Purdue University

# Purdue e-Pubs

2001

# Cost Optimal Record/Entity Matching

V. S. Verykios

Ahmed K. Elmagarmid
*Purdue University*, ake@cs.purdue.edu

G. V. Moustakides

Report Number:
01-014

# COST OPTIMAL RECORD/ENTITY MATCHING

V.S. Verykios
A.K. Elmagarmid
G.V. Moustakides

# Cost Optimal Record/Entity Matching

V. S. Verykios[*]     A. K. Elmagarmid[†]     G. V. Moustakides[‡]

March 12, 2001

## Abstract

Record (or entity) matching or linkage is the process of identifying records in one or more data sources, that refer to the same real world entity or object. In record linkage, the ultimate goal of a decision model is to provide the decision maker with a tool for making decisions upon the actual matching status of a pair of records (i.e., documents, events, persons, cases, etc.). Existing models of record linkage rely on decision rules that minimize the probability of subjecting a case to clerical review, conditional on the probabilities of erroneous matches and erroneous non-matches. In practice though, (a) the value of an erroneous match is, in many applications, quite different from the value of an erroneous non-match, and (b) the cost and the probability of a misclassification, which is associated with the clerical review, is ignored in this way. In this paper, we present a decision model which is optimal, based on the cost of the record linkage operation, and general enough to accommodate multi-class or multi-decision case studies. We also present a closed form decision model for a class of multivariate record comparison pairs with binomially distributed components along with an example and results from applying the proposed model to large comparison spaces.

## 1 Introduction

In today's competitive business environment, corporations in the private sector are being driven to focus on their customers in order to maintain and expand their market share. This shift is resulting in customer data and information about customers being viewed as a corporate asset. In the public sector, the very large expansion of the role of the government resulted in an unprecedented increase in the demand for detailed information. Only recently has the data analytical value of these administrative records been fully realized. Of primary concern is that, unlike a purposeful data collection effort, the coding of the data is not carefully controlled for quality. Likewise, data objects are not necessarily defined commonly across databases nor in the way data consumers would want. Two of the serious concerns

---

[*]College of Information Science and Technology, Drexel University, USA.

[†]Computer Sciences Department, Purdue University, USA.

[‡]Computer Engineering and Informatics Department, University of Patras, Greece.

which arise in this context are (a) how to identify records across different data stores that refer to the same entity and (b) how to identify duplicate records within the same data store.

If each record in a database or a file carried a unique, universal and error-free identification code, the only problem would be to find an optimal search sequence that would minimize the total number of record comparisons. In most cases, encountered in practice, the identification code of the record is neither unique nor error-free. In some of these cases, the evidence presented by the identification codes (i.e., primary key, object id, etc.) may possibly point out that the records correspond or that they do not correspond to the same entity. However, in the large majority of practical problems, the evidence may not clearly point to one or the other of these two decisions. Thus, it becomes necessary to make a decision as to whether or not a given pair of records must be treated as though it corresponds to the same real world entity. This is called the record matching or linking problem [10, 1, 5, 13, 6, 3].

The large volume of applications spanning the range of cases from (a) an epidemiologist, who wishes to evaluate the effect of a new cancer treatment by matching information from a collection of medical case studies against a death registry in order to obtain information about the cause and the date of death, to (b) an economist, who wishes to evaluate energy policy decisions by matching a database containing fuel and commodity information for a set of companies against a database containing the values and the types of goods produced by the companies, signifies the tremendous impact and applicability of the problem addressed in this paper.

The remaining of this paper is organized as follows. Section 2 provides some background information, and the notation that is used throughout this paper. Section 3 introduces the cost optimal model, along with the thresholds of the decision areas, and the probabilities of errors. Section 4 provides a detailed formulation of the model when the comparison vector components are conditionally independent binomially distributed random variables. In Section 5, we describe an approach which can be used to estimate the parameters of the model, if a set of matched records is available. An example is given in Section 6 to illustrate how the model can be applied. Section 7 provides some information about the experimental environment that we generated and the results of some experiments that we run by using it. Finally, Section 8 provides concluding remarks and guidelines for future extensions of this work.

## 2  Background

Record matching or linking is the process of identifying records, in a data store, that refer to the same real world entity or object. There are two types of record matching. The first one is called *exact* or *deterministic* and it is primarily used when there are unique identifiers for each record. The other type of record matching is called *approximate*. In this paper, we focus only on the second type of matching. The decision, as to the matching status of a pair of records, is based on the comparison of common characteristics between the corresponding pair of records. These common characteristics are related to the similarities in the schema of the corresponding records. For example, a customer table can have two different schema representations in two databases, both storing customer data. The first table may store information from the service department while the second one may store

2

information from the billing department. Despite the differences in the representation of the two tables, overlapping information (i.e., name, address, sex, marital status, etc.) if present, can be used for the identification of matches between records from different databases that refer to the same customer.

The two principal steps in the record matching process are the searching step where we search for potential linkable pairs of records and the matching step where we decide whether or not a given pair is correctly matched. The aim of the searching step must be to reduce the possibility of failing to bring linkable records together for comparison. For the matching step, the problem is how to enable the computer to decide whether or not a pair of records relates to the same entity, when some of the identifying information agrees and some disagrees. In the remaining of this section we provide information about the notation that we will use, we discuss existing techniques which have been deployed for the record matching process and we also review decision models which have been built for the matching step.

## 2.1  Notation

In the product space of two tables, a *match M* is a pair that represents the same entity and a *non-match U* is a pair that represents two different entities. Within a single table, a *duplicate* is a record that represents the same entity as another record in the same database. Common record identifiers such as names, addresses and code numbers (SSN, object identifier), are the matching variables that are used to identify matches. The vector, that keeps the values of all the attribute comparisons for a pair of records (comparison pair) is called *comparison vector $\underline{x}$*. The set of all possible vectors, is called *comparison space X*. A record matching rule is a decision rule that designates a comparison pair either as a *link $A_1$*, a *possible link $A_2$*, or a *non-link $A_3$*, based on the information contained in the comparison vector. Possible links are those pairs for which there is no sufficient identifying information to determine whether a pair is a match, or a non-match. Typically, manual review is required in order to decide upon the matching status of possible links. *False matches* (Type I errors) are those non-matches that are erroneously designated as links by a decision rule. *False non-matches* (Type II errors) are either (a) matches designated as non-links by the decision rule, or (b) matches that are not in the set of pairs to which the decision rule is applied.

For an arbitrary comparison vector $\underline{x} \in X$, we denote by $P(\underline{x} \in X | M)$ or $f_M(\underline{x})$ the frequency of the occurrence or the conditional probability of the particular agreement $\underline{x}$ among the comparison pairs that are matches. Similarly, we denote by $P(\underline{x} \in X | U)$ or $f_U(\underline{x})$ the conditional probability of $\underline{x}$ among the non-matches. Note that the agreement or comparison vector $\underline{x}$ can be defined as specifically as one wishes and this completely rests to the components of the comparison vector. Let $p_j$ be the probability that the $j$-th corresponding item on the records $a$ and $b$ is present when the outcome of the comparison $(a, b)$ is a match, and let $p_j^*$ be similarly defined when the outcome is a non-match. Likewise, let $q_j$ be the probability that the $j$-th corresponding item on the records $a$ and $b$ is identical when the outcome of the comparison $(a, b)$ is a match and let $q_j^*$ be similarly defined when the outcome is a true non linkage. Let us also denote by $P(d = A_i, r = j)$ and $P(d = A_i | r = j)$ correspondingly, the joint and the conditional probability that the decision $A_i$ is taken, when the actual matching status ($M$ or $U$) is $j$. We also denote by $c_{ij}$ the cost of making a decision $A_i$ when the comparison record corresponds to some pair of records with actual matching

3

status $j$. When the dependence on the comparison vector is obvious by the context, we eliminate the symbol $\underline{x}$ from the probabilities. Finally we denote the a-priori probability of $M$ or else $P(r = M)$ as $\pi_0$ and the a-priori probability of $U$ or else $P(r = U)$ as $1 - \pi_0$.

## 2.2 Decision Models for Record Matching

In 1950s, Newcombe et. al. [15, 16, 17] introduced concepts of record matching that were formalized in the mathematical model of Fellegi and Sunter [4]. Newcombe recognized that linkage is a statistical problem: in the presence of errors of identifying information to decide which record pair of potential comparisons should be regarded as linked. Fellegi and Sunter formalized this intuitive recognition by defining a linkage rule as a partitioning of the comparison space into the so-called "linked" subset, a second subset for which the inference is that the record pairs refer to different underlying units and a complementary third set where the inference cannot be made without further evidence.

Fellegi and Sunter in [4], making rigorous concepts introduced by Newcombe et. al. [16] considered ratios of probabilities of the form:

$$R = P(\underline{x} \in X | M) / P(\underline{x} \in X | U) \tag{1}$$

where $\underline{x}$ is an arbitrary agreement pattern in the comparison space $X$. The theoretical decision rule is given by:

(a) If $R >$ UPPER, then designate pair as link.
(b) If LOWER $\leq R \leq$ UPPER, then designate the pair as a possible link and hold for clerical review.
(c) If $R <$ LOWER, then designate the pair as non-link.

The UPPER and LOWER cutoff thresholds are determined by a-priori error bounds on false matches and false non-matches. Fellegi and Sunter [4] showed that the decision rule is optimal in the sense that for any pair of fixed upper bounds on the rates of false matches and false non-matches, the manual/clerical review region is minimized over all decision rules on the same comparison space $X$. If now, one considers the costs of the various actions, that might be taken, and the utilities associated with their possible outcomes, it is desirable to choose decision rules that will minimize the costs of the operation. Nathan in [14] proposes a model that involves minimization of a cost function, but restricts detailed discussion to cases in which the information used for matching appears in precisely the same form, whenever the item exists in either input source. Du Bois's [18] approach attempts to maximize the set of correct matches by minimizing the set of erroneous non-matches. Tepping in [20] provides a graphical representation of a solution methodology that minimizes the mean value of the cost under the condition that the expected value of the loss is a linear function of the conditional probability that the comparison pair is a match. The application of his mathematical model involves the estimation of the cost function for each action, as a function of the probability of a match, and the estimation of the probability that a comparison pair is a match. Pinheiro and Sun [19] present a text similarity measure based on dynamic programming for matching

verbatim text fields. Based on the similarity measures for each corresponding pair of fields, they build a classification model using logistic regression to predict whether any two records are matched or not.

## 2.3 Intelligent Search of the Comparison Space

Errors, in the form of failures to bring potentially linkable pairs of records together for comparison, could be reduced to zero simply by comparing each record with all the others. However, wherever the files are large, such a procedure would generally be regarded as excessively costly, if there are many wasted comparisons of pairs of records that are not matched. For this reason, it is usual to order the records in the database by using identifying information that is common to all of them. The ordering can be performed either on the key, or on some other combination of record fields, or even on parts of the fields. In the exact matching, sorting of the file or the database can be used to reduce the complexity of identifying duplicate records [2]. In the approximate record matching, various compression codes, i.e., phonetic codes, can be used to mask some of the errors that frequently appear in typical record fields such as names. There is a number of systems to do this and the most common one is known as the Russel Soundex code [17].

Often, we need to make a compromise between the number of record pairs that are compared, and the recall of the linkage process. The searching process must be intelligent enough to exclude from comparison, record pairs that completely disagree with each other. In order to do that, the searching process must identify only those record pairs which have high probability of matching (prospective matches) and leave uninspected those pairs that look very different (not prospective matches). Several techniques have been developed in the past for searching the space of record pairs. The first one, was presented early on in a paper by Newcombe [17] and is called, *blocking*. In this approach, the database is scanned by comparing only those records that agree on a user-defined key, which for example can be the key used to sort the records. The characteristics used for blocking purposes are known as *blocking variables*. Kelley in [9] presents results related to a method for determining the best blocking strategy.

Another technique for cutting down the number of unwanted comparisons in the approximate record matching, is to scan the database by using a fixed size window and check for matches by comparing every pair of records that falls inside that window, assuming that the records are already sorted. This approach is known as the *sorted-neighborhood* approach and has been proposed by Hernadez and Stolfo in [6]. Because of the various types of errors that exist in the data sets that are compared, it is very common that the information selected, for blocking or sorting the data sets, contains errors. If that happens, we expect that some records are being clustered far way from those records with which they should be compared to. In this case, a *multi-pass approach*, proposed in [6], can be used. In this approach, a number of different blocking variables, or sorting keys, can be used for clustering the records in different ways. The database is then scanned as many times as the number of the different keys. The results from independent passes are combined to give the final set of matching records. An extension to the multi-pass approach has also been implemented by Hernadez and Stolfo [6] where the transitive closure of independent passes's results is computed. A similar approach, that has been proposed independently by Monge and Elkan in [12], makes

5

use of an algorithmic technique that identifies the connected components of a graph. By considering each record cluster as a connected component, this process can be effectively used to select the records that belong to the same cluster. Both groups of researchers presented very similar results, regarding the accuracy and the cost of the searching process.

# 3 The Cost Optimal Decision Model

Fellegi and Sunter [4] were the first that proposed a model for clustering the decision space into three decision areas, namely link, non-link and possible link. The authors considered the link and non-link decisions as positive dispositions. Their model minimized the probability of failing to make a positive disposition, under certain user-defined error bounds on these probabilities. Tepping [20] proposed that the matching problem should be regarded as a problem of decision making, subject to a utility function that depends upon the state of nature. He criticizes the Fellegi-Sunter model by arguing that the minimization of the probability of subjecting a case to clerical review, conditional on bounds on the probabilities of erroneous matches and erroneous non-matches ignores important facts:

- the value of an erroneous match is, in many (or perhaps in most) applications, quite different from the value of an erroneous non-match.

- the cost and the probability of a mis-classification which is associated with the manual review

We do not necessarily want to minimize the number of clerical reviews but to maximize the value of the record linkage operation. This implies that one must not only determine the costs of the various components of the operation, but must also set values on the possible outcomes. In this regard, Tepping refines the third decision area proposed by Fellegi and Sunter, as the one where some kind of further investigation is required before deciding on a classification. The investigation may simply involves the personal scrutiny or the search for additional data. This is exactly the case, when records are matched in an iterative fashion. For example, record pairs, created within the file blocks, are best subjected to a simple first test, prior to initiating the full comparison sequence or maybe an incremental sequence of comparisons. The utility function would specify a gain or loss for each of the possible decisions, conditional on whether the pair is a match or a non-match.

A possible set of actions that should be taken for a record comparison pair is presented below:

- Treat the comparison pair as if it designated to the same individual of some population. This is equivalent to the "link" decision.

- Temporarily treat the comparison pair as a link but obtain additional information before classifying the pair as a link or a non-link.

- Take no action immediately but obtain additional information before classifying the pair as a link or non-link.

6

- Temporarily treat the pair as if it was associated with different individuals of the population, but obtain additional information before classifying the pair as link or non-link.

- Treat the pair as if it was associated with different individuals in the population (non-link).

Other actions may be added to the list, including for example the use of a randomizing device, to determine the treatment of the comparison pair. Note that the underlying assumption is that each comparison pair is either a match or a non-match. Thus the set of all comparison pairs is the sum of mutually exclusive sets $M$ (the "match" pairs) and $U$ (the "non-match" pairs).

In order to be able to make a decision, we assume that the distributions of the random comparison vectors are known. Determining the distributions of the random vectors requires a pre-processing phase of training. For the training phase, a set of classified random vectors, which can be used for determining these a-priori matching probabilities, is required. Human experience can also play a major role in the training phase. The important thing to note here is that for this model to work, we need to know (a) the a-priori matching probabilities of the random comparison vectors, and (b) the various costs that should be assigned to different classifications/misclassifications. The model that we are building will determine the necessary and sufficient criterion for testing the $M$ hypothesis against the $U$ hypothesis and vice versa, and the thresholds required for this reason.

Below, we propose a new cost optimal decision model for record matching. The model presented here is a generalization of the model that it was proposed in [22] in the sense that the number of decision areas (link, non-link, possible link) is not restricted to three but it can be any non-negative number $n$. In general, let us denote by $c_{ij}$ the cost of making a decision $A_i$ for a comparison pair in the state of nature $j$. Each one of the decisions that are made, based on the existing evidence, about the linking status of a comparison pair, is associated with a certain cost that has two aspects. The first aspect is related to the decision process itself and is associated with the cost of making a particular decision; for example, the number of value comparisons that are needed in order to decide, affects the cost of this decision. The second aspect is associated with the cost of the impact of a certain decision; for example, making a wrong decision should always cost more than making the correct decision. Table 1 illustrates the costs for all the various decisions that could be made during the record matching process.

A record linkage process assigns each one of the comparison pairs to one and only one decision area. In order to compute the mean cost of the record linkage process, we consider one by one the costs of all decision areas. Without loss of generality, let us consider the cost of the $i$-th decision area. What we know about this area is that it has been assigned a number of comparison vectors based on a decision process that we are trying to identify. It is also the case, that among the comparison vectors allocated to this area, there maybe both matched and non-matched comparison pairs. There is a certain probability measure about the fact that a comparison pair (matched or non-matched) is allocated to this decision area. This is denoted by the joint probability $P(d = A_i, r = M)$ and $P(d = A_i, r = U)$ correspondingly. For every matched comparison pair assigned to the decision area $i$ the associated cost is $c_i^M$

| Cost | Decision | State of Nature |
|------|----------|-----------------|
| $c_1^M$ | $A_1$ | $M$ |
| $c_1^U$ | $A_1$ | $U$ |
| $c_2^M$ | $A_2$ | $M$ |
| $c_2^U$ | $A_2$ | $U$ |
| $\cdots$ | $\cdots$ | $\cdots$ |
| $c_n^M$ | $A_n$ | $M$ |
| $c_n^U$ | $A_n$ | $U$ |

Table 1: Costs of the decisions.

and for every non-matched comparison pair assigned to this area, the cost is $c_i^U$. The mean cost over all decision areas can then be written as follows:

$$\bar{c} = \sum_{i=1}^{n} [c_i^M \cdot P(d = A_i, r = M) + c_i^U \cdot P(d = A_i, r = U)] \qquad (2)$$

We can express the joint probabilities in Eq. 2 as a function of the conditional probabilities by using the Bayes theorem. Based on this observation, we get:

$$P(d = A_i, r = j) = P(d = A_i | r = j) \cdot P(r = j), \text{ where } i = 1, 2, \cdots, n \text{ and } j = M, U. \qquad (3)$$

Let us also assume that $\underline{x}$ is a comparison vector drawn randomly from the space of the comparison vectors which is shown in Figure 1. Then the following equality holds for the conditional probability $P(d = A_i | r = j)$:

$$P(d = A_i | r = j) = \sum_{\underline{x} \in A_i} f_j(\underline{x}), \text{ where } i = 1, 2, \cdots, n \text{ and } j = M, U. \qquad (4)$$

where $f_j$ is the probability density of the comparison vectors when the state of nature is $j$. We also denote the a-priori probability of $M$ or else $P(r = M)$ by $\pi_0$ and the a-priori probability of $U$ or else $P(r = U)$ as $1 - \pi_0$.
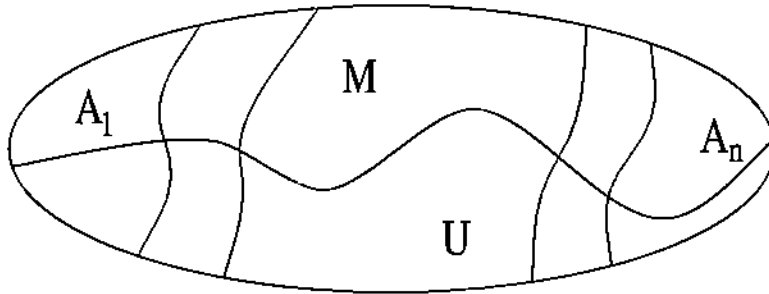


Figure 1: A partitioning of the decision space.

The mean cost $\bar{c}$ in Eq. 2 based on Eq. 3 is written as follows:

8

$$\bar{c} = \sum_{i=1}^{n} [c_i^M \cdot P(d = A_i | r = M) \cdot P(r = M) + c_i^U \cdot P(d = A_i | r = U) \cdot P(r = U)] \qquad (5)$$

By using Eq. 4, Eq. 5 becomes:

$$\bar{c} = \sum_{i=1}^{n} [c_i^M \cdot \sum_{\underline{x} \in A_i} f_M(\underline{x}) \cdot P(r = M) + c_i^U \cdot \sum_{\underline{x} \in A_i} f_U(\underline{x}) \cdot P(r = U)] \qquad (6)$$

By substituting the a-priori probabilities of $M$ and $U$ in Eq. 6, we get the following equation:

$$\bar{c} = \sum_{i=1}^{n} [c_i^M \cdot \pi_0 \cdot \sum_{\underline{x} \in A_i} f_M(\underline{x}) + c_i^U \cdot (1 - \pi_0) \cdot \sum_{\underline{x} \in A_i} f_U(\underline{x})] \qquad (7)$$

which by dropping the dependent vector variable $\underline{x}$, and combining the information for each part of the decision space, can be rewritten as follows:

$$\bar{c} = \sum_{i=1}^{n} \sum_{\underline{x} \in A_i} [f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_i^U \cdot (1 - \pi_0)] \qquad (8)$$

Every point $\underline{x}$ in the decision space $A$, belongs either in partition $A_1$, or in $A_2$, ..., or in $A_n$ and it contributes additively in the mean cost $\bar{c}$. We can thus assign each point independently either to $A_1$, or $A_2$, ..., or $A_n$ in such a way that its contribution to the mean cost is minimum. This will lead to the optimum selection for the sets which we denote by $A_1^o$, $A_2^o$, ..., and $A_n^o$. Based on this observation, a point $\underline{x}$ is assigned to the optimal decision area $A_i^o$ iff the following $n - 1$ inequalities hold:

$$
\begin{aligned}
f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_i^U \cdot (1 - \pi_0) &\leq f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_1^U \cdot (1 - \pi_0) \\
f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_i^U \cdot (1 - \pi_0) &\leq f_M \cdot c_2^M \cdot \pi_0 + f_U \cdot c_2^U \cdot (1 - \pi_0) \\
&\cdots \\
f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_i^U \cdot (1 - \pi_0) &\leq f_M \cdot c_n^M \cdot \pi_0 + f_U \cdot c_n^U \cdot (1 - \pi_0)
\end{aligned} \qquad (9)
$$

We thus conclude from the above that for any value of $i$, the corresponding decision area is given by the formula below:

$$A_i^0 = \{\underline{x} : \min_i (f_M \cdot c_i^M \cdot \pi_0 + f_U \cdot c_i^U \cdot (1 - \pi_0))\} \qquad (10)$$

In order for our model to define the decision areas, it makes use of $n$ systems of $n - 1$ linear inequalities. By solving for the likelihood ratio $f_M/f_U$ in each one of these systems:

$$f_M/f_U \leq (c_1^U - c_i^U)/(c_i^M - c_1^M) \cdot (1 - \pi_0)/\pi_0$$

$$\cdots$$

$$f_M/f_U \leq (c_{i-1}^U - c_i^U)/(c_i^M - c_{i-1}^M) \cdot (1 - \pi_0)/\pi_0$$

$$f_M/f_U \geq (c_{i+1}^U - c_i^U)/(c_i^M - c_{i+1}^M) \cdot (1 - \pi_0)/\pi_0$$

$$\cdots$$

$$f_M/f_U \geq (c_n^U - c_i^U)/(c_i^M - c_n^M) \cdot (1 - \pi_0)/\pi_0 \tag{11}$$

we get $n(n-1)$ values that the likelihood ratio should be compared with. These values denote the thresholds that explicitly define the decision areas. By inspecting these values closely, we observe that half of them are unique. Notice for example, that the last inequality in $A_1^o$ and the first in $A_n^o$ give raise to $f_M/f_U \geq (c_n^U - c_1^U)/(c_1^M - c_n^M) \cdot (1 - \pi_0)/\pi_0$ and $f_M/f_U \leq (c_1^U - c_n^U)/(c_n^M - c_1^M) \cdot (1 - \pi_0)/\pi_0$ correspondingly, where the thresholds are exactly the same. In general, the $n$ systems of $n-1$ equations generate $\binom{n}{2}$ unique thresholds. In order for all of the $n$ decision areas to exist, the following sufficient and necessary condition should hold for $n-1$ of these thresholds:

$$\frac{c_n^U - c_{n-1}^U}{c_{n-1}^M - c_n^M} \leq \frac{c_{n-1}^U - c_{n-2}^U}{c_{n-2}^M - c_{n-1}^M} \leq \cdots \leq \frac{c_4^U - c_3^U}{c_3^M - c_4^M} \leq \frac{c_3^U - c_2^U}{c_2^M - c_3^M} \leq \frac{c_2^U - c_1^U}{c_1^M - c_2^M} \tag{12}$$

Notice that for simplicity reasons, the ratio of prior probabilities have been eliminated from all the thresholds. For example, if for the likelihood ratio of a comparison vector the following inequality holds:

$$\frac{c_4^U - c_3^U}{c_3^M - c_4^M} \cdot \frac{1 - \pi_0}{\pi_0} \leq \frac{f_M}{f_U} \leq \frac{c_3^U - c_2^U}{c_2^M - c_3^M} \cdot \frac{1 - \pi_0}{\pi_0} \tag{13}$$

then the comparison vector belongs to $A_3^o$.

## 3.1 Optimality of the Decision Model

We can now prove that the decision model that we have proposed (i.e., the sets $A_1^o$, $A_2^o$, ..., $A_n^o$) is an optimal one. Based on the discussion above we know that $A = A_1 \bigcup A_2 \bigcup \cdots \bigcup A_n$, where $A_1$, $A_2$, $\cdots$, $A_n$ are pair-wise disjoint. Every point will be assigned to either one of these decision areas. We also introduce the indicator function $I_C$ of a set $C$, as the function which takes the value of 1 if the point $x$ belongs to $C$ and the value 0, otherwise. Note that we can formally write Eq. 8 as:

$$\bar{c} = \sum_{x \in A_1} z_1(x) + \sum_{x \in A_2} z_2(x) + \cdots + \sum_{x \in A_n} z_n(x) \tag{14}$$

where $z_i(x)$, $i = 1, 2, \ldots, n$ denote the expressions inside the corresponding sums in Eq. 8.

Using the indicator functions, we can write:

$$\bar{c} = \sum_{x \in A_1} z_1(x) + \sum_{x \in A_2} z_2(x) + \cdots + \sum_{x \in A_n} z_n(x) = \tag{15}$$

$$\sum_{x \in A} [z_1(x) \cdot I_{A_1}(x) + z_2(x) \cdot I_{A_2}(x) + \cdots + z_n(x) \cdot I_{A_n}(x)] \geq \tag{16}$$

10

$$\sum_{x \in A} \min\{z_1(x), z_2(x), \ldots, z_n(x)\} \overset{\text{def}}{=} \tag{17}$$

$$\sum_{x \in A_1^o} z_1(x) + \sum_{x \in A_2^o} z_2(x) + \cdots + \sum_{x \in A_n^o} z_n(x) \tag{18}$$

## 3.2 Error Estimation

The probability of errors can now be easily computed. There are two types of errors. The first one is called Type I error, and it occurs when a *non-link* action is taken although the two records are actually matched. The probability of this error can be estimated as follows:

$$P(d = A_n, r = M) = P(d = A_3 | r = M) \cdot P(r = M) = \pi_0 \cdot \sum_{x \in A_n} f_M(x). \tag{19}$$

The second type of error is called Type II error and it occurs when the *link* action is taken although the pair of records is actually non-matched. The probability of this error can be estimated as follows:

$$P(d = A_1, r = U) = P(d = A_1 | r = U) \cdot P(r = U) = (1 - \pi_0) \cdot \sum_{x \in A_1} f_U(x). \tag{20}$$

By computing these two types errors, we assume that all the other areas, in between these two, are not considered as definite decisions, and for this reason, we can use points assigned to them in either kind of error before further investigation.

# 4 The Formulation of the Multivariate Record Linkage Problem

Formally, a record is a finite collection of identifiers or items that describes a member of a given population. The notation used in Table 2 denote records as $L$-item information vectors $a_i$ and $b_k$, where the components $a_{ij}$ and $b_{kj}$ denote the $j$-th item recorded on the records $a_i$ and $b_k$ respectively. The components of these vectors are items of identifying information such as last name, middle initial, date of birth, etc. When a record $a_i$ from source $A$ is compared with a record $b_k$ from source $B$, a comparison vector:

$$(a_i, b_k) = \{(a_{i1}, b_{k1}), (a_{i2}, b_{k2}), \cdots, (a_{iL}, b_{kL})\} \tag{21}$$

is generated and the outcome is either a linkage (common or matching records), or a non-linkage (non-matching records). Since it is not known whether the outcome of a comparison $(a, b)$ is a linkage or a non linkage, a decision should be made based on the outcome of the comparison to whether a linkage or a non-linkage has occurred, where $(a, b)$ denotes records $a$ and $b$ from sources $A$ and $B$ respectively.

Define the indicator variable:

$$I_{(a,b)} = \begin{cases} 1 & \text{if the outcome of the comparison } (a, b) \text{ is a true linkage,} \\ 0 & \text{otherwise.} \end{cases} \tag{22}$$

11

|  | Source A | Source B |
|---|---|---|
|  | $a_1 = (a_{11}, a_{12}, \cdots, a_{1L})$ | $b_1 = (b_{11}, b_{12}, \cdots, b_{1L})$ |
|  | $a_2 = (a_{21}, a_{22}, \cdots, a_{2L})$ | $b_2 = (b_{21}, b_{22}, \cdots, b_{2L})$ |
|  | $\vdots$ | $\vdots$ |
|  | $a_M = (a_{M1}, a_{M2}, \cdots, a_{ML})$ | $b_N = (b_{N1}, b_{N2}, \cdots, b_{NL})$ |
| Number in the file | M | N |
| Common in both files | D | D |
| Proportion Common | $p = D/M$ | $P = D/N$ |

Table 2: Notation used for records and data sources.

then the number of records common to sources $A$ and $B$ is

$$D = \sum_{i=1}^{M} \sum_{k=1}^{N} I_{(a_i, b_k)} \tag{23}$$

where $M$ and $N$ are the number of records in sources $A$ and $B$ respectively. Let $p_j$ be the probability that the $j$-th corresponding item on the records $a$ and $b$ is present when the outcome of the comparison $(a, b)$ is a true linkage, and let $p_j^*$ be similarly defined when the outcome is a true non-linkage. Likewise, let $q_j$ be the probability that the $j$-th corresponding item on the records $a$ and $b$ is identical when the outcome of the comparison $(a, b)$ is a true linkage and let $q_j^*$ be similarly defined when the outcome is a true non linkage.

The advantage of the above formulation of the multivariate record linkage problem is the fact that the $p_j$'s and $q_j$'s can be used as parameters of a probability distribution over pairs of records belonging to the category of true linkages and similarly for the $p_j^*$'s and $q_j^*$'s with respect to pairs of documents belonging to the category of true non-linkages. The above parameters can be estimated by the method of maximum likelihood when these parameters are unknown and also when samples of pairs of records known to belong to the category of true linkages and the category of true non-linkages are available. This simplifies the complex problem of developing a set of weights to attach to corresponding items according to their degree of uniqueness.

The objective of this investigation is to develop a procedure for classifying record $a$ from source $A$ and record $b$ from source $B$ into one of two categories, namely, the category $M$ of potential linkages and the category $U$ of potential non-linkages. One approach to this problem is to observe whether corresponding items on both records are present and identical at the same time. We can formalize this approach by defining two variables $X_j$ and $Y_j$ in the following manner:

$$X_j = \begin{cases} 1 & \text{if the } j\text{-th corresponding item on both records is present,} \\ 0 & \text{otherwise.} \end{cases} \tag{24}$$

$$Y_j = \begin{cases} 1 & \text{if the } j\text{-th corresponding item on both records is identical,} \\ 0 & \text{otherwise.} \end{cases} \tag{25}$$

then

12

$$X_j Y_j = \begin{cases} 1 & \text{if the } j\text{-th corresponding item on both records} \\ & \text{is present and identical,} \\ 0 & \text{otherwise.} \end{cases} \qquad (26)$$

which implies that $Y_j$ is observable when $X_j = 1$. We use $X_j Y_j$ in order to distinguish between when $y_j = 0$ and when $Y_j$ is not observable due to $X_j = 0$. This distinction is useful when a computer is used to evaluate the multivariate probability mass function of the random comparison vector. Consider the $2L$-dimensional random vector

$$X^* = (X_1, X_2, \cdots, X_L; X_1 Y_1, X_2 Y_2, \cdots, X_L Y_L) \qquad (27)$$

defined so that the outcome of any comparison $(a, b)$ is a point $X^*(a, b) = \underline{x}^*$ where:

$$\underline{x}^* = (x_1, x_2, \cdots, x_L; x_1 y_1, x_2 y_2, \cdots, x_L y_L) \qquad (28)$$

In order to give to the record linkage problem a mathematical structure let the random vectors $X' = (x_1, x_2, \cdots, X_L)$ and $Y' = (Y_1, Y_2, \cdots, Y_L)$ have independent components and suppose the $X_j$'s and $Y_j$'s are distributed as point binomials $b(1, p_j)$ and $b(1, q_j)$ on $M$ and $b(1, p_j^*)$ and $b(1, q_j^*)$ on $U$. If the components $X_1, X_2, \cdots, X_L$ and $X_1 Y_1, X_2 Y_2, \cdots, X_L Y_L$ of the random vector $X^*$ are mutually independent within each sequence, then the multivariate probability mass functions of the random vector $X^*$ under $M$ and $U$ is given by:

$$f_1 = f(x^*|M) = \prod_{j=1}^{L} p_j^{x_j} (1 - p_j)^{1 - x_j} q_j^{x_j y_j} (1 - q_j)^{(1 - y_j) x_j} \qquad (29)$$

and

$$f_2 = f(x^*|U) = \prod_{j=1}^{L} p_j^{*x_j} (1 - p_j^*)^{1 - x_j} q_j^{*x_j y_j} (1 - q_j^*)^{(1 - y_j) x_j} \qquad (30)$$

# 5  Parameter Estimation

If the parameters of $f_1$ and $f_2$ are known, a linkage rule can be constructed in the following manner. The relevant information for making a decision is summarized in the set of 0's and 1's contained in the vector $x^*$. Then we calculate $f_1$ and $f_2$ from Eq. 29 and Eq. 30 respectively. Given a table of classification costs (i.e., Table 1), we decide that the outcome of the comparison $(a, b)$ must be assigned to a certain category based on the range of thresholds where the ratio $f_1/f_2$ belongs to, as it is determined by Eq. 12.

In actual applications of record linkage, the parameters $f_1$, $f_2$ and the prior probabilities used in the decision process, are rarely known. Generally these parameters can be estimated if verified matched and non-matched records are available. In some applications though, it is difficult to verify whether a particular pair of records is a true linkage. In this study, we assume that pairs of records belonging to $M$, and $U$ can be verified by additional information available or manual inspection. Let the $L$-item records in the samples be $a_1, a_2, \cdots, a_m$ from source $A$, and $b_1, b_2, \cdots, b_n$ from source $B$. For the $mn$ pairs of records from the sample, observe the $2L$-dimensional random vectors $X_1^*, X_2^*, \cdots, X_{mn}^*$ where

13

$$X_u^* = (X_{u1}, X_{u2}, \cdots, X_{uL}; X_{u1}Y_{u1}, X_{u2}Y_{u2}, \cdots, X_{uL}Y_{uL}) \qquad (31)$$

and the random variables $X_{uj}$ and $X_{uj}Y_{uj}$ assume the values 0 and 1. Assuming that the pairs of records belonging to $M$ and $U$ can be verified, the $mn$ random vectors $X_1^*, X_2^*, \cdots, X_{mn}^*$ can be decomposed into random samples of size $d$ and $mn - d$ from $f_1$ and $f_2$ respectively, after verification.

Using the method of maximum likelihood, we have:

$$
\begin{aligned}
\hat{\pi}_0 &= \frac{d}{mn} \\
\hat{p}_j &= \sum_{u=1}^{d} \frac{x_{uj}}{d} \\
\hat{p}_j^* &= \sum_{u=d+1}^{mn} \frac{x_{uj}}{mn - d} \\
\hat{q}_j &= \frac{\sum_{u=1}^{d} x_{uj}y_{uj}}{\sum_{u=1}^{d} x_{uj}} \\
\hat{q}_j^* &= \frac{\sum_{u=d+1}^{mn} x_{uj}y_{uj}}{\sum_{u=d+1}^{mn} x_{uj}}
\end{aligned}
\qquad (32)
$$

as estimates of $p_j$, $p_j^*$, $q_j$ and $q_j^*$, where $x_{uj} = 0, 1$ and $x_{uj}y_{uj} = 0, 1$ are values of the random variables $X_{uj}$ and $X_{uj}Y_{uj}$. The linkage rule based on the above estimators follows by replacing the parameters of $f_1$ and $f_2$ and $\pi_0$ in Eq. 12 by their respective estimates from Eq. 32.

# 6  Application

The previously presented model will be demonstrated in a file maintenance application, where the source data are lists of subscribers of two large magazine publishers. Table 3 shows tentative unit costs developed by the staff of the publishers on the basis of consideration of the character of the actions and the consequences of these actions. For example, based on the contents of this table, the cost $c_2^M$ is $0.41. The actions listed are roughly the same as those given above as examples in the description of the model.

In order to delineate the decision areas, we need to start with the test given in Eq. 12. By using this test we can find out whether all the areas are well defined, and if so, which are these areas for each action. In this example, the number of actions, or else decision areas, is 5. So, intuitively, four thresholds (the four rightmost ones in Eq. 12) can be checked. By substituting the values of the costs from Table 3 in Eq. 12 we get:

$$0.232 \leq 7.2 \leq 1 \leq 11.902 \qquad (33)$$

It is obvious that in Eq. 33 not all of the thresholds are in the right order. This means that not all the areas (5 decision areas) are defined by these costs so in order to define

| Action | True Status | |
|---|---|---|
| | Match | Non-match |
| 1 | $0.00 | $6.01 |
| 2 | 0.41 | 1.13 |
| 3 | 0.77 | 0.77 |
| 4 | 0.82 | 0.41 |
| 5 | 2.59 | 0.00 |

Table 3: Tentative Unit Costs

$$r_{11} \quad \geq \frac{c_2^U - c_1^U}{c_1^M - c_2^M} = r_{12} \quad \geq \frac{c_3^U - c_1^U}{c_1^M - c_3^M} = r_{13} \quad \geq \frac{c_4^U - c_1^U}{c_1^M - c_4^M} = r_{14} \quad \geq \frac{c_5^U - c_1^U}{c_1^M - c_5^U} = r_{15}$$

$$\leq \frac{c_1^U - c_2^U}{c_2^M - c_1^M} \quad r_{22} \quad \geq \frac{c_3^U - c_2^U}{c_2^M - c_3^M} = r_{23} \quad \geq \frac{c_4^U - c_2^U}{c_2^M - c_4^M} = r_{24} \quad \geq \frac{c_5^U - c_2^U}{c_2^M - c_5^M} = r_{25}$$

$$\leq \frac{c_1^U - c_3^U}{c_3^M - c_1^M} \quad \leq \frac{c_2^U - c_3^U}{c_3^M - c_2^M} \quad r_{33} \quad \geq \frac{c_4^U - c_3^U}{c_3^M - c_4^M} = r_{34} \quad \geq \frac{c_5^U - c_3^U}{c_3^M - c_5^M} = r_{35}$$

$$\leq \frac{c_1^U - c_4^U}{c_4^M - c_1^M} \quad \leq \frac{c_2^U - c_4^U}{c_4^M - c_2^M} \quad \leq \frac{c_3^U - c_4^U}{c_4^M - c_3^M} \quad r_{44} \quad \geq \frac{c_5^U - c_4^U}{c_4^M - c_5^M} = r_{45}$$

$$\leq \frac{c_1^U - c_5^U}{c_5^M - c_1^M} \quad \leq \frac{c_2^U - c_5^U}{c_5^M - c_2^M} \quad \leq \frac{c_3^U - c_5^U}{c_5^M - c_3^M} \quad \leq \frac{c_4^U - c_5^U}{c_5^M - c_4^M} \quad r_{55}$$

Table 4: A 5-by-5 system of inequalities for the file maintenance application.

them, we then need to consider the initial detailed model and the corresponding systems of equations. The system of inequalities for this application is depicted in Table 4.

Notice that the unique thresholds for $f_M/f_U$ in Table 4 are the $r_{ij}$'s, since the thresholds in the lower diagonal system are the same as their diagonal images. Also notice that $r_{ii} = 0$. Observe that the following two systems of inequalities should hold in order for all of the five areas to be well defined:

$$r_{12} \geq r_{13} \quad r_{12} \geq r_{14} \quad r_{12} \geq r_{15} \tag{34}$$
$$r_{23} \geq r_{24} \quad r_{23} \geq r_{25} \tag{35}$$
$$r_{34} \geq r_{35} \tag{36}$$

and

$$r_{35} \geq r_{45} \quad r_{25} \geq r_{45} \quad r_{15} \geq r_{45} \tag{37}$$
$$r_{24} \geq r_{34} \quad r_{14} \geq r_{34} \tag{38}$$
$$r_{13} \geq r_{23} \tag{39}$$

15

Notice, for example, that the system of inequalities in Eq. 34 holds because the threshold in the cell(2,1) in Table 4 (diagonal image $r_{12}$), needs to be the maximum threshold in the first row, otherwise there will be a gap between the first and the second decision areas.

By combining the inequalities above, we verify that $r_{12} \geq r_{23} \geq r_{34} \geq r_{45}$. In our case by substituting the values of the unit costs to the original system of inequalities, we get:

$$
\begin{array}{ccccc}
r_{11} & r_{12} = 11.09 & r_{13} = 6.805 & r_{14} = 6.82 & r_{15} = 2.32 \\
 & r_{22} & r_{23} = 1 & r_{24} = 1.75 & r_{25} = 0.51 \\
 & & r_{33} & r_{34} = 7.2 & r_{35} = 0.42 \\
 & & & r_{44} & r_{45} = 0.232 \\
 & & & & r_{55}
\end{array}
$$

By processing the information in the above system, we generate the decision areas:

$$
\text{Area} = \begin{cases}
1 & \text{if } f_M/f_U \geq 11.09 \\
2 & \text{if } 11.09 \geq f_M/f_U \geq 1.75 \\
4 & \text{if } 1.75 \geq f_M/f_U \geq 0.232 \\
5 & \text{if } 0.232 \geq f_M/f_U
\end{cases}
$$

Area 3 is not defined, since there is no region in the real axis in which $f_M/f_U \leq 1.75$ and at the same time $f_M/f_U \geq 7.2$. Notice also that the thresholds given above should be scaled by the ratio of prior probabilities $(1 - \pi_0)/\pi_0$. In the next section, we elaborate on this issue.

# 7 Experiments and Results

In order to validate and evaluate the proposed decision model, we have built an experimental evaluation system [21]. The evaluation system is built on top of a public domain system, the database generator [5], that automatically generates source data, with user-selected a-priori characteristics. The database generator allows us to perform controlled studies so as to establish the accuracy (or else the overall error), the percentage of comparison pairs which are assigned to the various decision areas and the overall cost of the record linkage process. The database generator provides a large number of parameters for selection such as the size of the generated database, the percentage of duplicate records in the database, and the percentage of the error in the duplicated records. Each one of the generated records, consists of the fields shown in Table 5. Some of the fields, as well, can be empty, affecting in this way the presence value. As it is reported in [5] the names were chosen randomly from a list of 63000 real names. The cities, the states and the zip codes (all from the USA) come from publicly available lists.

For each study, the evaluation system makes an external call to the database generator in order to generate two databases. The first database is used for training the decision model and the second database for testing the model. The training process includes the estimation of the required parameters by the decision model. Both databases are generated by using almost the same parameter settings. Only the number of records and the number of record clusters in each database can be different. A record cluster is a group of records in the same database that refers to the same person. All the records in the same cluster are considered as

duplicates. The training and the test databases are used correspondingly for generating the training comparison space and the test comparison space. As we mentioned earlier on, the comparison space is populated from comparison vectors which correspond to a component by component comparison of a pair of database records. In our system, we can explicitly select the type of the comparison, to be performed between each pair of values corresponding to the same attribute, and the type of the comparison result. In this study, the comparison vector has binary components and for this reason the result of a comparison can either be 0 or 1.

Some of the options that are provided to the users of the experimental system, for the generation of the training and test comparison spaces, include: (a) the pre-conditioning of the database records, (b) the selection of the sorting keys to be used for sorting the original database records, (c) the functions to be used for the comparison of each record attribute, (d) the searching strategy along with its parameters if applicable, and (e) the thresholds for the decision model. For the pre-conditioning of the database records, we may select to convert all the characters to uppercase or lowercase, and compute the Soundex code of the last name. Any subset or part of the record fields can be used as a sorting key. Among the functions to be selected for comparing pairs of field values, the most frequently used are the Hamming distance for numerical attributes, and the edit distance [11], the n-grams [7], the Jaro distance [8], and the Smith-Waterman algorithm [13] for character string attributes. For the searching strategy, the experimental system currently supports the blocking and the sorted-neighborhood approach. In the sorted-neighborhood approach the window size to be used should also be provided as an input parameter to the system. The last part of the parameters that are required by the system include the threshold values, which delimit the various decision areas in the proposed model.

In the set of experiments that we present, we make use of a comparison space of 10, 000 comparison records with known true matching status, as the training set, and a set of 1, 000, 000 records in the testing set. Notice that the size of the comparison space depends heavily on the searching technique used and is usually close to an order of magnitude larger than the number of actual database records compared. The estimated probabilities of presence and agreement are given in Table 5. These probabilities can be easily computed, based on the process that we describe in Section 5, by using the information in the training comparison space, since the actual matching status is considered known. This is possible, because each database record has been assigned a cluster identifier by the database generator, which is used for the identification of the cluster that each record belongs to.

The system also uses the costs of the various actions, in the decision process. Here, we make use of the costs presented in Table 3. In the experiments, we have generated pairs of training and testing record comparison sets with a variable size of cluster size. We have run many experiments in order to estimate the total cost of the linkage process for each testing comparison set by using a variable size of comparison fields from the Table 5. The results are shown in Table 6 and indicate that (a) the cost of the record linkage process decreases as the dimensionality of the comparison space increases and (b) for fixed dimensionality, there is no clear evidence whether the prior matching probability affects positively or negatively the total cost. The first observation is consistent with the intuition that more "reliable" comparison components help the model to minimize the total cost, while the second observation, is necessary so as our model to be independent of the data and so − to the degree this is

| Attribute | True Status | | | |
|---|---|---|---|---|
| | Match | | Non-match | |
| | $\hat{p}_j$ | $\hat{p}_j^*$ | $\hat{q}_j$ | $\hat{p}_j^*$ |
| SSN | 0.87 | 0.85 | 0.81 | 0.15 |
| First Name | 0.91 | 0.87 | 0.83 | 0.08 |
| Middle Initial | 0.76 | 0.64 | 0.93 | 0.05 |
| Last Name | 0.86 | 0.75 | 0.83 | 0.21 |
| Street Number | 0.90 | 0.57 | 0.81 | 0.10 |
| Street Address | 0.67 | 0.58 | 0.88 | 0.07 |
| Apartment Number | 0.45 | 0.47 | 0.89 | 0.05 |
| City | 0.56 | 0.59 | 0.91 | 0.12 |
| State | 0.78 | 0.81 | 0.86 | 0.16 |
| Zip Code | 0.89 | 0.91 | 0.92 | 0.06 |

Table 5: Estimated probabilities of presence and agreement in the training comparison space.

| $\hat{\pi}_0$ | Number of Vector Components | | | | |
|---|---|---|---|---|---|
| | 1 | 2 | 4 | 8 | 10 |
| 0.500 | \$222,700.000 | 113,359.125 | 18,790.750 | 6.300 | 1.900 |
| 0.250 | 200,206.250 | 109,679.312 | 19,582.665 | 9.038 | 2.700 |
| 0.200 | 173,105.000 | 101,050.850 | 18,540.950 | 9.544 | 2.880 |
| 0.125 | 125,140.625 | 80,936.280 | 15,839.407 | 10.230 | 3.090 |
| 0.100 | 109,152.500 | 74,231.425 | 14,563.803 | 10.420 | 3.150 |

Table 6: Total cost of record linkage for $1,000,000$ comparison records. The prior probability is estimated on a set of $10,000$ comparison records.

possible – unbiased. Other experiments performed, indicate that our model provides always the most cost efficient linkage.

# 8 Conclusions

This paper presents a new cost optimal decision model for the record matching process. The proposed model uses the ratio of the prior odds along with appropriate values of thresholds to partition the decision space to a number of decision areas. The model that we presented, is similar with the one proposed by Fellegi and Sunter [4] as it uses the same criterion for discriminating between matches and non-matches. The major difference between our model and all the other already existing models is that it minimizes the cost of making a decision rather than the probability of an erroneous decision. Our model is also much more efficient than other error-based models, as it does not resort to the sorting of the posterior odds in order to select the threshold values. The applicability of this model is independent of the

characteristics of the comparison fields, of the database fields, of the sorting techniques used and of the matching functions. This is strongly indicated in our work with the formulation of the multivariate record linkage problem for binomially distributed comparison fields. We should mention at this point that most of the researchers in this field are not considering presence independently of agreement and disagreement. We believe that this is a promising modeling approach that we plan to investigate further in the future.

In our future endeavors, we are also considering the design of a model for cost and time optimal record matching. By using such a model, it will be feasible not only to make a decision based on the entire comparison vector, but also to acquire as many comparison components as required, in order to make a certain decision. This will save computation time and at the same time it will facilitate the on-line decision making in the record matching context.

# References

[1] Wendy Alvey and Bettye Jamerson, *Record Linkage Techniques – 1997*, Proceedings of an International Workshop and Exposition, March 1997, Federal Committee on Statistical Methodology, Office of Management and Budget.

[2] Dina Bitton and David J. DeWitt, *Duplicate Record Elimination in Large Data Files*, ACM Transactions on Database Systems 8 (1983), no. 2, 255–265.

[3] Sunita Sarawagi (Ed.), *Special Issue on Data Cleaning*, IEEE Data Engineering Bulletin, December 2000.

[4] I. P. Fellegi and A. B. Sunter, *A Theory For Record Linkage*, Journal of the American Statistical Association 64 (1969), no. 328, 1183–1210.

[5] Mauricio Antonio Harnández-Sherrington, *A Generalization of Band Joins and the Merge/Purge Problem*, Ph.D. thesis, Department of Computer Sciences, Columbia University, 1996.

[6] Mauricio A. Hernadez and Salvatore J. Stolfo, *Real-World Data is Dirty: Data Cleansing and the Merge/Purge Problem*, Data Mining and Knowledge Discovery 2 (1998), no. 1, 9–37.

[7] Jeremy A. Hylton, *Identifying and Merging Related Bibliographic Records*, Master's thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, 1996.

[8] Matthew A. Jaro, *Advances in Record-Linkage Methodology as Applied to Matching the 1985 Census of Tampa, Florida*, Journal of the American Statistical Association 84 (1989), no. 406, 414–420.

[9] Patrick R. Kelley, *Advances in Record Linkage Methodology: A Method for Determining the Best Blocking Strategy*, Proceedings of the Workshop on Exact Matching Methodologies (1985), 199–203.

[10] Beth Kliss and Wendy Alvey, *Record Linkage Techniques – 1985*, Proceedings of the Workshop on Exact Matching Methodologies, May 1985, Department of the Treasury, Internal Revenue Service, Statistics Income Division.

[11] U. Manber, *Introduction to Algorithms*, Addison-Wesley Publishing Company, 1989.

[12] Alvaro E. Monge and Charles P. Elkan, *An Efficient Domain-Independent Algorithm for Detecting Approximately Duplicate Database Records*, Proceedings of the SIGMOD Workshop on Research Issues on DMKD, 1997, pp. 23–29.

[13] Alvaro Edmundo Monge, *Adaptive Detection of Approximately Duplicate Records and the Database Integration Approach to Information Discovery*, Ph.D. thesis, Department of Computer Science and Engineering, University of California, San Diego, 1997.

[14] Gad Nathan, *Outcome Probabilities for a Record Matching Process with Complete Invariant Information*, Journal of the American Statistical Association **62** (1967), no. 318, 454–469.

[15] H.B. Newcombe and J.M. Kennedy, *Record Linkage: Making Maximum Use of the Discriminating Power of Identifying Information*, Communications of the ACM **5** (1962), 563–566.

[16] H.B. Newcombe, J.M. Kennedy, S.J. Axford, and A.P. James, *Automatic Linkage of Vital Records*, Science **130** (1959), no. 3381, 954–959.

[17] Howard B. Newcombe, *Record Linking: The Design of Efficient Systems for Linking Records into Individual and Family Histories*, Americal Journal of Human Genetics **19** (1967), no. 3.

[18] Jr. N.S. D' Andrea Du Bois, *A Solution to the Problem of Linking Multivariate Documents*, Journal of the American Statistical Association **64** (1969), no. 325, 163–174.

[19] Jose C. Pinheiro and Don X. Sun, *Methods for Linking and Mining Heterogeneous Databases*, Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining, 1998, pp. 309–313.

[20] Benjamin J. Tepping, *A Model for Optimum Linkage of Records*, Journal of the American Statistical Association **63** (1968), 1321–1332.

[21] Vassilios S. Verykios, Mohamed G. Elfeky, Ahmed K. Elmagarmid, Munir Cochinwala, and Sid Dalal, *On the Accuracy and Completeness of the Record Matching Process*, 2000 Information Quality Conference (2000), 54–69.

[22] Vassilios S. Verykios and George V. Moustakides, *A Cost Optimal Decision Model for Record Matching*, Workshop on Data Quality: Challenges for Computer Science and Statistics (2001).