

ELIMINATION IN LINEAR EDITING AND ERROR LOCALIZATION

Stanley S. Weng, National Agricultural Statistics Service, USDA
3251 Old Lee Hwy, Fairfax, VA 22030

KEY WORDS: Automatic Editing and Imputation, Fellegi-Holt Methodology, Linear Edit, Implied Edit, Elimination by Equality Edit

1. Introduction

For the error localization (EL) problem in automatic data editing and imputation (E/I) with linear edits under the Fellegi-Holt (F-H) methodology (Fellegi and Holt, 1976), the linear programming approach provides proper methods for solution (Rubin, 1975; Sande, 1978; Schiopu-Kratina and Kovar, 1989). However, in practice, the computational efficiency of error localization has been an issue (Winkler, 1999; Winkler and Chen, 2002). Various efforts have been made to improve the efficiency, including using an algorithm other than Chernikova's for linear programming, e.g., one based on Duffin's (1974) analysis of a system of linear inequalities (Houbiers, 1999); a tree-search approach instead of a Chernikova's algorithm-like process (Quere, 2000; Quere and De Waal, 2000); and even an entirely different approach, while still in the spirit of F-H (Bankier, 2000; Bankier, et al., 2000).

One other consideration is to simplify the linear edit system by using its special structure and features, for example, to reduce the dimensionality of the system and thus the magnitude of computation for error localization.

Edits used in economic surveys and censuses, like those created by USDA/NASS for the U.S. Census of Agriculture, are primarily linear. They also contain a considerable number of equality edits, for example, balance edits in which an aggregate variable is equal to the sum of its component variables.

In the presence of equality edits in a linear edit system, it seems preferable to use the equality edits to eliminate fields (variables), leading to a simplified system in reduced dimension. However, until now, none of the automatic computer E/I systems for numerical data have distinguished conceptually between equality and inequality edits. Equality edits have generally been treated as a special case of inequality edits. Some algorithms adopted the representation of an equality edit by two inequalities of opposite direction. Such handling seems to ignore the more informative specification of an equality edit. The equality form defines a more restrictive relationship than that of an inequality. In linear theory, an equality represents a lower dimension hyperplane in the data

linear space. The contribution of an equality edit to an editing problem should therefore be more than that of an inequality edit.

From the point of view of F-H methodology, there is an important distinction between equality and inequality edits in their generation of implied edits. This paper identifies such a distinction and establishes a method of using equality edits to eliminate fields and reach an equivalent linear edit system, for which all the inequality edits form a linear edit system of lower dimension. The original linear editing problem, for example error localization, can be solved by first solving the problem with respect to this reduced system, and then determining the remaining fields by the specification of the equality edits.

Benefits in computational efficiency from this methodology can be significant. The magnitude of the editing problem is reduced through elimination, and the program needs only to handle inequality edits.

The outline of this paper is as follows. Section 2 describes the basic setting and concepts of linear editing. Section 3 reviews some basic concepts and results of the F-H theory in the context of linear editing, that are related to the topic of this paper. Section 4 presents our theoretical results on the methodology of elimination by equality edit. Section 5 gets back to the main editing problem, error localization, which motivated this research and now can be solved in reduced scale with improved efficiency. Section 6 briefly discusses the implementation issue.

2. Linear Editing

The editing problem of numerical data from a survey/census is generally defined by a set of *linear edits* in the following form:

$$e_i : a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n \leq b_i, \quad (1a)$$

$$i = 1, 2, \dots, m$$

with *positivity* constraints for the variables x_j :

$$x_j \geq 0, \quad j = 1, 2, \dots, n. \quad (1b)$$

Here in (1a) the inequality sign may represent either inequality or equality. In matrix notation, the above linear edit system is written as

$$\mathbf{Ax} \leq \mathbf{b} \quad (2a)$$

and

$$\mathbf{x} \geq \mathbf{0} \quad (2b)$$

where \mathbf{A} ($m \times n$) is the edit coefficient matrix of (1a), \mathbf{b} ($m \times 1$) is the right-hand-side vector of (1a), and $\mathbf{x} = (x_1, x_2, \dots, x_n)^\tau$ is the data record vector (where τ denotes the transpose of a vector). Data editing so specified is called *linear editing*. Additional constraints may be added to the above basic setting to define various linear editing problems, for example error localization that will be discussed in Section 5.

A data record is a *passing* record with respect to a linear edit system if the record satisfies all edits in the system. Otherwise, the record is a *failed* one. All data points that satisfy the linear edit system constitute the *feasible area* of the system. A passing record is also called feasible, and a failed record infeasible. A linear edit system is completely described by its feasible area. Two linear edit systems are considered equivalent if they have identical feasible areas. Geometrically, the feasible area of a linear system is a polyhedron in the data space, which can be described by the set of all its *extremal points*.

We are actually in the setting of *linear programming* (Gass, 1985; Luenberger, 1984; Nemhauser and Wolsey, 1988). Linear editing problems, such as error localization, are generally related to solutions of a linear program. A linear program can be solved by finding the set of all extremal points of its feasible area. *Chernikova's algorithm* (Chernikova, 1964, 1965) is used to find all extremal points of a linear system of nonnegative variables.

3. F-H Theorem on Linear Edits

Fellegi and Holt (1976) established the fundamental theory of *automatic editing and imputation* in the following criteria, widely referred to as the F-H principles:

- (1) The data in each record should be made to satisfy all edits by changing the fewest possible items of data (fields).
- (2) Imputation rules should be derived from the corresponding edit rules without explicit specification.
- (3) When imputation takes place, it should maintain, as far as possible, the frequency structure of the data file.

For a failed record, identifying the fewest possible fields that may be changed to make the resulting record satisfy all edits is the *error localization* problem.

To solve the error localization problem, F-H showed that both explicit (the original) edits, as

specified by subject-matter experts, and implied edits are needed. An *implied edit* is one that is logically implied by a set of explicit edits. An implied edit is said to be an *essentially new* edit if it does not involve all the fields (variables) explicitly involved in the edits that generated it. A field that is eliminated in generating an essentially new implied edit is called a *generating field* of the implied edit. A set of edits together with all essentially new implied edits that can be generated from the set of edits, forms a *complete set of edits*. The concept of a complete set of edits is crucial in F-H theory, which underlies their main theorem.

We focus on linear editing. For linear edits, the generation of essentially new edits and the derivation of a complete set of edits take an explicit form, as given by Theorem 3 of Fellegi and Holt (1976). The following is a restatement of the theorem.

Theorem (F-H, 1976). An essentially new implied edit e_t is generated from edits e_r and e_s , as in (1a), using field j as a generating field, if and only if a_{rj} and a_{sj} are both nonzero and of opposite sign.

The coefficients of the new edit, a_{tk} , are given by

$$a_{tk} = a_{sk}a_{rj} - a_{rk}a_{sj}, k = 1, 2, \dots, n$$

where r and s are so chosen that $a_{rj} > 0$ and $a_{sj} < 0$.

Repeated application of the above procedure will derive all essentially new implied edits.

The theorem simply states that from two linear inequalities where the inequality signs are in the same direction, a variable can be eliminated by taking their linear combination if and only if the variable has coefficients in the two inequalities which are of the opposite sign. The essence of generating an essentially new implied edit is elimination of a field.

In linear theory, the method as used in F-H Theorem 3 to generate essentially new implied edits, is called *Fourier elimination* (Duffin, 1974; Fourier, 1826; Schrijver, 1986). This approach was proposed by Fourier to solve linear programming problems by elimination of variables. A variable, say, x_h , can be eliminated by taking positive combinations of two inequalities which have opposite signs in the coefficient of x_h . By adding suitable combinations of all possible pairs of inequalities with a positive and a negative coefficient of x_h , and subsequently adding all inequalities that did not contain x_h in the first place,

one gets a new system of inequalities which does not contain variable x_h . This process can continue in successive elimination of other variables.

In a Fourier elimination process, the number of inequalities can grow excessively. Moreover, by taking all possible linear combinations of the original inequalities during the elimination process, it could easily occur that some inequalities become redundant. That is, an inequality can be written as a positive linear combination of some of the other inequalities. Duffin (1974) in his method of analyzing systems of linear inequalities proposed a “refined elimination” rule which deletes any inequality which has been generated by adding $t + 2$ or more of the original inequalities, when t variables have been eliminated. Houbiers (1999) applied Duffin’s method to error localization.

Fourier’s original problem of interest was whether a feasible solution to a specified set of linear inequalities exists. This can be restated, in the terminology of modern automatic data editing, as whether a set of fields can be imputed in such a way that a specified set of linear edits can be satisfied. Fourier’s method of successive elimination has fostered modern automatic data editing, as generalized in the F-H methodology.

4. Elimination by Equality Edit

In addressing linear editing problems, it seems that the role of equality edits hasn’t been fully explored. Equality edits have generally been treated as a special case of inequality edits, without using the defining feature, the deterministic aspect, of an equality edit. Actually, from the implied edit point of view, there is an important distinction between equality edits and inequality edits in their generation of implied edits, as shown by the two lemmas to be introduced below.

Before stating the lemmas, we introduce the concept of equivalent edits. Two sets of edits are *equivalent*, if they imply each other, that is, each edit in one set is implied by (some edits of) the other set. In the linear edit context, two sets of linear edits are equivalent if their feasible area (thus, the set of extremal points) are identical. Two sets of equivalent linear edits have the same contribution to a linear edit system; they may thus replace each other. Editing problems with respect to two equivalent sets of edits are considered the same.

The following two lemmas extend the statements of Fellegi and Holt (1976) Theorem 3 in the situation where one edit is an equality. They state that, in such situations, it is always possible to generate an essentially new implied edit when a common field is involved. Furthermore, the original inequality edit can be replaced by the essentially new implied edit

generated.

Lemma 1. An essentially new implied edit can always be generated from edits e_r and e_s , where e_s is an equality edit, using field j as a generating field, provided the coefficients of field j in the two edits are both nonzero.

Proof. The lemma is clearly true. Since we can always make the coefficient of the generating field in the equality edit to be opposite in sign to that in the other edit, the lemma is thus an immediate consequence of Fellegi and Holt (1976) Theorem 3.

Lemma 2. An inequality edit e_r can be replaced by an essentially new implied edit e_t generated from e_r and an equality edit e_s .

Proof. The set of edits e_r and e_s is equivalent to the set of edits e_t and e_s , since edit e_r can also be generated as an implied edit by edits e_t and e_s . Thus we may use the set of edits e_t and e_s to replace the original set of edits e_r and e_s ; or, equivalently, use the essentially new implied edit e_t to replace the original inequality edit e_r .

The above lemmas show how an equality edit can be used to simplify a linear edit system. Based on these two lemmas, our next two theorems show that, just as elimination of free variables can be made using equalities in the linear system, so can elimination of positively constrained variables using the equality edits present in the linear edit system. The theorems are stated in the context of linear editing through the F-H concept of implied edit.

Theorem 1 (Elimination by equality edit). Suppose a linear edit system contains m inequality edits and one equality edit, with n positivity constraints for the n fields involved. Then, one nonzero field of the equality edit can be eliminated from all other edits involving that field. The resulting new linear edit system contains $m + 1$ inequality edits involving $n - 1$ fields, with $n - 1$ corresponding positivity constraints, and the original equality edit. The new system is equivalent to the original one. The extremal points of the original linear system can thus be obtained by first obtaining the extremal points in the $n - 1$ fields of the new linear system excluding the equality edit, and then determining the remaining field by the equality edit.

Proof of Theorem 1 can be found in Weng (2002), which provides the elimination method by repeated application of Lemma 1 and Lemma 2.

Theorem 1 can be extended to linear edit systems containing multiple equality edits, as follows.

Theorem 2. Suppose a linear edit system contains m inequality edits and q equality edits, with n positivity constraints for the n fields involved ($q \leq n$). Assume the q equality edits are of full rank. Then, a new linear edit system, which is equivalent to the original one, can be formed through elimination using the q equality edits. The new system contains $m + q$ inequality edits involving $n - q$ fields, with $n - q$ corresponding positivity constraints, and the original q equality edits. The extremal points of the original linear system can thus be obtained by first obtaining the extremal points in the $n - q$ fields of the new linear system excluding the q equality edits, and then determining the remaining q fields using the q equality edits.

Theorem 2 is established by repeated application of the elimination method of Theorem 1. A formal proof can be found in Weng (2002).

5. Error Localization

The error localization problem is stated as: for a failed record, anticipating the F-H principles, which components of the record must be changed in order that, with as few as possible changes, the record can be made to pass the edit system?

In linear editing, the linear programming approach to solving the error localization problem (Sante, 1978, 1979; Schiopu-Kratina and Kovar, 1989) has adopted Rubin's (1975) version of Chernikova's algorithm in the formulation of a *cardinality constrained linear program* problem, expressed as:

$$\begin{aligned} & \max \mathbf{d}^T \mathbf{x} \\ & \text{subject to} \\ & \mathbf{A} \mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}, \\ & |\mathbf{x}|^+ \leq \boldsymbol{\eta}, \end{aligned} \tag{9}$$

where \mathbf{x} and \mathbf{d} are $n \times 1$, \mathbf{A} is $m \times n$, \mathbf{b} is $m \times 1$, $|\mathbf{x}|^+$ denotes the *cardinality* ((the number of strictly positive elements of a nonnegative vector) of \mathbf{x} , and $\boldsymbol{\eta}$ is a

positive integer less than m in $\{m, n\}$. The linear programming directly produces the extremal points of the feasible area $G = \{\mathbf{x} | \mathbf{A} \mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ that satisfy $|\mathbf{x}|^+ \leq \boldsymbol{\eta}$, and then the optimal extremal point is determined. As Tanahashi and Luenberger (1971) showed, an optimal solution to (9) can always be found in G .

Implementation of such an approach has included GEIS (Sande, 1978; Schiopu-Kratina and Kovar, 1989) and CherryPi (De Waal, 1996).

Houbiers (1999) applied Duffin's method on Fourier's analysis of linear inequality systems to error localization. He compared Duffin's method with Chernikova's algorithm - two similar algorithms with different control rules for excessive growth of the matrix, and showed that Duffin's method is expected to be more efficient. Quere (2000) developed a new algorithm which performs Fourier elimination in a tree search process, instead of a Chernikova's algorithm-like process, to determine all optimal solutions to the error localization (see also Quere and De Waal, 2000).

In the presence of equality edits in the linear edit system, by the elimination methodology provided in last section, we can solve the error localization problem with respect to a simplified system in reduced dimension, as described below.

Through elimination by the equality edits, the linear edit system is restructured into the following form :

$$\begin{aligned} L_1 : \quad & \mathbf{A}_1 \mathbf{x}^{(1)} \leq \mathbf{b}^{(1)}, \\ & \mathbf{x}^{(1)} \geq \mathbf{0}, \end{aligned}$$

and

$$L_2 : \quad \mathbf{A}_2 \mathbf{x} = \mathbf{b}^{(2)},$$

where $\mathbf{x} = \begin{pmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{pmatrix} (n \times 1)$, $\mathbf{x}^{(1)} (n - q) \times 1$

consisting of the fields involved in the inequality edits in L_1 , $\mathbf{x}^{(2)} (q \times 1)$ consisting of the fields eliminated from the inequality edits; \mathbf{A}_1 $m_1 \times (n - q)$, $\mathbf{A}_2 (q \times n)$ of full rank, $\mathbf{b}^{(1)} (m_1 \times 1)$, and $\mathbf{b}^{(2)} (q \times 1)$.

Let $\mathbf{x}_0 = \begin{pmatrix} \mathbf{x}_0^{(1)} \\ \mathbf{x}_0^{(2)} \end{pmatrix}$ be a failed record. The correction

procedure is: if the subrecord $\mathbf{x}_0^{(1)}$ fails L_1 , perform error localization and imputation for $\mathbf{x}_0^{(1)}$ with respect to the system L_1 . And then correct $\mathbf{x}_0^{(2)}$ by the imputed $\mathbf{x}_0^{(1)}$ using the equality edits of L_2 . If $\mathbf{x}_0^{(1)}$ is feasible with respect to L_1 , but \mathbf{x}_0 fails L_2 , we only need to correct $\mathbf{x}_0^{(2)}$, again, by $\mathbf{x}_0^{(1)}$ using the equality edits of L_2 , a deterministic imputation.

Benefits in computational efficiency for error localization can be significant from application of the elimination methodology. In processing a row with Chernikova's algorithm, excessive growth of the number of columns depends on the number of fields, which causes the storage problem. Reduction of the number of fields reduces the magnitude of computation. Also, the computer code does not need to handle equality edits, which also simplifies the computation.

6. Implementation

In linear editing, elimination of fields by equality edits restructures the linear edit system. This restructuring is conducted prior to data editing, since data are not involved. A separate module can be created to perform the elimination.

Generally, when q (linearly independent) equality edits are present in the linear edit system, any subset of q fields may be selected for elimination from the inequality edits, provided the elimination process is valid according to Theorems 1 and 2. That is, the q variables are linearly independent. When performing a successive elimination, at each stage, there is no additional theoretical criterion for choosing a field for elimination, besides the general requirement of a nonzero field. Practically, some strategies may be developed for choosing the fields for elimination. At each stage of elimination, maximizing the number of zeros in the coefficients of inequality edits, appears a practical criterion. Aggregate variables seem natural candidates for elimination. Other strategies may be developed based on the structure of the edit system.

In computer implementation of the elimination process, either successive elimination or simultaneous elimination can be performed.

Acknowledgments: The author wishes to thank Dr. Ton de Waal, Statistics Netherlands, for his very helpful information and comments, provided in their correspondences, regarding the issues addressed in this paper. Thanks also to Dale Atkinson and other reviewers at USDA/NASS for their suggestions which

have helped to improve the presentation of this paper.

References

- Bankier, M. (2000), "Imputing Numeric and Qualitative Variables Simultaneously," Technical Report, Social Survey Methods Division, Statistics Canada.
- Bankier, M., Lachance, M., and Poirier, P. (2000), "2001 Canadian Census Minimum Change Donor Imputation Methodology - Extended Version of Report," Technical Report, Social Survey Methods Division, Statistics Canada.
- Chernikova, N.V. (1964), "Algorithm for Finding a General Formula for the Non-negative Solutions of a System of Linear Equations," *USSR Computational Mathematics and Mathematical Physics*, **4**, 151-158.
- (1965), "Algorithm for Finding a General Formula for the Non-negative Solutions of a System of Linear Inequalities," *USSR Computational Mathematics and Mathematical Physics*, **5**, 228-233.
- De Waal, T. (1996), "CherryPi: A Computer Program for Automatic Edit and Imputation," Statistics Netherlands, Voorburg.
- (2000a), New Developments in Automatic Edit and Imputation at Statistics Netherlands. Report, Statistics Netherlands, Voorburg.
- (2000b), "An Optimality Proof of Statistics Netherlands' New Algorithm for Automatic Editing of Mixed Data," Report, BPA 3295-00-RSM, Statistics Netherlands, Voorburg.
- (2002), Personal correspondences.
- Duffin, R.J. (1974), "On Fourier's Analysis of Linear Inequality Systems," *Mathematical Programming Studies*, Vol. I, 71-95, New York: North-Holland.
- Fellegi, I.P. and Holt, D. (1976), "A Systematic Approach to Automatic Edit and Imputation," *Journal of the American Statistical Association*, **71**, 17-35.
- Fourier, J.B.J. (1826), "Solution d'une Question Particuliere du Calcul des Inegalites," *Oeuvres II*, Paris.
- Houbiers, M.(1999), "Application of Duffin's Analysis of Linear Inequality Systems to the Error Localization Problem and Chernikova's Algorithm," Report, BPA 3107-99-RSM, Statistics Netherlands, Voorburg.
- Gass, S.I. (1985), *Linear Programming*, Fifth Edition. New York: McGraw-Hill.
- Luenberger, D.G. (1984), *Linear and Nonlinear Programming*, Second Edition. Reading, MA: Addison-Wesley.

- Nemhauser, G.L., and Wolsey, L.A. (1988), *Integer and Combinatorial Optimization*, New York: Wiley.
- Rubin, D.S. (1975), "Vertex Generation and Cardinality Constrained Linear Programs," *Operations Research*, **23**, 555-565.
- Quere, R. (2000), "Automatic Editing of Numerical Data," Report, BPA 2284-00-RSM, Statistics Netherlands, Voorburg.
- Quere, R., and De Waal, T. (2000), "Error Localization in Mixed Data Sets," Report, BPA 2285-00-RSM, Statistics Netherlands, Voorburg.
- Sande, G. (1978), "An Algorithm for the Fields to Impute Problems of Numerical and Coded Data," Technical Report, Statistics Canada.
- (1979), "Numerical Edit and Imputation," *Proceedings of the 42nd Session of the International Statistical Institute*, Manila, Philippines.
- Schiopu-Kratina, I., and Kovar, J.G. (1989), "Use of Chernikova's Algorithm in the Generalized Edit and Imputation System," Working Paper No. BSMD-89-001E, Statistics Canada. Ottawa, Ontario.
- Schrijver, A. (1986), *Theory of Linear and Integer Programming*, New York: John Wiley.
- Tanahashi, K. and Luenberger, D. (1971), "Cardinality-Constrained Linear Programming," Stanford University.
- Weng, S.S. (2002), "Elimination in Linear Editing and Error Localization," Research Report RDD-02-06, NASS, U.S. Department of Agriculture, Washington, DC.
- Winkler, W.E. (1999), "State of Statistical Data Editing and Current Research Problems," Working Paper No. 29, Conference of European Statisticians, Rome, Italy.
- Winkler, W.E. and Chen, B.-C. (2002), "Extending the Fellegi-Holt Model of Statistical Data Editing," Research Report, U.S. Bureau of the Census, Washington, D.C.