# CHAPTER II
# BACKGROUND AND LITERATURE SURVEY

## 2.1 Overview of Data Reconciliation (DR)

In a modern chemical plant, process measurements are used in a variety of activities such as planning, process control, optimization, and performance evaluation. The presence of random and nonrandom errors (gross errors) in "raw" measurement data, i.e. measurement data collected directly from plant instruments, leads to inaccurate process data, which do not even satisfy the steady state material and energy balances of the process. Such erroneous process data easily lead to deterioration in plant performance. The problem of improving the accuracy of process data so that they are consistent with material and energy balances of the system is known as data reconciliation. Process data after being treated by data reconciliation technique is called reconciled data. Simultaneously, there is also the problem of estimating unmeasured process variables, which is known as coaptation.

Data reconciliation is the technique to improve the accuracy of process data by making use of process constraints (typically material and energy balances).

## 2.2 Formulation of Data Reconciliation Problem

The essence of data reconciliation is that given the process measurements $y$ from the plant, we want to estimate the process state $x$, which satisfies the process constraints. We denote these estimated (or reconciled) values of process data as $\hat{x}$. In general, at steady state and in the absence of gross errors, the model for the measurements can be described as:

$$y = x + \varepsilon \tag{2.1}$$

Where $y$ is a vector of $n$ measurements, $x$ is the corresponding vector of the true value of the measured variables and $\varepsilon$ is the vector of unknown random errors.

Now we want to obtain $\hat{x}$ : vector of estimated values of measured variables (the estimates). In a probabilistic framework, this problem can be approached with

the laws of probability and maximum likelihood principle by maximizing the probability function of the difference between the measurements $y$ and the estimations $\hat{x}$, subject to the constraints. In fact, data reconciliation problem is an optimization problem. If random errors are assumed to follow a multivariate normal distribution with zero mean and known variance (which is usually the case), then we have the conventional weighted least-squares objective function, which is the most widely used form of objective function for data reconciliation problem, as follows:

$$\underset{x}{\text{Min}}(y-x)^\mathrm{T}S^{-1}(y-x) \tag{2.2}$$

Subject to process constraints.

When random errors are assumed to follow a multivariate normal distribution, measurements $y$ also follow multivariate normal distribution with the expected value $E[y] = E[x] + E[\varepsilon] = x$ and the variance given by matrix $S$. $S$ is the variance-covariance matrix of measurements, and usually is a diagonal matrix: $S = \text{diag}$ $(\alpha_1, \alpha_2 ... \alpha_n)$ with $\alpha_i$: measurement noise standard deviation

Process constraints can be equality constraints $G(x) = 0$ (typically material and energy balances) or both equality and inequality constraints in which inequality constraints are usually bounded limitations on values of process variables: lower limits $\leq x \leq$ higher limits (for example, $0 \leq$ molar fraction $\leq 1$) or physical/thermodynamics limitations such as temperature of hot stream $\geq$ temperature of cold stream (in heat exchangers). The most widely used process constraints are simple material balances. The reason is that we want to improve the accuracy of process data by making use of process constraints, so every parameter in process constraints model should be as accurate as possible since inaccurate parameters leads to inaccurate estimates of process variables and that makes the problem even worse.

## 2.3 Redundancy and Observability

Redundancy and observability of a measurement are inherently associated with data reconciliation problem since data reconciliation can not be conducted without redundancy of measurements. Besides, in process plants, there are hundreds

of variables and, for technical and economics reasons, it is not possible to measure all of them. It is thus important to know for the given process and a set of measured variable, which of the unmeasured variable can be estimated. The concept of observability deals with this issue.

Definition of observability and redundancy (Narasimhan and Jordache, 2000):

Observability: a variable is said to be observable if it can be estimated by using the measurements and steady-state process constraints.

Redundancy: a measured variable is said to be redundant if it is observable even when its measurement is removed.

From the above definition of observability, it is obvious that a measured variable is observable, since its measurement provides an estimate of the variable. However, an unmeasured variable is observable if it can be indirectly estimated by exploiting process constraint relationships and measurements in other variables. Measured variables are redundant if they can also be estimated indirectly through other measurements and constraints even when their measured values are eliminated.

## 2.4    Data Reconciliation In Linear Steady State System With All Variables Measured

Now if the system is linear and in the absence of inequality constraints, the process constraints model is:

$$Ax = 0 \qquad (2.3)$$

Where $A$ is $m \times n$ matrix containing parameters for process constraints model,

$m$: number of process constraint equations. $A$ is often called the constraint matrix.

The data reconciliation problem is given by (2.2) and the process constraints are given by (2.3). This case is the simplest case which is hardly found in real cases because the process constraints model contains only measured variables that are present in the objective function. The analytical solution to the above problem can be obtained using the method of Lagrange multipliers (Mah, 1990):

$$\hat{x} = y - SA^{T}(ASA^{T})^{-1}Ay \qquad (2.4)$$

or $\hat{x} = [I - SA^{T}(ASA^{T})^{-1}A]y = By \qquad (2.5)$

Equation 2.4 shows that the estimates are obtained using a linear transformation of the measurements. The estimates, therefore, are also normally distributed, with expected value and variance-covariance matrix given by (Narasimhan and Jordache, 2000):

$$E[\hat{x}] = BE(y) = Bx = x \tag{2.6}$$

$$Cov[\hat{x}] = E\{(By)(By)^T\} = BSB^T \tag{2.7}$$

Equation 2.6 implies that the estimates are unbiased, which is a property of maximum likelihood estimates for the linear systems. Equation 2.7 gives a measure of the accuracy of the estimates.

## 2.5 Data Reconciliation In Linear Steady State System With Both Measured And Unmeasured Variables

For partially measured system, the reconciliation problem is usually solved by decomposing it into two subproblems. In the first one, the redundant measured variables are reconciled, followed by the coaptation problem in which the observable unmeasured variables are estimated. To reconcile redundant measured variables, we need to derive constraints involving only measured variables that are present in the objective function. This can be accomplished in two ways. In graph theory, the reduced process graph obtained by pairwise aggregation of nodes linked by edges of unmeasured streams renders us constraints involving only measured variables (Mah, 1990). On the other hand, this can be done mathematically by making use of projection matrix developed by Crowe (1983). After we have obtained estimates of measured variables through data reconciliation, the observable unmeasured variables can then be estimated from measured variables through process constraints.

To handle this problem mathematically, we divide the variables into two sets: the vector $x$ of measured variables and the vector $u$ of unmeasured variables, and the process constraints are recast in terms of both the measured and unmeasured variables as follows:

$$A_x x + A_u u = 0 \tag{2.8}$$

Where $x$ is vector ($n \times 1$) of measured variables, $u$ is vector ($p \times 1$) of unmeasured variables, $A_x$ and $A_u$ are of dimensions $m \times n$ and $m \times p$, respectively.

We want to eliminate unmeasured variables from equation 2.8. This is done by premultiplying the constraints by a matrix $P$ called projection matrix. The matrix $P$ should satisfy the property:

$$PA_u = 0 \tag{2.9}$$

After premultiplying equation 2.8 by matrix $P$, we get the reduced set of constraints involving only measured variables as:

$$PA_x x = 0 \tag{2.10}$$

The number of constraints in the reduced set is known as degrees of redundancy.

The data reconciliation problem is also given by (2.2), the process constraints are given by (2.10). Similar to equation 2.4, the analytical solution to the above problem is given as (the matrix $A$ being replaced by the reduced matrix $PA_x$):

$$\hat{x} = y - S(PA_x)^T[(PA_x)S(PA_x)^T]^{-1}(PA_x)y \tag{2.11}$$

Substituting $\hat{x}$ into equation 2.8, the estimates of unmeasured variables $\hat{u}$ are then obtained. If all the unmeasured variables are observable (or the columns of $A_u$ are linearly independent), then unique estimates for the unmeasured variables $u$ exist and are obtained by least-square approximation solution (Narasimhan and Jordache, 2000):

$$\hat{u} = -(A_u^T A_u)^{-1}(A_x \hat{x}) \tag{2.12}$$

There are many ways to find the projection matrix $P$ and the most efficient one is the QR factorization of the matrix $A_u$. Properties and construction of projection matrix $P$ were discussed by Narasimhan and Jordache (2000).

## 2.6 Importance of Gross Error Detection (GED)

Gross errors are systematic errors that can exist in measurements (measurement biases) and process model (process leaks). Measurement bias relates to malfunction of instruments caused by miscalibration, improper installation, instrument degradation, etc, and is the more prevalent form of gross error. Leaks in units (tanks, heat exchangers...) in the process system make material flow imbalanced. Even when only one gross error exists, it deteriorates accuracy of all measurements in the process system through "*smearing effect*" of data reconciliation.

The reason is that a large deviation from true value in one measurement (i.e. gross error) will cause a series of small "corrections" made to other measurements through data reconciliation treatment. Thus the presence of gross errors can give rise to bad data, and invalidate the statistical basis of data reconciliation. Therefore it is crucial that gross errors are detected, identified and eliminated.

## 2.7 Hypothesis Testing For Gross Error Detection

There are various techniques for detecting and identifying gross errors, the most common used ones are based on statistical hypothesis testing, some are capable of detecting measurement bias only while some are capable of detecting both measurement bias and process leaks. The basis idea is to test the measured data against alternative hypotheses. They are the *null hypothesis* $H_0$, i.e. no gross error is present, and the *alternative hypothesis* $H_1$, i.e. one or more gross errors are present in the measurements. All statistical techniques for choosing between these two hypotheses make use of a *test statistic* which is the function of the measurements and constraint model. The test statistic is compared with a prespecified threshold value and the null hypothesis is rejected or accepted, respectively, depending on weather the statistic exceeds the threshold or not. The threshold value is also known as the test criterion or the critical value of the test.

The outcome of hypothesis testing is not perfect. A statistical test may declare the presence of gross errors, when in fact there is no gross error ($H_0$ is true). In this case, the test commits a *Type I error* or give rise to a false alarm. On the other hand, the test may declare the measurements to be free of error, when in fact one or more gross error exists (*type II error*). The *power* of a statistical test, which is the probability of correct detection, is equal to 1 − Type II error probability. The power and Type I error probability of any statistical test are intimately related. By allowing a larger Type I error probability, the power of a statistical test can be increased (more aggressive way) and vice versa (more conservative way). The test criterion can be selected so that the probability of Type I error is less than or equal to a specified value $\alpha$ which is also referred to as the *level of significance* for the statistical test.

We consider here the cases of *steady-state systems* and *linear constraint models*. The measurement model, equation 2.1, is modified to allow for the possible presence of gross errors:

$$y = x + \varepsilon + \delta \tag{2.13}$$

where $\delta$ is the gross error vector whose elements are the magnitudes of the gross errors.

The linear constraint model is given by:

$$Ax = c \tag{2.14}$$

Where $A$ is the linear constraint matrix and the vector $c$ contains known coefficients (typically $c$ is a zero vector).

The first two tests are based on the *vector of balance residuals*, $r$, which is given by:

$$r = Ay - c$$

As random errors $\varepsilon$ are assumed to follow normal distribution, the measurements $y$ and residuals $r$ are also normally distributed.

Under $H_0$, the expected value of $r$ is given by:

$$E[r] = E[Ay - c] = AE(y) - c = Ax - c = 0. \tag{2.15}$$

And the covariance matrix of $r$ is given by:

$$V = \text{cov}[r] = ASA^{\text{T}} \tag{2.16}$$

Where $S$ is the variance-covariance matrix of measurements $y$

It has been shown that if $\varepsilon$ is normally distributed; $r$ follows a t-variate normal distribution under $H_0$ where $t$ is the rank of matrix $A$ (Mah, 1990).

### 2.7.1 The Global Test

This test is capable of detecting but not identifying gross errors. It uses the test statistic given by:

$$\gamma = r^T V^{-1} r \tag{2.17}$$

It has been pointed out that the global test statistic is actually the optimal data reconciliation objective function (Narasimhan and Jordache, 2000). It follows chi-square distribution ($\chi^2$-distribution) with t degrees of freedom under $H_0$. If the test criterion is chosen as $\chi^2_{1-\alpha, t}$ (recall that $\alpha$ is the level of significance, $t$ is the rank of matrix $A$) then $H_0$ is rejected and a gross error is detected if $\gamma \geq \chi^2_{1-\alpha, t}$.

The chief drawback of the global test is that the test statistic applies to the whole process flowsheet. Once the presence of gross errors is detected, a separate procedure is required to identify them.

### 2.7.2 The Constraint Or Nodal Test

The test statistics are given by:

$$z_{r,i} = \frac{|r_i|}{\sqrt{V_{ii}}} \qquad i = 1,2,...m \qquad (m \text{ is number of constraints}) \qquad (2.18)$$

They follow a standard normal distribution $N(0,1)$ under $H_0$

The test statistic proposed by Crowe (1989), which is given by:

$$z_{r,i}^{*} = \frac{\left|[V^{-1}r]_i\right|}{\sqrt{[V^{-1}]_{ii}}} \qquad i = 1,2,...m \qquad (2.19)$$

is called the maximum power (MP) constraint test since it has the maximum power.

Unlike the global test, the constraint test processes each constraint residual separately and gives rise to $m$ univariate tests. Since multiple tests are performed using the same critical value, the probability of Type I error is supposed to be more than the specified value of $\alpha$. To control the type I error probability, the following modified level of significance $\beta$ was proposed (Mah and Tamhane, 1982):

$$\beta = 1 - (1-\alpha)^{1/m} \qquad (2.20)$$

Alternatively, Rollins and Davis (1992) proposed the use of a critical value based on the Bonferroni confidence interval given by:

$$\beta = \alpha/m \qquad (2.21)$$

The test criterion for all the constraint tests can be chosen as $Z_{1-\beta/2}$.

### 2.7.3 The Measurement Test

The third test is based on the vector of measurement adjustments:

$$a = y - \hat{x} = SA^T(ASA^T)^{-1}Ay \qquad (2.22)$$

where $\hat{x}$ are reconciled estimates of the process variables.

Under $H_0$, measurement adjustments $a$ follows a multivariate normal distribution $N(0, \bar{W})$, where:

$$\bar{W} = \text{cov}(\mathbf{a}) = SA^T(ASA^T)^{-1}AS \tag{2.23}$$

The measurement test statistics are given as:

$$z_{a,j} = \frac{|a_j|}{\sqrt{\bar{W}_{jj}}} \qquad j = 1,2...n \tag{2.24}$$

They follow a standard normal distribution $N(0,1)$ under $H_0$.

The maximum power (MP) measurement test proposed by Mah and Tamhane (1982) is obtained by premultiplying $\mathbf{a}$ by $S^{-1}$:

$$d = S^{-1}a \tag{2.25}$$

Under $H_0$, $d$ is also normally distributed with zero mean and a covariance matrix

$$W = \text{cov}(d) = A^T(ASA^T)^{-1}A \tag{2.26}$$

The test statistics given as:

$$z_{d,j} = \frac{|d_j|}{\sqrt{W_{jj}}} \qquad j = 1,2...n \tag{2.27}$$

have been shown to possess maximum power if S is a nondiagonal matrix.

Similar to the nodal test, the measurement test also involves multiple univariate tests. The type I error probability will be less than or equal to $\alpha$ if the test criterion is chosen as $Z_{1-\beta/2}$, where $\beta$ is given by equation 2.20 or 2.21 with $m$ being replaced by $n$, the number of univariate measurement tests.

Among the three tests already shown, the measurement test can directly identify location of the gross error, but measurement bias only. The measurement $y_j$ that corresponds to the test statistic $z_{a,j}$ or $z_{d,j}$ that exceeds the critical value is suspected of containing bias. However, the measurement test requires data reconciliation first. Whereas the global test does not require data reconciliation but it requires additional identification. The nodal test also does not require data reconciliation and identification problem is easier than global test, moreover, it can detect both measurement bias and process leak. This is based on the principle: one or more measurements involving in the constraint (node) that fails the test should contain gross error or there should be leak in that node (unit). But it also has problem of error cancellations, i.e. gross errors in measurements compensate for one another in such a way that constraint equations are satisfied.

### 2.7.4 Other Tests

The three tests discussed above basically treat only measurement or sensor biases. Overcoming this shortcoming, the two tests presented here can detect different types of gross errors (i.e. measurement biases and process leaks).

The generalized likelihood ratio (GLR) proposed by Narasimhan and Mah (1987) is based on the maximum likelihood ratio principle used in test statistic. This test requires a model of the process in the presence of a gross error, also known as gross error model. The GLR approach provides a framework for identifying any types of gross errors that can be mathematically modeled.

The other test is principal component (PC) test proposed by Tong and Crowe (1995). The PC test was claimed to be able to detect more subtle gross errors and have greater power to correctly identify the variables in error than the first three tests. However, it has been shown that the PC tests do not significantly enhance the ability in gross errors identification. Furthermore, the PC tests involve intensive computations in calculating eigenvalues and eigenvectors. It was discussed elsewhere by Tong and Crowe (1995) or Jiang *et al.* (1999).

## 2.8 Equivalency Theory

This theory was recently presented by Bagajewicz and Jiang (1998). It basically states that two sets of gross errors are equivalent when they have the same effect in data reconciliation, that is, when simulating either one in a compensation model, leads to the same value of objective. Therefore, the equivalent sets of gross errors are theoretically indistinguishable. In other words, when a set of gross errors are identified, there exists an equal possibility that the true location of gross errors are in one of its equivalent sets. From the view of graph theory, equivalent sets exist when candidate streams/leaks form a loop in an augmented graph consisting of the original graph representing the flowsheet with the additional of environment node. For the case of measurement biases, they proved that if a set of $k$ variables forms a cycle of the process graph, then gross errors in any combination of $(k-1)$ measurements from this set can not be distinguished from any other such

combination. By applying this theory, any proposed set of gross errors candidates cannot form a loop. Otherwise the size of these gross errors is indeterminate.

## 2.9 Simultaneous Strategies For Data Reconciliation and Gross Error Detection

Hypothesis tests as shown above are useful for detecting and identifying single gross error. To identity multiple gross errors, we need tailored strategies to be used in combination with one of these tests. Usually these strategies perform data reconciliation and gross error detection at the same time. This means that these strategies can render us the types, locations and magnitudes of gross errors together with the reconciled estimates which are free of errors of the process variables.

There are three types of strategies that make use of one of hypothesis tests in their procedure. They are serial elimination strategy, serial compensation strategy and simultaneous or collective compensation strategies. Another approach for simultaneous data reconciliation and gross error detection is a pure mathematical approach. Like other techniques (DR, GED) discussed so far, all these strategies are based on the assumption that process system is at steady-state and random errors are normally distributed, and are basically developed for linear systems.

### 2.9.1 Serial Elimination Strategy

This strategy is useful in identifying gross errors caused by measurement biases only, because it replies on eliminating measurements suspected of containing a bias. At each stage of the serial procedure, a gross error is identified in one measurement (based on some criteria) and the corresponding measurement is eliminated and treated as unmeasured variable before proceeding to the next stage. The major advantage of serial elimination is that it does not require any prior knowledge about the existence or location of gross errors, but it has the drawback of losing redundancy and is not applicable to process leaks.

The Modified Iterative Measurement Test (MIMT) proposed by Serth and Heenan (1986) is a popular and effective strategy that belongs to this kind. It uses the measurement test to detect and identify gross error and incorporates

information on bounds of variables into criterion to identify measurement suspected of containing a bias. The simplified procedure for this strategy is shown below:

Step 1: Perform the initial data reconciliation problem and obtain vectors $\tilde{x}$, $a$ and $d$.

Step 2: Compute the measurement test statistic $z_{d_j}$ for all measured variables (set T)

Step 3: Find $z_{max}$ = Max $\{z_{d_j}\}$, if $z_{max} \leq$ critical value $z_c$, proceed to step 6. Otherwise, select the measurement corresponding to $z_{max}$ and temporarily add it to set C.

Step 4: Remove the measurements contained in set C from the original set of measured variables and treat them as unmeasured variables, perform data reconciliation again and obtain new values of vectors $\tilde{x}$, $a$ and $d$.

Step 5: Check whether the estimates $\tilde{x}$ of variables are within their bounds. If so, store the current solution and return to step 2. Otherwise, delete the last entry in set C and, replace it with the measured variable corresponding to the next largest value of $z_{d_j} > z_c$, and return to step 4. If $z_{d_j} \leq z_c$ for all remaining variables, delete the last entry in set C and proceed to step 6.

Step 6: the measurements belong to set C are suspected of containing gross errors. The reconciled values after gross errors have been eliminated are those obtained in step 4 of the last iteration.

The MIMT was later modified by Kim *et al.* (1997) to handle nonlinear systems. They called it the *modified MIMT using NLP* in which data reconciliation was performed using nonlinear programming techniques and showed that it handles nonlinear systems better than the original one, especially when the number of gross errors is increased.

### 2.9.2 Serial Compensation Strategy (SCS)

This strategy was proposed by Narasimhan and Mah (1987) in conjunction with the use of the GLR test. In this strategy, one measurement bias is identified at a time, then estimated, and mathematically removed, before attempting to identify another bias. This strategy can identify all types of gross errors and can keep redundancy during the procedure. However, Rollins and Davis (1992) pointed out that it has two drawbacks: (1) it can have a large probability of making a wrong conclusion for measured variables that are unbiased (i.e. high Type I error) when at

least one variable is biased; and (2) estimates for measurement biases can be inaccurate.

### 2.9.3 Simultaneous Or Collective Compensation Strategies

Compared with the above two strategies, collective compensation strategies have some advantages: they are applicable to all types of gross errors, can maintain redundancy during the procedure, and provide better estimates thanks to the collective estimation (Rollins and Davis, 1992). There are three strategies that are considered to be the most efficient ones: UBET, SICC and MSEGE (Bagajewicz, 2000).

#### 2.9.3.1   UBET (Unbiased Estimation Technique)

This strategy was proposed by Rollins and Davis (1992). It was developed based on the following measurement model that takes into account measurements biases and process leaks:

$$y = \mu + \varepsilon \tag{2.28}$$

such that: $A\mu = M\gamma$ (2.29)

where $\mu = x$ (vector of true values); $\delta$: vector of measurements biases, $\varepsilon$: errors, which are assumed to follow multivariate normal distribution with mean $\delta$ and covariance matrix $S$: $\varepsilon \sim N_p(\delta, S)$; $M = [m_1,...m_m]$ where $m_j$ is chosen differently depending on the nature of constraints. If only total flow balances are involved, $m_j$ is identical to $e_j$, i.e. unit vector with one in position $j$ and zero elsewhere.

$\gamma$: vector of magnitudes of process leaks; $A$: constraint matrix.

Under $H_0$: $\mu_r = E[r] = 0$;

Under $H_1$: $\mu_r = E[r] = AE[y] = A\mu + AE[\varepsilon] = M\gamma + A\delta$ (2.30)

Where $r = Ay = A\mu + A\varepsilon$ : vector of balance residuals.

Partitioning $A$, $M$, $\delta$, $\gamma$ and rewriting equation 2.30, we get:

$$\mu_r = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix} + \begin{bmatrix} M_{12} & 0 \\ 0 & M_{22} \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}$$

$$\mu_r = \begin{bmatrix} A_{11} & 0 \\ A_{21} & M_{22} \end{bmatrix} \begin{bmatrix} \delta_1 \\ \gamma_2 \end{bmatrix} + \begin{bmatrix} A_{12} & M_{11} \\ A_{22} & 0 \end{bmatrix} \begin{bmatrix} \delta_2 \\ \gamma_1 \end{bmatrix} \qquad (2.31)$$

$$\mu_r = B_1\theta_1 + B_2\theta_2$$

It has been pointed out that the maximum number of gross errors that can be identified is equal to the number of process constraints $m$. Moreover, signature vectors of these gross errors are linearly independent (i.e. these gross errors do not form a loop) in other that their magnitudes can be uniquely estimated. In the UBET method, in order to obtain unbiased estimations of gross errors, two assumptions were made. The first one was that initially there are $m$ gross errors, moreover, the types and locations of these gross errors are also specified: there are $u$ measurements biases (denoted by vector $\delta_1$) and $v$ leaks (denoted by vector $\gamma_2$). The second one was that rank of $B_1$ is equal to $m$, in order words; these candidate gross errors do not form a loop. Therefore vector $\theta_1$ of gross errors is a $m \times 1$ vector and vector $\theta_2$ (dimension $n \times 1$) is a zero vector. Then equation (2.31) reduces to:

$$\mu_r = B_1\theta_1 \qquad (2.32)$$

By introducing:

$$l_i^T = e_i^T B_1^{-1} \qquad (2.33)$$

We get:

$$l_i^T \mu_r = e_i^T B_1^{-1} B_1 \theta_1 = e_i^T \theta_1 = \theta_i \ (= \delta_i \text{ or } \gamma_i) \qquad (2.34)$$

Therefore $l_i^T r$ ($i = 1, ..., m$) are unbiased estimators of the components of $\delta$ and $\gamma$ contained in $\theta_1$, there comes the name unbiased estimation technique (UBET).

The primary application of equation 2.33 is estimations of gross errors after some identification procedure has identified the required $\theta_1$ (or equivalently the required $\theta_2$). Additionally, it is used to construct a hypothesis test called Bonferroni test which can be used to identify the required $\theta_2$, or to determine if estimates for gross errors $\theta$s are statistically significant.

The test is given as: reject $H_0$: $l^T \mu_r = 0$ in favor of $H_1$: $l^T \mu_r \neq 0$, if:

$$\frac{\left| l^T r \right|}{\sqrt{l^T V l}} \geq z_{\alpha/2k} \qquad (2.35)$$

where $k$ equals the number of Bonferroni confidence intervals, $V$: covariance matrix of balance residuals $r$, given by equation 2.16

After unbiased estimates of gross errors have been found using equation 2.33, the unbiased reconciled estimates for process variables (that have known distribution and satisfy physical constraints) can be determined as described by Rollins and Davis (1992). The modified version of this strategy (MUBET) has been presented by Bagajewicz *et al.* (1999) in which they addressed singularities and uncertainties of the original method in view of the equivalency theory.

### 2.9.3.2 MSEGE

Sanchez and Romagnoli (1994) proposed a combinational approach called simultaneous estimation of gross errors (SEGE) to pick candidate gross errors and use them in a compensation model based on the use of the global test. This strategy was developed for systems with all redundant variables. In this technique, a statistical test based on the vector of adjustments $a$ (given by equation 2.22) is selected to detect the presence of gross errors, which is actually the global test statistic: $a^T S^{-1} a$ (recall that it is also the optimal data reconciliation objective function).

The test is given as: reject $H_0$: $E[a] = 0$ in favor of $H_1$: $E[a] \neq 0$, if:

$$a^T S^{-1} a \geq \chi^2_{g, \alpha} \tag{2.36}$$

where $g$: rank of constraint matrix A.

In this procedure, equations (process constraints) are added one by one to the least-squares estimation problem of the measured variables $x$. After each addition, the objective function (obv $= a^T S^{-1} a$ ) of the least-squares estimation technique is calculated and compared with the critical value $\tau_c = \chi^2_{g, \alpha}$ to detect gross errors. For updating the test statistic, the following expressions are applied:

$$\Sigma_c^{new} = \Sigma_c^{old} - \Sigma_c^{old} B_i^T (B_i \Sigma_c^{old} B_i^T )^{-1} B_i \Sigma_c^{old} \tag{2.37}$$

$$\hat{x} = \Sigma_c^{new} S^{-1} y \tag{2.38}$$

$$obv = (y - \hat{x})^T S^{-1} (y - \hat{x}) \tag{2.39}$$

where $\Sigma_c^{new}$ and $\Sigma_c^{old}$ represent the covariance matrices of the measurement estimates $\hat{x}$ after and before equation addition, and $B_i$ stands for the added equation.

Sanchez *et al.* (1999) have modified this strategy (called MSEGE) to address singularities and uncertainties of the original method in view of the equivalency theory, and to handle systems with both measured and unmeasured variables.

### 2.9.3.3 SICC

The MSEGE strategy has been shown to be highly accurate, however, it is still not suitable for large system, because it requires intensive calculation. Bagajewicz and Jiang (1998) proposed another strategy called SICC which was originally developed for linear dynamic data reconciliation problem. The steady-state version of this strategy was later presented by Bagajewicz and Jiang (1999). This strategy has been shown to be comparably accurate and requires less calculation. It relies on the measurement test for gross error detection. It uses the MT to make a list of suspected gross errors and identifies from the list one gross error using a compensation model. This error is put in a list of confirmed gross errors. Next a new list of suspected gross errors is constructed, and the compensation model is run using the confirmed gross errors and a candidate gross error (from the new list) at a time to determine which should be added to the confirmed gross errors list. The procedure is repeated until no gross errors are detected. Leaks are identified using the equivalency theory.

### 2.9.4 Mathematical Approach

Tjoa and Biegler (1991) proposed a method in which they used an objective function which was constructed using maximum likelihood principles on a combined distribution function. The function takes into account contributions from random and gross errors known as contaminated Gaussian distribution. The tailored hybrid SQP algorithm was developed to solve the contaminated Normal (Gaussian) objective function. Soderstrom *et at.* (2001) combined the gross error detection and identification problem with the data reconciliation within a mixed integer optimization framework. The advantage of these two mathematical approaches is that they do not require iterative procedure. However, they are significantly more computationally intensive.

## 2.10 Concept of Software Accuracy

Accuracy of measurements (hardware accuracy) has been defined as the sum of precision (standard deviation) and bias. Unfortunately, this definition is of little use unless bias is independently assessed. This leads to a new definition of more practical use. Bagajewicz (2004a) has introduced the concept of software accuracy (to distinguish it with the existing definition of hardware accuracy), which is based on the notion that data reconciliation with some test statistics is used to detect biases. In other words, report on accuracy should be made in relation to the ability of a gross error detection technique to detect and eliminate gross errors and contingent on the number of gross errors present in the measurements.

### 2.10.1 Induced Bias

From equation 2.13 and equation 2.4, the expected value of the vector of estimators can be derived as follows:

$$E[\hat{x}] = x + \delta - SA^T(ASA^T)^{-1}A\delta \qquad (2.40)$$

Note that $E[\hat{x}] = x$ only when $\delta = 0$.

The difference between $E[\hat{x}]$ when gross errors are present and the true value $x$ is defined as induced bias:

$$\hat{\delta} = [I - SW]\delta \qquad (2.41)$$

where $W = A^T(ASA^T)^{-1}A$ (variance-covariance matrix of $d = S^{-1}a$)

Various hypothesis tests can be used to detect bias of a size that is above a certain threshold. Below such threshold, the bias goes undetected and smears all the estimators, including those of the variables for which the corresponding instruments have no bias, called induced bias. Causing induced biases, an undetected bias in one measurement will corrupt accuracy of all measurements in the system through *smearing effect* of data reconciliation. Note that induced bias $\hat{\delta}$ is always smaller than the actual measurement bias $\delta$. Therefore, the (software) accuracy of an estimator is defined as the sum of precision (standard deviation) of the estimator plus the maximum possible undetected induced bias in that variable due to sensor biases anywhere in the system, including the instrument measuring the variable itself:

$$\hat{a}_i = \hat{\sigma}_i + \delta_i^* \qquad (2.42)$$

where $\hat{a}_i, \hat{\sigma}_i, \delta_i^*$ are the accuracy, precision (square root of variance $S_{ii}$) and the maximum undetected induced bias of the estimator, respectively.

## 2.11 Economic Value of Precision In The Monitoring of Linear Systems

Performing a statistical analysis, Bagajewicz *et al.* (2003) were able to obtain expressions for assessing the economic value of precision. A formula was developed for such value based on the downside expected loss that occurs when an operator adjusts the throughput of a plant when the measurements or estimators obtained through data reconciliation suggest that the targeted production is met or surpassed. However, there is a finite probability that the measurement or estimator is above the target when in fact the real flow is below it, hence the expected financial loss calculation. The calculation procedure to obtain the formula is briefly described as follows:

Bagajewicz *et al.* (2003) argued that a typical refinery consists of several tank units that receive the crude, several processing units, and several tanks where products are stored, summarized in three blocks as in figure 2.
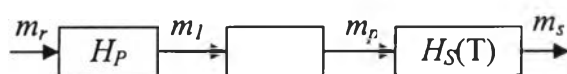


**Figure 2.1** Material balance in a refinery

In figure 2.1, H represents hold ups and $m$ flowrates.

They argued that the probability of not meeting the targeted production is $P\left\{H_S(T) \le H_s^*\right\}$, which in turn can be rewritten as $P\left\{m_p(t) \le m_p^*\right\}$, that is, it is equal to the probability of the true value of $m_p$ being smaller than the targeted value $m_p^*$. Let $\hat{m}_p$ be the estimate one has of the true value of $m_p$ and consider that production is adjusted to meet the targeted value, based on the estimate. In other words, if $\hat{m}_p < m_p^*$, production is increased and vice versa, if $\hat{m}_p > m_p^*$, production is

decreased. They assumed that, when $\hat{m}_p > m_p^*$, that is, the measurement indicates that the target has been met, the operator would not do any correction to the set points. They argued that the probability of being wrong is given by the conditional probability $P\{m_p \leq m_p^* | \hat{m}_p \geq m_p^*\}$, that is, the probability of having missed the target given that the estimator is larger than the target. Because these are independent, the above probability is equal to $P\{m_p \leq m_p^*\}P\{\hat{m}_p > m_p^*\}$. A statistical analysis was performed to derive the following expression:

$$P\left(m_p \leq m_p^* \middle| \hat{m}_p \geq m_p^*\right) = \int_{-\infty}^{m_p^*} \left\{ \int_{m_p^*}^{\infty} g_M(\xi, m_p, \hat{\sigma}_p)d\xi \right\} g_P(m_p, m_p^*, \sigma_p)dm_p \qquad (2.43)$$

where $g_P(m_p; m_p^*, \sigma_p)$ and $g_M(\xi; m_p, \hat{\sigma}_p)$ are probability distributions of the process values $m_p$ around the mean (targeted value) $m_p^*$ with variance $\sigma_p$ and of the estimators $\hat{m}_p$ around the true value $m_p$ with variance $\hat{\sigma}_p$, respectively.

For both distributions being normal, we obtain:

$$P\{m_p \leq m_p^* \mid \hat{m}_p \geq m_p^*\} = 0.5P\{\hat{m}_p \geq m_p^*\} = \frac{1}{4} + \frac{1}{2\sqrt{\pi}} \int_0^\infty erfc\left(z\sigma_p / \hat{\sigma}_p\right)e^{-z^2}dz \qquad (2.44)$$

The downside expected financial loss (loss incurred for not meeting the production target) was derived as follows:

$$DEFL(\hat{\sigma}_p, \sigma_p) = \int_{-\infty}^{m_p^*} g_P(m_p, m_p^*, \sigma_p)\left\{ K_s T \int_{-\infty}^{m_p^*} (m_p^* - \hat{m}_p)g_M(\hat{m}_p, m_p, \hat{\sigma}_p)d\hat{m}_p \right\}dm_p \qquad (2.45)$$

where $K_s$ is the value of the products sold and $T$ is the period of time under consideration.

When both distributions are normal, the downside expected financial loss (DEFL) is given by:

$$DEFL(\hat{\sigma}_p, \sigma_p) = \gamma K_s T \hat{\sigma}_p \left\{ \frac{1}{\sqrt{\left(\dfrac{\sigma_p}{\hat{\sigma}_p}\right)^2 + 1}} + \frac{\sigma_p}{\hat{\sigma}_p}\sqrt{\frac{1}{\left(1/\left(\dfrac{\sigma_p}{\hat{\sigma}_p}\right) + 1\right)}} \right\} \qquad (2.46)$$

Where $\gamma = 0.19947$.

The probability given by equation 2.43 is viewed as the confidence with which the expected loss given by equation 2.45 is known. Under simplified assumptions of negligible process variations (i.e., $\sigma_p / \hat{\sigma}_p \ll 1$) and normal distributions, we have: $P\{\hat{m}_p \geq m_p^* \mid m_p \leq m_p^*\} \rightarrow 0.25$ and $DEFL(\hat{\sigma}_p, \sigma_p) \rightarrow \gamma K_s T \hat{\sigma}_p$, in other words, under almost no process variations, there is a 25% chance that the downside expected financial loss of $\gamma K_s T \hat{\sigma}_p$ is achieved.

## 2.12 Economic Value of Accuracy In Linear Systems

While precision is important, most instruments present biases and therefore the theory of economic value of precision needs to be extended to include them. This gives rise to the theory of economic value of accuracy which was presented by Bagajewicz (2004b). This theory is briefly described as follows:

Consider the same problem (material balance in a refinery) as shown above. When there is a bias, induced or not, it could go undetected, which means it has an absolute value size smaller than $\hat{\delta}_{p,\max}^{i1,\ldots,i n_T}$, which is the maximum induced bias that goes undetected by the Maximum Power Measurement test when there are $n_T$ gross errors. As we know, this value is a function of the existing instrumentation. We therefore concentrate in redefining $g_M(\xi; m_p, \hat{\sigma}_p)$ to include the possibility of biases. Assuming one gross error ($\delta_i$) in variable $i$ and none in the others, we have:

$$g_M = g_M^{R,i}(\xi; m_p, \hat{\sigma}_{p,i}^R), \text{ with } |\hat{\delta}_p^i| > \hat{\delta}_{p,\max}^i$$
$$g_M = g_M(\xi; m_p + \hat{\delta}_p^i, \hat{\sigma}_p), \text{ with } |\hat{\delta}_p^i| \leq \hat{\delta}_{p,\max}^i \tag{2.47}$$

where $\hat{\sigma}_{p,i}^R$ is the residual precision left after the measurement of variable $i$ has been eliminated, $\hat{\delta}_p^i = \hat{\delta}_p^i(\delta_i)$ is the induced bias in the estimator of $m_p$.

Let us assume that, when an instrument fails, which happens according to a certain probability $f_i(t)$ (a function of time), the size of the bias follows a certain distribution $h_i(\theta; \bar{\delta}_i, \rho_i)$ with mean $\bar{\delta}_i$ and variance $\rho_i^2$. Note that depending on the value of the measurement in the range of the instrument, the mean could be nonzero.

For simplicity, we assume here that $\overline{\delta}_i = 0$. We are also assuming here that the gross error size distribution is independent of time. Thus, we now need to integrate over all possible values of the gross error and multiply by the probability of such bias to develop. Therefore, if we assume that one instrument fails at a time, then, the probability of instrument $i$ failing and the others not is given by: $\Phi_i^1 = f_i(t) \prod_{s \neq i} [1 - f_s(t)]$. Thus, the probability of the estimate to be higher than the targeted value (with the underlying assumption of the true value to be lower than the targeted value), given a bias in measurement $i$, is given by:

$$P\{\hat{m}_p \geq m_p^* | i\} = \Phi_i^1 \int_{-\infty}^{\infty} P\{\hat{m}_p \geq m_p^* | \theta\} h_i(\theta, \overline{\delta}_i, \rho_i) d\theta$$

$$= \Phi_i^1 \int_{-\infty}^{-\tilde{\delta}_i^P} \left[ \int_{-\infty}^{m_p} \left\{ \int_{-m_p}^{\infty} g_M(\xi, m_p, \hat{\sigma}_{p,i}^R) d\xi \right\} g_p(m_p, m_p^*, \sigma_p) dm_p \right] h_i(\theta, \overline{\delta}_i, \rho_i) d\theta$$

$$+ \Phi_i^1 \int_{-\tilde{\delta}_i^P}^{+\tilde{\delta}_i^P} \left[ \int_{-\infty}^{m_p} \left\{ \int_{-m_p}^{\infty} g_M(\xi, m_p + \hat{\delta}_p^i(\theta), \hat{\sigma}_p) d\xi \right\} g_p(m_p, m_p^*, \sigma_p) dm_p \right] h_i(\theta, \overline{\delta}_i, \rho_i) d\theta \qquad (2.48)$$

$$+ \Phi_i^1 \int_{+\tilde{\delta}_i^P}^{+\infty} \left[ \int_{-\infty}^{m_p} \left\{ \int_{-m_p}^{\infty} g_M(\xi, m_p, \hat{\sigma}_{p,i}^R) d\xi \right\} g_p(m_p, m_p^*, \sigma_p) dm_p \right] h_i(\theta, \overline{\delta}_i, \rho_i) d\theta$$

where $P\{\hat{m}_p \geq m_p^* | i\}$ indicates the probability being conditional to the presence of one gross error in stream $i$ and $\tilde{\delta}_i^P = \hat{\delta}_{p,max}^i$. The first and third term correspond to the detection of the gross error (when gross error magnitude is larger than the critical value), the second term corresponds to the presence of undetected gross error.

Since all the events are assumed independent, the probability of the estimate to be higher than the targeted value conditional to the presence of one gross error is:

$$P\left(\hat{m}_p \geq m_p^* | n_T = 1\right) = \sum_{\forall i} P\left(\hat{m}_p \geq m_p^* | i\right) \qquad (2.49)$$

The same procedure is applied to derive expressions for the cases of two and more gross errors. However, the set of critical values (limits of detectability) for a particular set of multiple gross errors is not unique. We therefore define the following function:

$$G\left(\delta_{i1}, \dots \delta_{ip} | \delta_{i1}, \dots \delta_{inb}\right) = \begin{cases} 1 & \text{if the MP test flag positive for } \delta_{i1}, \dots \delta_{ip} \\ 0 & \end{cases} \qquad (2.50)$$

One important assumption is made at this point: we assume consistency in gross error detection, i.e. it points to the correct locations of gross errors so that the right measurements are eliminated. This means that the set $\{\delta_{i1}, \ldots \delta_{ip}\}$ is a subset of $\{\delta_{i1}, \ldots \delta_{inb}\}$, that is, some errors are too small to be detected, but those detected are in fact true biases.

When two gross errors are present in the system, the following expression is derived:

$$P\left(\hat{m}_p \geq m_p^{\bullet} \mid i1, i2\right) = \Phi_{i1,i2}^2 \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} P\left(\hat{m}_p \geq m_p^{\bullet} \mid \theta_1, \theta_2\right) h_{i1}(\theta_1; \bar{\delta}_{i1}, \rho_{i}) h_{i2}(\theta_2; \bar{\delta}_{i2}, \rho_{i2}) d\theta_1 d\theta_2 \quad \text{and}$$

$$P\left(\hat{m}_p \geq m_p^{\bullet} \mid n_b = 2\right) = \sum_{\forall i,k} P\left(\hat{m}_p \geq m_p^{\bullet} \mid i, k\right) \quad \text{since all the events are assumed independent}$$

$$(2.51)$$

where:

$$P\left(\hat{m}_p \geq m_p^{\bullet} \mid \theta_1, \theta_2\right) = \int_{\infty}^{m_p^{\bullet}} \left| \begin{array}{l} G\left(\theta_1,\theta_2 \mid \theta_1,\theta_2\right) \int_{m_p}^{\infty} g_M(\xi, m_p, \hat{\sigma}_{p,i1,i2}^R) d\xi + \\[2mm] [1 - G\left(\theta_1,\theta_2 \mid \theta_1,\theta_2\right)] \int_{m_p}^{\infty} g_M(\xi, m_p + \hat{\delta}_p^{i1,i2}(\theta_1,\theta_2), \hat{\sigma}_{p,m}) d\xi \\[2mm] + G\left(\theta_1 \mid \theta_1,\theta_2\right) \int_{m_p}^{\infty} g_M(\xi, m_p + \hat{\delta}_{p,m}^{i2}(\theta_2), \hat{\sigma}_{p,i1}^R) d\xi \\[2mm] + G\left(\theta_2 \mid \theta_1,\theta_2\right) \int_{m_p}^{\infty} g_M(\xi, m_p + \hat{\delta}_{p,m}^{i1}(\theta_1), \hat{\sigma}_{p,i2}^R) d\xi \end{array} \right| g_p(m_p, m_p^{\bullet}, \sigma_p) dm_p \quad (2.52)$$

The first term inside the bracket in equation 2.52 corresponds to the detection of both gross errors, the second term: none of gross errors is detected, the third term & fourth term: detection of gross error in measurement i1 and in measurement i2, respectively. If detected, gross errors are eliminated by treating the corresponding measurements as unmeasured variables. Therefore when there is detection of gross errors, precision is replaced by the corresponding residual precision. When there are undetected gross errors, the estimate of the variable distributes around its true value plus the induced bias caused by the undetected gross errors $(m_p + \hat{\delta}_p^i(\theta))$ instead of its true value $(m_p)$ only. Note that we have the condition $G\left(\theta_1,\theta_2 \mid \theta_1,\theta_2\right) + G(\theta_1 \mid \theta_1,\theta_2) + G(\theta_2 \mid \theta_1,\theta_2) \leq 1$ holding naturally, that is, either the two gross errors flag, or only one of them flags, or none flags.

Expression for the probability of the estimate to be higher than the targeted value conditional to the presence of $n_T$ gross errors $P\left(\hat{m}_p \geq m_p^{\bullet} \mid n_T\right)$ is derived similarly:

$$P\left(\hat{m}_p \geq m_p^{\bullet}\Big|i1,...in_b\right) = \Phi_{i1...in_b}^{n_T} \times \int_{-\infty}^{+\infty}...\int_{-\infty}^{+\infty} P\left(\hat{m}_p \geq m_p^{\bullet}\Big|\theta_1,...\theta_{in_b}\right)h_{i1}(\theta_1;\bar{\delta}_{i1},\rho_{i1})...h_{in_b}(\theta_{in_b};\bar{\delta}_{in_b},\rho_{in_b})d\theta_1...d\theta_{in_b}$$

$$and \quad P\left(\hat{m}_p \geq m_p^{\bullet}\Big|n_T\right) = \sum_{\forall i1,i2,...in_b} P\left(\hat{m}_p \geq m_p^{\bullet}\Big|i1,i2,...in_b\right) \tag{2.53}$$

where $P\{\hat{m}_p \geq m_p^{\bullet}|\theta_1,..,\theta_{in_b}\}$ is the probability of the estimate being larger than the target in the presence of a set of particular $n_T$ gross errors and is determined at the same fashion as in equation 2.52 by making use of the function $G\left(\delta_{i1},...\delta_{ip}\Big|\delta_{i1},...\delta_{inb}\right)$;

$\Phi_{i1,i2,...in_b}^{n_T} = f_{i1}(t)...f_{in_b}(t) \prod_{s\neq i1,\ s\neq in_b}[1 - f_s(t)]$ is the probability of these gross errors being present.

Finally, since all events are mutually excusive, the probability of the estimator being larger than the target is the sum of all the possible cases:

$$P\left(\hat{m}_p \geq m_p^{\bullet}\right) = \sum_{r=0}^{n} P\left(\hat{m}_p \geq m_p^{\bullet}\Big|r\right) \tag{2.54}$$

$$P\left(\hat{m}_p \geq m_p^{\bullet}\right) = \Phi^0\left[\frac{1}{4} + \frac{1}{2\sqrt{\pi}}\int_0^{\infty} erfc(z\sigma_p/\hat{\sigma}_p)e^{-z^2}dz\right] + \sum_{r=1}^{n} P\left(\hat{m}_p \geq m_p^{\bullet}\Big|r\right)$$

The first term is the result obtained by Bagajewicz and Markowski (2003) for linear systems in the absence of biases.

### 2.12.1 Downside Expected Financial Loss

When there is one gross error present, incorporating equation 2.47 into equation 2.45, we obtain expression for downside expected financial loss as follows:

$$DEFL^1\Big|i = \int_{-\infty}^{-\delta_i^p}\left(\int_{-\infty}^{m_p^{\bullet}}\left\{K_sT\int_{-\infty}^{m_p^{\bullet}}(m_p^{\bullet} - \xi)g_M(\xi,m_p,\hat{\sigma}_{p,i}^R)d\xi\right\}g_p(m_p,m_p^{\bullet},\sigma_p)dm_p\right)h_i(\theta,\delta_i,\rho_i)d\theta$$

$$+ \int_{-\delta_i^p}^{\delta_i^p}\left(\int_{-\infty}^{m_p^{\bullet}}K_sT\left\{\int_{-\infty}^{m_p^{\bullet}}(m_p^{\bullet} - \xi)g_M(\xi,m_p + \hat{\delta}_p^i(\theta_i),\hat{\sigma}_p)d\xi\right\}g_p(m_p,m_p^{\bullet},\sigma_p)dm_p\right)h_i(\theta;\bar{\delta}_i,\rho_i)d\theta$$

$$+ \int_{\delta_i^p}^{\infty}\left(\int_{-\infty}^{m_p^{\bullet}}\left\{K_sT\int_{-\infty}^{m_p^{\bullet}}(m_p^{\bullet} - \xi)g_M(\xi,m_p,\hat{\sigma}_{p,i}^R)d\xi\right\}g_p(m_p,m_p^{\bullet},\sigma_p)dm_p\right)h_i(\theta;\bar{\delta}_i,\rho_i)d\theta \tag{2.55}$$

The meaning of the terms in equation 2.55 is same as in equation 2.48
When two gross errors are present, the financial loss is given by:

$$DEFL^2\big|_{i1,i2} = \int_{-\infty}^{m_p} K_s T \begin{vmatrix} \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} G(\theta_1,\theta_2|\theta_1,\theta_2)\left[\int_{-\infty}^{m_p}(m_p^\bullet - \xi)g_M(\xi,m_p,\hat{\sigma}_{p,i}^R)d\xi\right] \\ h_{i1}(\theta_1;\bar{\delta}_{i1},\rho_{i1})h_{i2}(\theta_2;\bar{\delta}_{i2},\rho_{i2})d\theta_1 d\theta_2 \\ + \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty}\left[\int_{-\infty}^{m_p}(m_p^\bullet - \xi)g_M(\xi,m_p + \hat{\delta}_p^{i1,i2}(\theta_1,\theta_2),\hat{\sigma}_{p,m})d\xi\right] \\ [1 - G(\theta_1,\theta_2|\theta_1,\theta_2)]h_{i1}(\theta_1;\bar{\delta}_{i1},\rho_{i1})h_{i2}(\theta_2;\bar{\delta}_{i2},\rho_{i2})d\theta_1 d\theta_2 \\ + \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty}\left[\int_{-\infty}^{m_p}(m_p^\bullet - \xi)g_M(\xi,m_p + \hat{\delta}_{p,m}^{i2}(\theta_2),\hat{\sigma}_{p,i1}^R)d\xi\right] \\ G(\theta_1|\theta_1,\theta_2)h_{i1}(\theta_1;\bar{\delta}_{i1},\rho_{i1})h_{i2}(\theta_2;\bar{\delta}_{i2},\rho_{i2})d\theta_1 d\theta_2 \\ + \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty}\left[\int_{-\infty}^{m_p}(m_p^\bullet - \xi)g_M(\xi,m_p + \hat{\delta}_{p,m}^{i1}(\theta_1),\hat{\sigma}_{p,i2}^R)d\xi\right] \\ G(\theta_2|\theta_1,\theta_2)h_{i1}(\theta_1;\bar{\delta}_{i1},\rho_{i1})h_{i2}(\theta_2;\bar{\delta}_{i2},\rho_{i2})d\theta_1 d\theta_2 \end{vmatrix} g_p(m_p,\dot{m}_p,\bar{\sigma}_p)d\dot{m}_p \quad (2.56)$$

The meaning of the terms inside the bracket in equation 2.56 is same as in equation 2.52. Downside expected financial loss for more than two gross errors present can be derived at the same fashion. Detail expressions for the financial loss and the associated probability when more than two gross errors are present in the system were given by Bagajewicz (2004b). He was also able to derive analytical form for the expressions in the presence of one bias, but did not provide analytical forms for the expressions for the presence of more than one bias because the integrals involved require integrating a discontinuous function that changes form in different regions. Finally, we write:

$$DEFL = \Psi^0 DEFL^0 + \sum_i \Psi_i^1 DEFL^1\big|i + \sum_{i1,i2} \Psi_{i1,i2}^2 . DEFL^2\big|i1,i2 + .... \quad (2.57)$$

where $\Psi^0$ are the average fraction of time the system is in the state without biases, $\Psi_i^1$ the average fraction of time the system has only one undetected bias only in stream $i$, etc. These values are in fact equal to the probabilities of each state.

### 2.12.2 Trade Off Between Value And Cost

In the case of buying a data reconciliation package, we have:

$$NPV = d_n \{\text{Change in } DEFL\} - Cost\ of\ license \quad (2.58)$$

where $d_n$ is the sum of discount factors for $n$ years; NPV: net present value. The change in DEFL is the economical benefit of this investment.

And in the case of adding a new instrumentation, similarly we have:

$NPV = d_n$ {Change in $DEFL$} $-Cost\ of\ new\ instrumentation$ (2.59)

## 2.13 Problem Statement

Financial loss calculation help engineers determine economical benefit of instrumentation upgrade investments such as buying a data reconciliation package or adding a new instrumentation by using (2.58) or (2.59). Knowing economical benefit of the investments, they can decide whether to implement the instrumentation upgrade investments or not. Moreover, because financial loss associated to the accuracy of measurement is the function of plant instrumentation, the financial loss calculation can be used in the problem of sensor network design or retrofit subjected to economical objectives such as minimizing financial loss or maximizing economical profit (net present value). However, as mentioned above, analytical forms for the expressions for the financial loss in the presence of more than one bias do not exist. This work aims at developing methods to calculate the financial loss and the associated probability when multiple gross errors are present in the system.

## 2.14 Literature Survey

The problem of data reconciliation was first introduced in 1961 and during the past four decades more than 200 research publications in the two areas of data reconciliation and gross error detection have appear. Steady-state data reconciliation and gross error detection are well-established techniques and many applications in chemical and mining industry have been reported, especially since the late 1980s when powerful computers were available. Currently, most integrated systems for process simulation, optimization and control include a data reconciliation (in couple with gross error detection) system, which precedes all applications that make uses of process data. The simplest industrial applications are for reconciliation of data around single units, especially for distillation or separation columns. In these applications, the flows and compositions of feed and products

streams are reconciled using overall material balances around the column. Applications of data reconciliation to chemical reactors have also been reported and were mentioned in the book of Narasimhan & Jordache (2000). When data reconciliation is used to process data for on-line optimization application, it is more appropriate to perform data reconciliation for a set of interconnected process units constituting a subsystem, for example, a subsystem comprising the crude distillation tower together with the crude preheat train of exchangers in a refinery. Pierucci *et al.* (1996) and Chiari *et at.* (1997) reported implementation of online reconciliation and optimization (ORO) package in olefin plants and in hydrogen & sulfur plants of Italian refineries. Christiansen *et al.* (1997) apply data reconciliation to evaluate performance of catalytic processes. Data reconciliation has also been applied in a vinyl acetate plant and a ketene plant (Dempf & List, 1998), an industrial utility plant (Lee *et al.*, 1998), a beverage alcohol distillation plant (Meyer *et al.*, 1993), to name a few of its successful implementations.

The problem of sensor network design has been explored by several authors with different approach based on different criteria: optimal parameter estimation, minimum cost, maximum overall precision, desired level of observability & monitoring, reliability, error detectablity, error robustness, or multicriteria. Many methods for designing sensor network are described by Narasimhan and Jordache (2000). The problem of upgrading instrumentation draws less attraction. Alhéritière *et al.* (1998) described a refinery case study of a method that can quantify contribution of process data to estimation of key variables and use this result to optimize the economic trade-off between enhanced parameter accuracy and increased measurement costs on the refinery but they offered few details. This method can help make a decision on upgrading/investment in instrumentation. Bagajewicz and Sánchez (2000) presented models to perform the upgrading of instrumentation at minimum cost to achieve maximum precision of selected parameters.

Data reconciliation undoubtedly helps improve performance of chemical plants. However, the problem of assessing quantitatively economical benefit of data reconciliation wasn't touched. This problem has been first tackled recently by Bagajewicz when he introduced the concept of software accuracy and the theory of economic value of precision and the theory of economic value of accuracy

(Bagajewicz, 2003, Bagajewicz *et al.*, 2004 and Bagajewicz, 2004b). These research works are the only works dealing with this problem that have appeared so far.

We notice that the integral expressions (for assessing the DEFL and the associated probability) involve probability distribution integrand functions. Therefore, these integral expressions can be readily evaluated using the Monte Carlo numerical integration method. The Monte Carlo numerical integration method was discussed by Evans and Swartz (2000). Some examples of the application of Monte Carlo numerical integration method for evaluating this kind of integrals were given by Mori and Kato, 2003 and Lu and Zhang, 2003.