# On homogeneous least-squares problems and the inconsistency introduced by mis-constraining

Arie Yeredor[a,*], Bart De Moor[b]

[a] *Department of Electrical Engineering-Systems, Tel-Aviv University, P.O. Box 39040, Tel-Aviv 69978, Israel*
[b] *ESAT-SCD, Katholieke Universiteit Leuven, Kasteelpark Arenberg 10 Leuven B-3001, Belgium*

## Abstract

The term "homogeneous least-squares" refers to models of the form $Ya \approx 0$, where $Y$ is some data matrix, and $a$ is an unknown parameter vector to be estimated. Such problems are encountered, e.g., when modeling auto-regressive (AR) processes. Naturally, in order to apply a least-squares (LS) solution to such models, the parameter vector $a$ has to be somehow constrained in order to avoid the trivial solution $a = 0$. Usually, the problem at hand leads to a "natural" constraint on $a$. However, it will be shown that the use of some commonly applied constraints, such as a quadratic constraint, can lead to inconsistent estimates of $a$. An explanation to this apparent discrepancy is provided, and the remedy is shown to lie with a necessary modification of the LS criterion, which is specified for the case of Gaussian model-errors. As a result, the modified LS minimization becomes a highly non-linear problem. For the case of quadratic constraints in the context of AR modeling, the resulting minimization involves the solution of an equation reminiscent of a "secular equation". Numerically appealing solutions to this equation are discussed.
© 2003 Elsevier B.V. All rights reserved.

*Keywords:* Homogeneous least squares; Constraints; Inconsistency; Maximum likelihood

## 1. Introduction

In many estimation problems, e.g., in the context of the identification (parameter estimation) of linear systems (Söderström and Stoica, 1989), both the inputs and outputs

---

* Corresponding author. Tel.: +972-36405314; fax: +972-36407095.
  *E-mail addresses:* arie@eng.tau.ac.il (A. Yeredor), bart.demoor@esat.kuleuven.ac.be (B. De Moor).

of a system are observed (possibly in the presence of additive noise), from which it is desired to estimate the system's parameters. The standard least-squares (LS) approach for this estimation problem is conceptually and computationally appealing. It consists of seeking the set of parameters with which the linear difference equations relating each output sample to past output and input samples are most closely satisfied (in the sense of a possibly weighted $L2$ norm of the errors vector). Unfortunately, however, the LS estimate is well-known to be biased and inconsistent whenever the past samples used in the regression equations involve noise. Therefore, for such problems LS is only used in cases of very high signal-to-noise ratio.

Extensive research has been addressed over the past two decades towards attempts to modify the LS estimate in such problems, so as to eliminate its bias and regain consistency (De Moor et al., 1994; Stoica and Söderström, 1982; Söderström and Stoica, 1983; Fernando and Nicholson, 1985; Zheng, 1988, 2002a,b; Van Pelt and Bernstein, 2001). Some of the well-known approaches, which have become nearly common practice in system identification, are, e.g., the instrumental variable (Stoica and Söderström, 1982; Söderström and Stoica, 1983) or the Koopmans–Levin (Fernando and Nicholson, 1985) methods.

There are, however, two exceptional cases (except for the trivial noiseless case) in which the LS estimate is generally unbiased and consistent:

- When the observed input is noiseless and the system has a finite impulse response (also termed a "zeros only" system). In this case the past samples involved in the regression equations are only the noiseless input samples. The LS model errors are exactly the output noise, so that for zero-mean noise the resulting LS estimate is unbiased. Moreover, if the output noise is (or can be uniquely transformed into) a sequence of independent, identically distributed random variables, the LS estimate is also consistent. Additionally, if the noise is Gaussian, then the properly weighted LS estimate coincides with the maximum-likelihood (ML) estimate.
- When the system identification problem is actually the identification of an auto-regressive (AR) process, and the observed "output" (namely the process to be identified) is noiseless. The "input", or the process' "driving noise" in this case, is unobserved, which is equivalent to observations of zero, such that the "input observation noise" is actually (minus) the "driving noise". Thus, in such cases, the LS model errors are also the input noise, so that the same conditions (as mentioned above) for unbiased, consistent and ML-equivalent LS estimation prevail.

Although the ordinary LS estimate in these cases is unbiased and consistent, a problematic aspect thereof lies in the possible formulation of such LS problem as constrained homogeneous least-squares (HLS) problems. More generally, HLS problems are problems in which an observed data matrix $Y$ can be modeled approximately by

$$Ya \approx 0, \tag{1}$$

where $a$ is an unknown parameters vector, to be estimated from the observations. The inequality often implies the presence of some "driving noise", which can also be regarded as "modeling errors" in terms of the deviation of $Ya$ from $0$.

A straightforward LS approach would be to estimate $\boldsymbol{a}$ as the minimizer of the norm of $\boldsymbol{Ya}$. However, to avoid the trivial minimizer $\boldsymbol{a} = \boldsymbol{0}$, $\boldsymbol{a}$ has to be properly constrained. Usually in such problems, some "natural" constraint (or set of constraints) on $\boldsymbol{a}$ is dictated by the problem at hand (De Moor et al., 1994; Ninness, 1996; Van Pelt and Bernstein, 2001; Ysebaert et al., 2001). For example, the unbiased, consistent LS estimate for the two estimation problems mentioned above would be obtained by constraining the respective element of $\boldsymbol{a}$ to be 1.

However, it turns out that in general, the solution to the constrained minimization

$$\min_{\boldsymbol{a}} \boldsymbol{a}^{\mathrm{T}} \boldsymbol{Y}^{\mathrm{T}} \boldsymbol{Ya}, \quad \text{s.t. } \boldsymbol{f}(\boldsymbol{a}) = \boldsymbol{0}, \tag{2}$$

where $\boldsymbol{f}(\boldsymbol{a}) = \boldsymbol{0}$ is the set of constraints, can often lead to a biased, inconsistent, and, in a sense, useless estimate of $\boldsymbol{a}$. We shall show that in order to obtain a consistent estimate (subject to a certain pre-specified constraint), the LS criterion $\boldsymbol{a}^{\mathrm{T}} \boldsymbol{Y}^{\mathrm{T}} \boldsymbol{Ya}$ has to be modified, taking into consideration the distribution of the model errors (or "driving noise"). We shall explicitly specify the modification for the case of Gaussian errors.

The paper is structured as follows. We begin with a simple example in the next section, illustrating the problematic aspects of mis-constraining. In Section 3 we present a possible remedy in the form of a modified LS criterion, based on equivalence to ML estimation, for the identification of a first-order AR process. In Section 4 we generalize the criterion to a general-order AR process. A possible computational approach to the minimization of the proposed modified LS criterion is presented in Section 5. Conclusions and summary appear in the closing section.

## 2. An example

We begin by considering an example in which we illustrate the problems induced by choosing a "wrong" constraint. The problem originates from the identification of an AR process as treated, e.g., in Lemmerling and De Moor (2001) (see also Söderström and Stoica (1989) or Yeredor (2000)).

Let $y_n$ be a first-order AR (AR(1)) process satisfying the difference equation:

$$y_n = -a_1 y_{n-1} + e_n, \quad n = 1, 2, \ldots, N, \tag{3}$$

where $e_n$ is a white Gaussian noise process with zero mean and known variance $\sigma_e^2$. It is desired to estimate $a_1$ from the observations $y_1, y_2, \ldots y_N$. We further assume, for convenience, that $y_0$ is deterministically known to be zero.

We can formulate the model equations in matrix form as $\boldsymbol{Ya} = \boldsymbol{e}$, where

$$\boldsymbol{Y} \triangleq \begin{bmatrix} y_1 & 0 \\ y_2 & y_1 \\ \vdots & \vdots \\ y_N & y_{N-1} \end{bmatrix}, \quad \boldsymbol{a} \triangleq \begin{bmatrix} a_0 \\ a_1 \end{bmatrix}, \quad \boldsymbol{e} \triangleq \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix}. \tag{4}$$

Any estimate $\hat{a}$ of $a$ in which $\hat{a}_0 = a_0 = 1$, leads to an estimate $\hat{e}$ of $e$, via $\hat{e} = Y\hat{a}$. Our goal is then to choose $\hat{a}$ such that the norm of $\hat{e}$ is minimized, subject to the *linear* constraint $\hat{a}_0 = 1$:

$$\min_{\hat{a}} \hat{a}^{\mathrm{T}} Y^{\mathrm{T}} Y \hat{a}, \quad \text{s.t. } [1\,0] \cdot \hat{a} = 1. \tag{5}$$

Denoting $\hat{R} \triangleq 1/N\, Y^{\mathrm{T}} Y$, we obtain the well-known LS solution

$$\hat{a}_1 = -\frac{\hat{R}_{2,1}}{\hat{R}_{2,2}}, \tag{6}$$

where $\hat{R}_{i,j}$ denotes the $(i,j)$th element of $\hat{R}$. If $|a_1| < 1$, then $y_n$ is (asymptotically) stationary with autocorrelation satisfying

$$R[0] \triangleq E[y_n^2] = \frac{\sigma_e^2}{1 - a_1^2}, \quad R[1] \triangleq E[y_n y_{n-1}] = -a_1 \cdot R[0]. \tag{7}$$

Moreover, we also have (asymptotically)

$$\hat{R} \overset{N \to \infty}{\Rightarrow} \begin{bmatrix} R[0] & R[1] \\ R[1] & R[0] \end{bmatrix}, \tag{8}$$

so that $\hat{a}_1 \to -R[1]/R[0] = a_1$ is a consistent estimator. Its consistency can be attributed to the fact that it is essentially the ML estimate, whose consistency is guaranteed in this problem setup.

Suppose now, that we want to use a quadratic constraint on $\hat{a}$, and later "normalize" $\hat{a}_0$ to 1. Solving

$$\min_{\hat{a}} \hat{a}^{\mathrm{T}} Y^{\mathrm{T}} Y \hat{a}, \quad \text{s.t. } \hat{a}^{\mathrm{T}} \hat{a} = 1, \tag{9}$$

reduces to an eigenvalue problem, and asymptotically, due to (8), we would eventually always get either $\hat{a}_1 = 1$ (if $R[1] < 0$), or $\hat{a}_1 = -1$ (if $R[1] > 0$), which is (almost) always inconsistent.

In the context of the original problem, the straightforward explanation is that the quadratic constraint is inappropriate (even when followed by normalization), and therefore there is no reason to expect a consistent estimator. We note in this context, that in Van Pelt and Bernstein (2001) it is shown, that it is generally possible to obtain a consistent estimate by using an alternative quadratic constraint of the form $\hat{a}^{\mathrm{T}} N \hat{a} = 1$, where $N$ is some symmetric (not necessarily positive-definite) matrix (that generally depends on the noise statistics). In fact, this is one of the proposed remedies for the LS estimator's bias and inconsistency (as mentioned in the Introduction) in the general case. For our case it is evident, that using $N = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ is equivalent to the linear (sign-ambiguous) constraint $\hat{a}_0 = \pm 1$.

However, an interesting question is—what if the quadratic constraint were indeed part of the problem formulation—would such unreasonable estimates still be obtained?

In order to clarify this, consider an alternative formulation, in which we assume that the process $y_n$ satisfies

$$a_0 y_n = -a_1 y_{n-1} + e_n, \quad n = 1, 2, \ldots, \tag{10}$$

(with the same characteristics of $e_n$ as before), where it is now known that $a_0^2 + a_1^2 = 1$. It is desired to estimate $a_0$ and $a_1$ from $y_1, y_2, \ldots y_N$. Again we assume, for convenience, that $y_0$ is deterministically known to be zero.

Apparently, it is now legitimate to use the quadratically constrained minimization (9)—but then we would get the same highly inconsistent, nearly data-independent, totally unreasonable estimate.

Although this time the constraint is valid, the problem here lies with the objective function. As we shall show immediately, (9) is not the ML criterion, and therefore there is indeed no claim for consistency.

## 3. The correct (ML) criterion

As the matrix-vector product $\mathbf{Y}\mathbf{a}$ is bilinear in the data $\mathbf{y}$ and the vector $\mathbf{a}$, we can also rewrite the model equations as $\mathbf{e} = \mathbf{Y}\mathbf{a} = \mathbf{T}(\mathbf{a})\mathbf{y}$, where

$$\mathbf{T}(\mathbf{a}) \triangleq \begin{bmatrix} a_0 & & & \\ a_1 & a_0 & & \\ & \ddots & \ddots & \\ & & a_1 & a_0 \end{bmatrix}, \quad \mathbf{y} \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}. \tag{11}$$

Note that the matrix $\mathbf{T}(\mathbf{a})$ defined above is square, due to the zero initial conditions ($y_0 = 0$) in (4). For higher-order AR processes this would generalize to assuming further zero initial conditions, namely $y_0 = y_{-1} = \cdots = y_{-p+1} = 0$ for an AR($p$) process. However, often in practice the available data $\mathbf{y}$ are part of a stationary process, and then the "initial conditions" are not zeros, but must be treated as additional (random) unknowns. To avoid complications, it is common practice in these situations to ignore any equations involving "initial conditions" data, and then the respective first rows of $\mathbf{Y}$ in (4) are eliminated, resulting in a rectangular $\mathbf{T}(\mathbf{a})$. A "proper" way of incorporating non-zero initial conditions while maintaining $\mathbf{T}$ square can be found in Yeredor (2000).

Thus, if $\mathbf{e}$ is a zero-mean Gaussian vector with covariance $\sigma_e^2 \mathbf{I}$, then $\mathbf{y} = \mathbf{T}^{-1}(\mathbf{a})\mathbf{e}$ is also a zero-mean Gaussian vector, but its covariance is $\sigma_e^2 \mathbf{T}^{-1}(\mathbf{a})\mathbf{T}^{-T}(\mathbf{a})$. Therefore its distribution is given by

$$f(\mathbf{y}; \mathbf{a}) = \frac{1}{|2\pi\sigma_e^2 \mathbf{T}^{-1}(\mathbf{a})\mathbf{T}^{-T}(\mathbf{a})|^{1/2}} \, e^{-\frac{1}{2\sigma_e^2} \mathbf{y}^T \mathbf{T}^T(\mathbf{a})\mathbf{T}(\mathbf{a})\mathbf{y}}, \tag{12}$$

the logarithm of which is given by

$$L(y; a) = \log f(y; a)$$

$$= c + \log |T(a)| - \frac{1}{2\sigma_e^2} y^T T^T(a) T(a) y$$

$$= c + \log |T(a)| - \frac{1}{2\sigma_e^2} a^T Y^T Y a, \tag{13}$$

where $c$ is an irrelevant constant and $|\cdot|$ denotes the determinant. Evidently, in an AR problem, it is easy to observe from (11), that $|T(a)| = a_0^N$, so that the maximization of the likelihood $L(y; a)$ is equivalent to the following minimization problem:

$$\min_{\hat{a}} \{\hat{a}^T Y^T Y \hat{a} - N\sigma_e^2 \log \hat{a}_0^2\}, \quad \text{s.t. } \hat{a}^T \hat{a} = 1, \tag{14}$$

which would yield the consistent ML estimate of $a$, in contrast to the "wrong" minimization problem of (9).

In the appendix we verify the consistency of the resulting estimate by deriving the closed-form solution to this simple, two-dimensional problem.

## 4. Generalization and discussion

Straightforward generalization to the general-order AR(q) ($q \geqslant 1$) process with general constraints $f(a) = 0$, maintains the same objective function:

$$\min_{\hat{a}} \{\hat{a}^T Y^T Y \hat{a} - N\sigma_e^2 \log \hat{a}_0^2\} \quad \text{s.t. } f(\hat{a}) = 0.$$

The general constraints can be any (linear or nonlinear) constraints that are justified by the model, such as (but not limited to) linear constraints reflecting known coefficients (most commonly $\hat{a}_0 = 1$) or known poles. Evidently, if (and only if) one or more of the constraints impose $\hat{a}_0 = 1$, then the second term of the objective function is zeroed out, and we obtain the classical objective function of (9). It is interesting to note, however, that no artificial constraints are actually needed (unless the available a-priori information would dictate so), because the "trivial" solution $\hat{a} = 0$ is no longer a minimizer (due to the log term).

For the more general case of an HLS (not necessarily AR) problem, the entire matrix $T(a)$ has to be incorporated into the LS criterion. If, in addition, the model errors are assumed to have a general (known) covariance structure $C_e$, then the resulting minimization assumes the form

$$\min_{\hat{a}} \{\hat{a}^T Y^T C_e^{-1} Y \hat{a} - 2 \log |T(a)|\} \quad \text{s.t. } f(\hat{a}) = 0.$$

A potentially weak point of this approach, is that prior knowledge of the noise variance $\sigma_e^2$ (or the noise covariance $C_e$) is required. In some cases the noise statistics are

indeed known a-priori, e.g., through knowledge of a physical model or of technical specifications, such as a receiver's input noise-figure. In other cases it is sometimes possible to estimate the noise statistics "off-line" when the signal of interest is muted. When such means are not available, it may still be possible to employ an iterative strategy in which, following an intelligent initial guess, the noise level is re-estimated from the implied residual $\hat{e}_n$, and the parameter estimates are refined accordingly in each iteration; however, such a strategy should be applied with caution, so as to avoid a misleading feedback of error between iterations.

## 5. Minimization of the modified LS criterion

In this section, we discuss the minimization of the modified LS criterion (14) for the general-order AR(p) model with a quadratic constraint. Given the estimated (symmetric, positive definite) $M \times M$ correlation matrix $\hat{R}$ and the model error variance $\sigma_e^2$, we wish to minimize (with respect to (w.r.t.) $a$)

$$\min_{a}\{a^{\mathrm{T}}\hat{R}a - \sigma_e^2 \log(a_0^2)\} \quad \text{s.t. } a^{\mathrm{T}}a = 1, \tag{15}$$

where $a_0$ denotes the first element of $a$.

We form the Lagrangian,

$$L(a, \mu) = a^{\mathrm{T}}\hat{R}a - \sigma_e^2 \log(a_0^2) - \mu(a^{\mathrm{T}}a - 1), \tag{16}$$

differentiating w.r.t. $a$ (taking advantage of the symmetry of $\hat{R}$), and equating zero, we obtain

$$(\hat{R} - \mu I)a = \frac{\sigma_e^2}{a_0} \cdot i_1, \tag{17}$$

where $I$ denotes the $M \times M$ identity matrix and $i_1$ denotes its first column. Naturally, differentiation of $L(a, \mu)$ w.r.t. $\mu$ further yields the constraint $a^{\mathrm{T}}a = 1$.

Using the eigenvalue decomposition $\hat{R} = U\Lambda U^{\mathrm{T}}$ (with $U$ a unitary matrix and $\Lambda$ diagonal) and defining $\tilde{a} \triangleq a_0 \cdot a$, we may rewrite (17) as

$$U(\Lambda - \mu I)U^{\mathrm{T}}\tilde{a} = \sigma_e^2 i_1, \tag{18}$$

hence

$$\tilde{a} = \sigma_e^2 U(\Lambda - \mu I)^{-1}U^{\mathrm{T}}i_1, \tag{19}$$

or, defining $v \triangleq U^{\mathrm{T}}i_1$,

$$\tilde{a} = \sigma_e^2 U(\Lambda - \mu I)^{-1}v. \tag{20}$$

We now need to address the constraint $a^{\mathrm{T}}a = 1$. Noting that

$$\|\tilde{a}\|_2^2 = a_0^2 \cdot \|a\|_2^2, \tag{21}$$

we conclude that the constraint is satisfied if and only if $\|\tilde{a}\|_2^2 = a_0^2$. From (20) we have

$$\|\tilde{a}\|_2^2 = \sigma_e^4 v^{\mathrm{T}}(\Lambda - \mu I)^{-2}v. \tag{22}$$

On the other hand, $a_0^2$ is simply the first element of $\tilde{a}$, given by

$$a_0^2 = i_1^T \tilde{a} = \sigma_e^2 v^T (A - \mu I)^{-1} v. \tag{23}$$

Our constraint can thus be expressed as

$$\sigma_e^2 v^T (A - \mu I)^{-2} v = v^T (A - \mu I)^{-1} v, \tag{24}$$

or

$$\sigma_e^2 \sum_{m=1}^{M} \frac{v_m^2}{(\lambda_m - \mu)^2} = \sum_{m=1}^{M} \frac{v_m^2}{\lambda_m - \mu}, \tag{25}$$

where $v_m$ and $\lambda_m$ denote the $m$th elements of $v$ and the $(m, m)$th element of $A$, respectively. This expression leads to a polynomial (of degree $2M - 1$ at most) in $\mu$,

$$\sum_{m=1}^{M} \left\{ v_m^2 (\mu + \sigma_e^2 - \lambda_m) \prod_{n \neq m} (\mu - \lambda_n)^2 \right\} = 0, \tag{26}$$

whose rooting would yield at most $2M - 1$ possible real-valued solutions in $\mu$. Each of the candidate (real-valued) solution can be plugged into (20), yielding a candidate $\tilde{a}$. Dividing each element of $\tilde{a}$ by $\sqrt{\tilde{a}_0}$ would yield a candidate unit-norm solution for $a$. Of the resulting $2M - 1$ (at most) solutions, the one that yields the smallest objective-function value $a^T \hat{R} a - \sigma_e^2 \log(a_0^2)$ is to be chosen as the global minimizer.

To avoid the need for general polynomial rooting, one may observe the resemblance of the left-hand side (LHS) or right-hand side (RHS) of (25) to the form known as a "secular equation" (see, e.g. (Gu and Eisenstat, 1995a,b)). Then, with vertical asymptotes at the locations of the eigenvalues $\lambda_m$ (a typical situation is illustrated in Fig. 1), the graph for the LHS (as a function of $\mu$) would resemble "parabolic" curves between these asymptotes, while the graph for the RHS would resemble monotonic "cubic power" curves between the asymptotes (each extending from $-\infty$ at the left asymptote to $+\infty$ at the right asymptote). The solutions in that case are particularly easy to compute numerically, because they interlace with the eigenvalues, and can be found by simple bisections; additionally, denoting the smallest and largest eigenvalues by $\lambda_{\min}$ and $\lambda_{\max}$ (respectively), there are no solutions larger than $\lambda_{\max}$ (since for all $\mu > \lambda_{\max}$, the LHS is positive and the RHS is negative), and there is always at least one real-valued solution, smaller than $\lambda_{\min}$ (since when $\mu \to \lambda_{\min}$ from below, the LHS is larger than the RHS, whereas when $\mu \to -\infty$ the LHS is smaller than the RHS- and both are continuous in $(-\infty, \lambda_{\min})$, so they must intersect).

Additionally, although in general any solution of (25) is merely a stationary point of the Lagrangian, and can therefore be either a minimum, a maximum or a saddle point, it can be observed that $\mu$ values below $\lambda_{\min}$ are guaranteed to be associated with (at least local) minima. To show that, we examine the second derivative matrix (Hessian) of the Lagrangian (16) w.r.t. $a$:

$$H \triangleq \frac{\partial^2 L(a, \mu)}{\partial a^2} = \hat{R} - \mu I + \frac{\sigma_e^2}{a_0^2} I_{11} = U(A - \mu I) U^T + \frac{\sigma_e^2}{a_0^2} I_{11}, \tag{27}$$
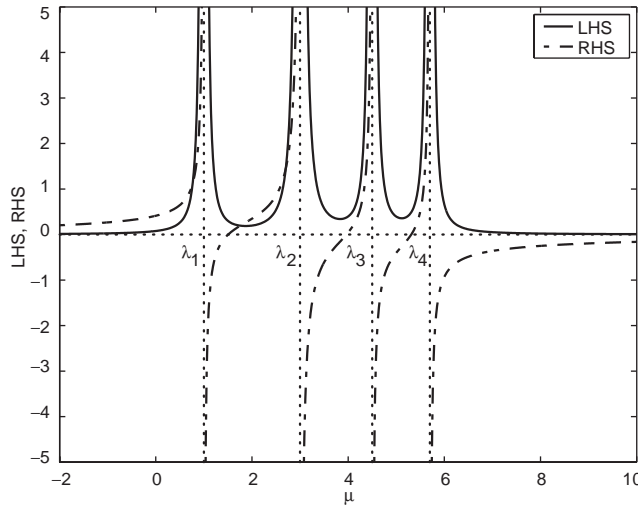
Fig. 1. A typical pattern of the LHS vs. the RHS of (25).

where $I_{11}$ denotes an $M \times M$ matrix with all-zeros entries, except for its $(1, 1)$ element, which is 1. Indeed, for $\mu < \lambda_{\min}$, the first term is positive-definite, hence (since $I_{11}$ is positive-semidefinite) the Hessian is positive-definite and the associated solution is guaranteed to be a minimum.

Although the converse is not guaranteed in general,[1] we conjecture that the solution associated with the *smallest* $\mu$ (which is always below $\lambda_{\min}$) is always the global minimizer of criterion (15). This conjecture has been supported by extensive experimentation, yet we were unable to provide a rigorous proof. It may be interesting to note, in this context, that when $\mu = 0$, the solution to (17) coincides with the maximum entropy (ME) estimate of $a$, and for different values of $\mu$ the associated solutions of (17) can be regarded as "modified ME" estimates of $a$. Thus, the smallest (absolute) value of $\mu$ implies "the least modification" of the ME estimate.

## 6. Conclusion

We have presented and explained the observation, that the constrained HLS approach can sometimes yield useless estimates. The reason is that the LS criterion in these cases has to be supplemented with an additional term in order to yield consistent estimates in a statistical framework. Fortunately, in several common applications this term is automatically zeroed-out by the constraint; however, when using constraints that do not guarantee zero value to this term, one has to take precaution not to exclude the term from the minimization. Thus, from an algebraic point of view, monic constraints

---

[1] For $\mu > \lambda_{\min}$ the Hessian may still be positive-definite.

may often be the most 'manageable', albeit theoretically not the only ones possible (even not from the point of view of ML).

We specified this additional term for the case of Gaussian "driving noise" (or "model errors"), and discussed possible numerical approaches for the resulting minimization. Note, however, that for the problems illustrated, the consistent estimators of the parameters are based on consistent estimates of second-order statistics of the data, so for ergodic processes with finite second-order moments, the consistency is maintained regardless of the distribution, which may be non-Gaussian.

### Appendix A. Closed-form solution for the AR(1) example

We shall show that the solution of the "correct" minimization problem (14) yields a consistent estimate. Dividing by $N$, we obtain the equivalent problem

$$\min_{\hat{a}}\{\hat{a}^{\mathrm{T}}\hat{R}\hat{a} - \sigma_e^2 \log \hat{a}_0^2\} \quad \text{s.t. } \hat{a}^{\mathrm{T}}\hat{a} = 1. \tag{A.1}$$

By parameterizing the constraint as $\hat{a}_0 = \cos(\hat{\theta})$ and $\hat{a}_1 = \sin(\hat{\theta})$ for some single parameter $\hat{\theta}$, we obtain the following equivalent *unconstrained* minimization:

$$\min_{\hat{\theta}}\{\cos^2(\hat{\theta})\hat{R}_{1,1} + 2\cos(\hat{\theta})\sin(\hat{\theta})\hat{R}_{1,2} + \sin^2(\hat{\theta})\hat{R}_{2,2} - \sigma_e^2 \log(\cos^2(\hat{\theta}))\}. \tag{A.2}$$

As already mentioned, assuming that $|a_1/a_0| < 1$, $y_n$ is (asymptotically) a stationary process, and $\hat{R}$ tends (as $N \to \infty$) to the true $R$ (as in (8)). Thus, with $\hat{R}_{1,1} = \hat{R}_{2,2} = R[0]$, and $\hat{R}_{1,2} = R[1]$, this minimization problem reduces to

$$\min_{\hat{\theta}}\{\sin(2\hat{\theta})R[1] - \sigma_e^2 \log(\cos^2(\hat{\theta}))\}. \tag{A.3}$$

(Note that without the second term, a minimum is always obtained either at $\hat{\theta} = \pi/4$ or at $\hat{\theta} = 3\pi/4$, depending only on the sign of $R[1]$, as already observed for the "wrong"

minimization problem earlier). Differentiating and equating to zero we obtain that the minimizing $\hat{\theta}$ should satisfy

$$\frac{\tan(\hat{\theta})}{\cos(2\hat{\theta})} = -\frac{R[1]}{\sigma_e^2}. \tag{A.4}$$

Indeed, for the stationary process

$$y_n = -\frac{\sin(\theta)}{\cos(\theta)}\,y_{n-1} + \frac{1}{\cos(\theta)}\,e_n, \tag{A.5}$$

we have

$$R[0] = \frac{\sigma_e^2}{\cos^2(\theta)}\,\frac{1}{1 - \tan^2(\theta)} = \frac{\sigma_e^2}{\cos(2\theta)} \tag{A.6}$$

and

$$R[1] = -\tan(\theta)R[0] = -\sigma_e^2\,\frac{\tan(\theta)}{\cos(2\theta)}, \tag{A.7}$$

from which the consistency of $\hat{\theta}$ is evident in view of (A.4).

# References

De Moor, B., Gevers, M., Goodwin, G., 1994. L2-Overbiased, L2-underbiased and L2-unbiased estimation of transfer functions. Automatica 30 (5), 893–898.

Fernando, K.V., Nicholson, H., 1985. Identification of linear systems with input and output noise: the Koopmans–Levin method. IEE Proc. D 132, 30–36.

Gu, M., Eisenstat, S.C., 1995a. A Divide-and-conquer algorithm for the bidiagonal SVD. SIAM J. Matrix Anal. Appl. 16, 79–92.

Gu, M., Eisenstat, S.C., 1995b. A Divide-and-conquer algorithm for the symmetric tridiagonal eigenvalue problem. SIAM J. Matrix Anal. Appl. 16, 172–191.

Lemmerling, P., De Moor, B., 2001. Misfit versus latency. Automatica 37, 2057–2067.

Ninness, B., 1996. Integral constraints on the accuracy of least-squares estimation. Automatica 32 (3), 391–397.

Söderström, T., Stoica, P., 1983. Instrumental Variable Methods for System Identification. Springer, New York.

Söderström, T., Stoica, P., 1989. System Identification. Prentice-Hall, Englewood Cliffs, NJ.

Stoica, P., Söderström, T., 1982. Bias correction in least-squares identification. Internal J. Control 35, 449–457.

Van Pelt, T.H., Bernstein, D.S., 2001. Quadratically constrained least-squares identification. Proceedings of the American Control Conference, Arlington, VA, USA, 25–27 June 2001, pp. 3684–3689.

Yeredor, A., 2000. The joint MAP-ML criterion and its relation to ML and to extended least-squares. IEEE Trans. Signal Process. 48 (12), 3484–3492.

Ysebaert, G., Van Acker, K., Moonen, M., De Moor, B., 2001. Constraints in channel shortening equalizer design for DMT-based systems. Internal Report 01–27, ESAT-SISTA, K.U. Leuven, Leuven, Belgium, 2001.

Zheng, W.-X., 1988. Consistent estimation of parameters of stochastic feedback systems in the presence of correlated disturbances. Adv. Modelling Simulation 14, 15–26.

Zheng, W.-X., 2002a. Noisy input–output system identification—the Koopmans–Levin method revisited. Proceedings of the Conference on Decision and Control, Las Vegas, NV, December 2002, pp. 636–637.

Zheng, W.-X., 2002b. A bias correction method for identification of linear dynamic errors-in-variables models. IEEE Trans. Automat. Control 47, 1142–1147.