parameters of Example 1, we have using (30)

$$\begin{aligned} \psi^T(t) \tilde{\theta} - \varepsilon(t) &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} \frac{q^{-1} + cq^{-2}}{C} [u_0(t) + \tilde{u}(t)] - e(t) + b_1 q^{-1} \tilde{u}(t) \\ &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} [u_0(t-1) + \tilde{u}(t-1)] + b_1 \tilde{u}(t-1) - e(t) \\ &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} u_0(t-1) + b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \tilde{u}(t-1) - e(t). \end{aligned} \quad (32)$$

Further,

$$\begin{aligned} &E \left\{ \psi(t) \left[\psi^T(t) \tilde{\theta} - \varepsilon(t) \right] \right\} \\ &= E \left\{ \begin{pmatrix} \frac{b_1 q^{-2}}{A} u_0(t) + \frac{1}{A} e(t-1) \\ -\frac{1}{C} [u_0(t-1) + \tilde{u}(t-1)] \\ -\frac{1}{C} [u_0(t-2) + \tilde{u}(t-2)] \\ -\frac{1}{C} e(t-1) + \frac{b_1}{C} \tilde{u}(t-2) \end{pmatrix} \right. \\ &\quad \left. \times \left[-b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} u_0(t-1) + b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \tilde{u}(t-1) \right] \right\} \\ &= \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \left[b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} \sigma^2 - b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \lambda_u^2 \right] = 0, \end{aligned} \quad (33)$$

which proves (30).

The example can in fact be generalized a bit, with the same behaviour. Assume A to be arbitrary, $B = b_1 q^{-1} C$, $u_0(t)$ white noise of variance σ^2 .

Then an informed guess is that the bias satisfies, cf. (30),

$$\tilde{\theta} = b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} \begin{pmatrix} O_{n_a \times 1} \\ 1 \\ c_1 \\ \vdots \\ c_{n_c} \\ O_{n_c \times 1} \end{pmatrix}. \quad (34)$$

To prove (34) we proceed as earlier. First we establish

$$\begin{aligned} \psi^T(t)\tilde{\theta} - \varepsilon(t) &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} q^{-1} \frac{C}{C} u(t) - e(t) + \frac{b_1 q^{-1} C}{C} \tilde{u}(t) \\ &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} [u_0(t-1) + \tilde{u}(t-1)] - e(t) + b_1 \tilde{u}(t-1) \\ &= -b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} u_0(t-1) + b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \tilde{u}(t-1) - e(t), \quad (35) \end{aligned}$$

and

$$\begin{aligned} &E \left\{ \psi(t) \left[\psi^T(t)\tilde{\theta} - \varepsilon(t) \right] \right\} \\ &= E \left\{ \begin{pmatrix} O_{n_a \times 1} \\ -\frac{q^{-1}}{C} (u_0(t) + \tilde{u}(t)) \\ \vdots \\ -\frac{q^{-n_b}}{C} (u_0(t) + \tilde{u}(t)) \\ O_{n_c \times 1} \end{pmatrix} \right. \\ &\quad \left. \times \left[-b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} u_0(t-1) + b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \tilde{u}(t-1) \right] \right\} \\ &= \begin{pmatrix} O_{n_a \times 1} \\ 1 \\ O_{(n_b-1+n_c) \times 1} \end{pmatrix} \left[-b_1 \frac{\lambda_u^2}{\sigma^2 + \lambda_u^2} \sigma^2 + b_1 \frac{\sigma^2}{\sigma^2 + \lambda_u^2} \lambda_u^2 \right] \\ &= O_{(n_a+n_b+n_c) \times 1}, \quad (36) \end{aligned}$$

which proves the correctness of the hypothesis (34).

5 Conclusions

When an output error model is used, for a model with a small pole-zero separation δ , the obtained parameter bias was shown in [1] to be of order

$O(1/\delta)$. The size of the bias when an ARMAX model is applied was examined in this study. It was shown that no simple and general result applies. In many cases the bias is of order $O(1)$ for small values of δ . In a few cases, the bias is even independent of δ . The bias of individual parameter estimates may be of order $O(1/\delta)$ or $O(\delta)$.

References

- [1] T. Söderström and U. Soverini. When are errors-in-variables aspects particularly important to consider in system identification? Technical Report 2021-006, Department of Information Technology, Uppsala University, Uppsala, Sweden, 2021. Available as <http://www.it.uu.se/research/publications/reports/2021-006>.
- [2] T. Söderström and U. Soverini. Aspects of errors-in-variables identification: When would a standard PEM give a large bias? In *22nd IFAC World Congress, submitted*, Yokohama, Japan, July 12-15 2023.
- [3] T. Söderström and P. Stoica. *System Identification*. Prentice Hall International, Hemel Hempstead, UK, 1989.