

Gross Error Detection When Constraints Are Bilinear

Derrick K. Rollins and Shonda D. Roelfs

Depts. of Chemical Engineering and Statistics, Iowa State University, Ames, IA 50011

Rollins and Davis (1992) introduced an approach that can obtain unbiased estimators for true values of process variables when biased process measurements and leaks exist. They addressed linear physical constraints only. The purpose of this work is to discuss ways to extend their approach to cases where constraints are bilinear. Bilinear terms exist when two measured variables (each to the power of 1) are multiplied by each other. They are statistically more complex than linear constraints. This complexity is due to product of variables not having normal distributions when their individual variables do have normal distributions. For simplicity, this approach is described in the context of bilinear energy balances; however, it is applicable to any kind of bilinear physical constraint such as component mass balances. Also, for simplicity, only measurement biases will be treated although this approach is also applicable to process leaks.

The approach of Rollins and Davis (1992) to determine unbiased estimators of process variables consists of the following steps: 1. identification of a specific number of unbiased measurements (δ 's = 0); 2. estimation of the δ 's not identified in step 1 (These estimates are statistically tested for significance. Thus, this step is used to identify nonzero δ 's.); and 3. estimation of process variables. The success of step 1 allows the unbiased estimation of the δ 's in step 2, which allows the unbiased estimation of process variables in step 3. The identification procedure in step 1 involves the *selected* testing of hypotheses for material or energy balance closure. Thus, there are two critical elements to identification: the development of accurate and powerful test statistics (that are likely to make correct conclusions when δ 's are zero and nonzero) and the development of a strategy for optimal hypothesis test selection [needed because the ability not to misidentify *biased* measurements by hypothesis testing is maximized as the number of superfluous hypotheses are minimized (Rollins, 1990)]. We see the first element, development of accurate and powerful test statistics, as the most critical and the most challenging. We take this position because with accurate and powerful tests even a poor test selection strategy (for example, one that tests a very large number of hypotheses) can identify zero δ 's accurately. Even with the best test selection strategy, however,

inaccurate test statistics will lead to inaccurate identification. Hence, in this work we describe three ways to test hypotheses when constraints are bilinear. Two approaches approximate the true statistical distribution and use normal theory. Consequently, they can be used now. We also give equations to obtain estimates of δ 's (step 2) for both approaches. The third approach proposes testing hypotheses by using the exact distribution. From a search in the statistical literature, we have concluded that this distribution has not been developed. Hence, we discuss our plans to develop this distribution and use it to test hypotheses. Note that optimal test selection is not addressed here not only due to the limited space but because it can be treated as a general topic independent of the types of constraints present. The interested reader may consult Rollins (1990) for a hypothesis testing strategy using an optimal test selection criteria.

Model

For a chemical process network of streams and nodes (interconnecting units), the following set of equations can be used to describe physical and statistical relationships for measured variables in the absence of chemical reactions.

$$F = \mu_F + \delta_F + \epsilon_F \quad (1)$$

$$T = \mu_T + \delta_T + \epsilon_T \quad (2)$$

such that

$$A \cdot \mu_F = 0 \quad (3)$$

$$A \cdot \text{Diag}(\mu_F) \cdot \mu_T + Q = 0 \quad (4)$$

with

$$\epsilon_F \sim N_p(0, \Sigma_F) \quad (5)$$

$$\epsilon_T \sim N_p(0, \Sigma_T) \quad (6)$$

$$\text{COV}(\epsilon_T, \epsilon_F) = 0. \quad (7)$$

Correspondence concerning this work should be addressed to D. K. Rollins.

For a case with several components, each component would have an equation like Eq. 1 or 2. In addition note that by Eq. 7, $\text{COV}(F, T) = \mathbf{0}$. Equations 3 and 4 represent the true total mass and energy balances, respectively. They also indicate that process leaks are assumed to be zero (for simplicity) and that steady-state conditions apply. Equation 3 is said to be linear in flow variables, and Eq. 4 is bilinear because it is a function of products of flow and temperature variables. Equation 7 simply means that errors made in flow and temperature measurements are unrelated.

The test statistics and estimators used here will be developed from the following transformations of the above equations. Let

$$r = AF. \quad (8)$$

Then, using Eq. 3, the mean and variance of r are given as:

$$E[r] = \mu_r = A(\mu_F + \delta_F) = A\delta_F \quad (9)$$

$$\begin{aligned} \text{Var}(r) &= \Sigma_r = A \cdot \text{Var}(F) \cdot A^T \\ &= A\Sigma_F A^T. \end{aligned} \quad (10)$$

Thus,

$$r \sim N_q(\mu_r, \Sigma_r). \quad (11)$$

Additionally, let

$$s = A \cdot \text{Diag}(F) \cdot T + Q. \quad (12)$$

Then, from Eq. 4, the $E[s]$ and $\text{Var}[s]$ are:

$$\begin{aligned} E[s] &= \mu_s = A \cdot E[\text{Diag}(F)] \cdot E[T] + Q \\ &= A \cdot \text{Diag}(E[F]) \cdot (\mu_T + \delta_T) + Q \\ &= A \cdot \text{Diag}(\mu_F + \delta_F) \cdot (\mu_T + \delta_T) + Q \\ &= A \cdot \text{Diag}(\mu_F + \delta_F) \cdot \delta_T + A \cdot \text{Diag}(\delta_F) \cdot \mu_T \end{aligned} \quad (13)$$

$$\begin{aligned} \text{Var}(s) &= \Sigma_s = A \cdot [\Sigma_F \Sigma_T + \Sigma_F \cdot \text{Diag}(\mu_T) \cdot \text{Diag}(\mu_T) \\ &\quad + \Sigma_T \cdot \text{Diag}(\mu_F) \cdot \text{Diag}(\mu_F)] \cdot A^T = A \cdot D \cdot A^T. \end{aligned} \quad (14)$$

Thus,

$$s \sim \mathcal{N}_q(\mu_s, \Sigma_s), \quad (15)$$

where \mathcal{N} is used to represent a distribution that is a sum of terms, and each term is a product of two normal random variables. The development for the distribution of r (Eq. 11) is the same as the one given by Rollins and Davis (1992) for linear constraints. However, the distribution for s is new to this work. We now consider three ways of using s to help identify and estimate measurement biases.

Two-Stage Approach

The first approach that we will introduce consists of two

steps. In the first step we attempt to obtain unbiased estimates for flow variables using the method of Rollins and Davis (1992). In the second step, the flow estimates are assumed to be their true values, which means that s is assumed to be linear in the temperature variables. This approach will be reasonably valid as long as the flow estimates are accurate. The accuracy of the flow estimates will depend on the variances of the measured variables and on the accuracy of the identification and estimation of δ_F . Additionally, as noted by Rollins and Davis (1992), it may not always be possible to estimate δ_F . The mathematical details of both steps will now be given.

The first step, obtaining the flow estimator, follows Rollins and Davis (1992) very closely except for some notational changes and the assumed existence of process leaks. Let r be partitioned as follows (although δ_{F1} is shown to be a $q \times 1$ vector, it may be any vector with dimension $\leq q$),

$$\begin{aligned} r &= A\delta_F + A\epsilon_F = [A_{11}^{q \times q} A_{12}^{q \times (p-q)}] \begin{bmatrix} \delta_{F1}^{q \times 1} \\ \delta_{F2}^{(p-q) \times 1} \end{bmatrix} + A\epsilon_F \\ &= A_{11}\delta_{F1} + A\epsilon_F \end{aligned} \quad (16)$$

with $\delta_{F2} = \mathbf{0}$ and the rank $(A_{11}) = q$. Now let

$$\hat{\delta}_{F1} = A_{11}^{-1}r, \quad (17)$$

$$\hat{\delta}_F = \begin{bmatrix} \hat{\delta}_{F1} \\ \mathbf{0} \end{bmatrix} \quad (18)$$

and the estimator for μ_F is:

$$\hat{\mu}_F = F - \hat{\delta}_F \quad (19)$$

Therefore, if $\delta_{F2} = \mathbf{0}$, $E[\hat{\delta}_F]$ and $E[\hat{\mu}_F]$ are given as:

$$E[\hat{\delta}_F] = \begin{bmatrix} E[\hat{\delta}_{F1}] \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \delta_{F1} \\ \mathbf{0} \end{bmatrix} = \delta_F \quad (20)$$

$$E[\hat{\mu}_F] = \mu_F + \delta_F - \delta_F = \mu_F. \quad (21)$$

Hence, if $\delta_{F2} = \mathbf{0}$ and the rank $(A_{11}) = q$, then $\hat{\delta}_F$ is an unbiased estimator for δ_F . $\hat{\delta}_F$ is also a consistent estimator of δ_F , which means as n (the sample size) increases, $\hat{\delta}_F$ approaches δ_F . Thus, for a significantly large n , $\hat{\delta}_F$ will be close to δ_F , and thus, $\hat{\mu}_F$ will be close to μ_F .

When $\hat{\mu}_F$ is very close to μ_F and substituted into the energy balance constraints, s will statistically behave as its linear in temperature constraints. When this substitution is valid:

$$s \approx A \cdot \text{Diag}(\hat{\mu}_F) \cdot T + Q. \quad (22)$$

From Eq. 22, the $E[s]$ and the $\text{Var}(s)$ are computed as follows:

$$\begin{aligned} E[s] &= \mu_s \approx A \cdot \text{Diag}(\hat{\mu}_F) \cdot E[T] + Q \\ &= A \cdot \text{Diag}(\hat{\mu}_F) \cdot [\mu_T + \delta_T] + Q \\ &= A \cdot \text{Diag}(\hat{\mu}_F) \cdot \delta_T + Q \end{aligned} \quad (23)$$

$$\text{Var}(s) = \Sigma_s \approx A \cdot \text{Diag}(\hat{\mu}_F) \cdot \Sigma_T \cdot [A \cdot \text{Diag}(\hat{\mu}_F)]^T \quad (24)$$

Therefore, approximately,

$$s \sim N_q(\mu_s, \Sigma_s) \quad (25)$$

Proceeding as in Rollins and Davis (1992), let

$$M = A \cdot \text{Diag}(\hat{\mu}_F) = [M_{11}^{q \times q} M_{12}^{q \times (p-q)}] \quad (26)$$

with rank $(M_{11}) = q$. Therefore,

$$\begin{aligned} s &\approx M \cdot T + Q = M\mu_T + M\delta_T + M\epsilon_T + Q \\ &= M\delta_T + M\epsilon_T \\ &= [M_{11} M_{12}] \begin{bmatrix} \delta_{T1}^{q \times 1} \\ \delta_{T2}^{(p-q) \times 1} \end{bmatrix} + M\epsilon_T \\ &= M_{11}\delta_{T1} + M\epsilon_T \end{aligned} \quad (27)$$

if $\delta_{T2} = 0$. Now let

$$\hat{\delta}_T = \begin{bmatrix} \hat{\delta}_{T1} \\ 0 \end{bmatrix} = \begin{bmatrix} M_{11}^{-1} s \\ 0 \end{bmatrix}, \quad (28)$$

which gives $E[\hat{\delta}_T] = \delta_T$, if the rank $(M_{11}) = q$ and $\delta_{T2} = 0$ as before. The proposed estimator for μ_T in this case is:

$$\hat{\mu}_T = T - \hat{\delta}_T \quad (29)$$

Therefore, if this approach is valid, then the $E[\hat{\mu}_T] \approx \mu_T$ and approximate $100(1 - \alpha)\%$ simultaneous confidence intervals may be obtained following the procedure described by Rollins and Davis (1992). Additionally, if Eq. 25 is reasonably accurate, then Eq. 24 may be used to test hypotheses of the form $H_0: l^T \mu_s = 0$ vs. $H_a: l^T \mu_s \neq 0$, where $l^T \mu_s$ is some linear combination of the components of μ_s . The test statistics to test these hypotheses would be similar to those described by Rollins and Davis (1992).

Linearization

The second approach consists of linearizing the error terms by taking a Taylor series expansion about the means. Here, all δ 's can be simultaneously estimated which should provide better accuracy, but this approach will not be valid when bilinear effects are significantly large. Before linearizing, r and s will be rearranged into one vector, u , and described as follows:

$$\begin{aligned} u^{2q \times 1} &= \begin{bmatrix} r \\ s \end{bmatrix} = \begin{bmatrix} AF \\ A \cdot \text{Diag}(F) \cdot T + Q \end{bmatrix} \\ &= \begin{bmatrix} A & 0 \\ 0 & A \end{bmatrix} \begin{bmatrix} F \\ \text{Diag}(F) \cdot T \end{bmatrix} + \begin{bmatrix} 0 \\ Q \end{bmatrix} = \dot{A} \begin{bmatrix} F \\ \text{Diag}(T) \cdot F \end{bmatrix} + \Omega \\ &= \dot{A} \begin{bmatrix} I & 0 \\ 0 & \text{Diag}(T) \end{bmatrix} \begin{bmatrix} F \\ F \end{bmatrix} + \Omega = \dot{A} \begin{bmatrix} I & 0 \\ 0 & \text{Diag}(F) \end{bmatrix} \begin{bmatrix} F \\ T \end{bmatrix} + \Omega \\ &= \dot{A} B b + \Omega. \end{aligned} \quad (30)$$

where b contains the flow and temperature variables modeled with bias, and B contains the variables that are modeled as being unbiased. Expanding and linearizing u gives:

$$\begin{aligned} u &\approx \dot{A} E[B] E[b] + \dot{A} E[B] (b - E[b]) + \dot{A} (B - E[B]) E[b] + \Omega \\ &= \dot{A} E[B] b + \dot{A} B E[b] - \dot{A} E[B] E[b] + \Omega \end{aligned} \quad (31)$$

Therefore,

$$\begin{aligned} E[u] &\approx \dot{A} E[B] E[b] + \Omega \\ &= \dot{A} \mu_B (\mu_b + \delta_b) + \Omega \\ &= \dot{A} \mu_B \delta_b \end{aligned} \quad (32)$$

Equation 32 motivates the following estimator for δ_b :

$$\hat{\delta}_b = \begin{bmatrix} \hat{\delta}_{b1} \\ \hat{\delta}_{b2} \end{bmatrix} = \begin{bmatrix} C_{11}^{-1} u \\ 0 \end{bmatrix} \quad (33)$$

where

$$C = [C_{11}^{2q \times 2q} C_{12}^{2q \times (2p-2q)}] = \dot{A} B, \quad (34)$$

$$\delta_b = \begin{bmatrix} \delta_{b1}^{2q \times 1} \\ \delta_{b2}^{(2p-2q) \times 1} \end{bmatrix}, \quad (35)$$

δ_{b2} is assumed to be zero, and the rank (C_{11}) is assumed to be $2q$. Note that since Eq. 31 is statistically linear in normal random variables, normal theory may be used to test hypotheses about $E[u]$. Specifically, the testing procedures used by Rollins and Davis (1992) may be applied with the proper substitutions.

Product of Normals

The true distribution of s is \mathfrak{N} . It is \mathfrak{N} because each term s_i , an element of s , is the product of a temperature variable that is normally distributed and a flow variable that is normally distributed. The probability distribution function for the product of two normally distributed random variables was given by Craig (1936). In this work, Craig (1936) specifically addressed random variables of the following type: $Z = X/\sigma_x \cdot Y/\sigma_y$, where σ_x^2 and σ_y^2 are the variances of the random variables X and Y , respectively. Later Aroian (1947) showed that Z approaches a normal distribution as X/σ_x and Y/σ_y approach ∞ . Meeker et al. (1981) prepared tabled values for the distribution of Z by numerically integrating the probability distribution function. However, the works of Craig (1936), Aroian (1947), and Meeker et al. (1981) do not directly apply to our setting which is of the form:

$$Z = X_1 X_2 + X_3 X_4 + \dots X_t X_{t+1} \quad (36)$$

such that

$$\begin{aligned} X_i &\sim N(\mu_{X_i}, \sigma_{X_i}^2) \\ i &= 1, \dots, t+1. \end{aligned} \quad (37)$$

To develop identification tests that use the \mathcal{N} distribution, its probability distribution function will need to be found and applied. We are considering two approaches; a likelihood ratio test and the method of Buehler (1957) which was originally developed for obtaining confidence intervals for the product of two binomial parameters. Jobe and David (1992) have successfully extended Buehler's method to obtain upper confidence bounds in reliability/maintainability problems where the distributions are exponential.

Closing Remarks

Three approaches presented use bilinear constraints for gross error identification and estimation. Each technique has merits and limitations of its own. The two-stage and linearization approaches can be used now. However, we will be conducting simulation studies in the near future to better understand their best applications. In contrast, the product of normals approach will require distribution and test statistic development. This work is also planned for the near future.

Notation

A = $q \times p$ matrix of constants with negative entries for output streams and positive entries for input streams
 A_{11} = $q \times q$ partitioned matrix of A
 A_{12} = $q \times (p - q)$ partitioned matrix of the remaining elements of A
 A = $2q \times 2p$ matrix consisting of 0 's and A 's
 b = $2q \times 1$ vector of flow and temperature variables modeled with bias
 b_i = i th element of b
 B = $2p \times 2p$ matrix of flow and temperature values modeled without bias
 C = $2q \times 2p$ matrix defined by Eq. 34
 C_{11} = $2q \times 2q$ partitioned matrix of C
 C_{12} = $2q \times (2p - 2q)$ partitioned matrix of the remaining elements of C
 F = $p \times 1$ random vector of flow measurements
 I = identity matrix
 l = vector of zero's and one's used for making linear combinations of measurements
 M = $q \times p$ matrix defined by Eq. 26
 M_{11} = $q \times q$ partitioned matrix of M
 M_{12} = $q \times (p - q)$ partitioned matrix of the remaining elements of M
 n = number of samples
 N_p = multivariate normal p distribution
 N_q = multivariate normal q distribution
 p = number of measured variables
 q = number of nodes
 Q = $q \times 1$ vector of *known* constants representing all other significant forms of energy transfer; that is, heat transfer and work; other extensive forms of energy transfer are assumed to be negligible
 r = vector defined by Eq. 8
 s = vector defined by Eq. 12
 T = $p \times 1$ random vector of temperature measurements
 u = $2q \times 1$ vector consisting of r and s
 Z = random variable with a product of normals distribution

Greek letters

δ_b = $2p \times 1$ vector of biased measurements

$\hat{\delta}_b$ = estimate for δ_b
 δ_{b1} = $2q \times 1$ vector consisting of the first $2q$ elements in δ_b
 $\hat{\delta}_{b1}$ = estimate for δ_{b1}
 δ_{b2} = $(2p - 2q) \times 1$ vector of the remaining elements of δ_b
 $\hat{\delta}_{b2}$ = estimate for δ_{b2}
 δ_F = $p \times 1$ vector of flow measurements
 $\hat{\delta}_F$ = estimate for δ_F
 δ_{F1} = $q \times 1$ vector consisting of the first q elements in δ_F
 $\hat{\delta}_{F1}$ = estimate for δ_{F1}
 δ_{F2} = $(p - q) \times 1$ vector of the remaining elements of δ_F
 δ_T = $p \times 1$ vector of unknown biases for temperature measurements
 $\hat{\delta}_T$ = estimate for δ_T
 δ_{T1} = $q \times 1$ vector consisting of the first q elements in δ_T
 $\hat{\delta}_{T1}$ = estimate for δ_{T1}
 δ_{T2} = $(p - q) \times 1$ vector of the remaining elements of δ_T
 ϵ_F = $p \times 1$ vector of flow measurement errors
 ϵ_T = $p \times 1$ vector of temperature measurement errors
 μ_B = $2p \times 2p$ matrix of unknown true values of B
 μ_b = $2p \times 1$ vector of unknown true values of b
 μ_F = $p \times 1$ vector of unknown true flow rates
 $\hat{\mu}_F$ = estimate for μ_F
 μ_T = $p \times 1$ vector of unknown true temperatures
 μ_r = $q \times 1$ vector of unknown true values of r
 μ_s = $q \times 1$ vector of unknown true values of s
 Ω = $2q \times 1$ vector consisting of 0 's and Q
 Σ_F = $p \times p$ known variance-covariance matrix of ϵ_F
 Σ_T = $p \times p$ known variance-covariance matrix of ϵ_T
 Σ_r = $q \times q$ known variance-covariance matrix for r
 Σ_s = $p \times p$ known variance-covariance matrix for s

Other symbols

\mathcal{N} = distribution of the sum of products of normal random variables
 \mathcal{N}_q = multivariate q distribution of \mathcal{N}
 \sim = is distributed
 \approx = is approximately

Superscripts

T = transpose

Literature Cited

- Aroian, L. A., "The Probability of a Product of Two Normally Distributed Variables," *Ann. Math. Statist.*, **18**, 265 (1947).
 Buehler, R. J., "Confidence Bounds For The Product Of Two Binomial Parameters," *J. of Amer. Statist. Assoc.*, **52**, 482 (1957).
 Craig, C. C., "On the Frequency Function of XY ," *Ann. Math. Statist.*, **7**, 1 (1936).
 Jobe, J. M., and H. T. David, "Buehler Confidence Bounds For A Reliability-Maintainability Measure," *Technometrics*, accepted (1992).
 Meeker, W. M., Jr., L. W. Cornwell, and L. A. Aroian, "The Product of Two Normally Distributed Random Variables," *Selected Tables In Mathematical Statistics*, Vol. VII (1981).
 Rollins, D. K., "Unbiased Estimates of Measured Process Variables When Measurement Biases And Process Leaks Are Present," PhD Diss., Ohio State Univ., Columbus (1990).
 Rollins, D. K., and J. F. Davis, "Unbiased Estimation of Gross Errors In Process Measurements," *AIChE J.*, **38**, 563 (1992).

Manuscript received Feb. 3, 1992, and revision received June 8, 1992.