# Electre Tri-Machine Learning Approach to the Record Linkage Problem

Renato De Leone\* Valentina Minnetti<sup>†‡</sup>
May 26, 2015

#### Abstract

In this short paper, the Electre Tri-Machine Learning Method, generally used to solve ordinal classification problems, is proposed for solving the Record Linkage problem. Preliminary experimental results show that, using the Electre Tri method, high accuracy can be achieved and more than 99% of the matches and nonmatches were correctly identified by the procedure.

#### 1 Introduction

Machine Learning is a scientific discipline that is concerned with the design and development of algorithms that allow computers to "learn data". More precisely, "learn" is here intended as the possibility to automatically recognize complex patterns and make "intelligent" decisions, based on information data. Hence, machine learning is closely related to fields such as statistics, probability theory, data mining, pattern recognition, artificial intelligence, adaptive control and theoretical computer science.

Machine learning algorithms can be classified in the following types:

- supervised learning algorithms: a function/classifier is generated, that maps outputs on the training inputs, based on labeled examples inputoutput;
- unsupervised learning algorithms: patterns in the input are recognized, the examples have no labels;
- semi-supervised learning algorithms: supervised and unsupervised learning information is combined;

<sup>\*</sup>School of Science and Technology, University of Camerino, Italy renato.deleone@unicam.it 

†Faculty of Information Engineering, Informatics and Statistics, Sapienza University of Rome, Italy valentina.minnetti@uniroma1.it

<sup>&</sup>lt;sup>‡</sup>Italian National Institute of Statistics, Rome, Italy minnetti@istat.it

• reinforcement learning: actions from observation of the world are generated. Every action has some impact in the environment and the environment provides feedbacks that are translated into a score that guide the learning process.

The principal supervised learning techniques currently applied or under consideration at statistical agencies worldwide to solve the record linkage matching problem are: classification tree [4, 7], support vector machine [1, 2, 3] and neural network [15]. In this short paper, another machine learning technique is proposed to solve the record linkage problem: the multi-criteria classification method Electre Tri. It is the first time that multi-criteria machine learning technique is used to solve the record linkage problem.

This application answers to one of "many challenges in applying supervised machine learning to record linkage matching" [10], showing that the use of multi-criteria classification method Electre Tri to solve the record linkage problem provides good results in term of classification model performances. The importance of this application is in light of the increasing development of the use of administrative sources data. In this context, an important problem is that of finding matching pairs of records from heterogeneous databases, while maintaining privacy of the databases parties. To this purpose secure computation of distance metrics is important for secure record linkage [5].

The paper is organized as follows. Section 2 describes an introduction to the Record Linkage problem; then the next Section 3 describes the method Electre Tri, used to solved the Record Linkage and in the last Section 4 a preliminary experiment is conducted on simulated data. The paper closes with some final remarks and conclusions.

## 2 Linked Data: the Record Linkage

Generally speaking, in integration of two data sets the objective is the detection of those records, in the different data sets, that belong to the same statistical unit. This action allows the reconstruction of a unique record of data that contains all the unit information collected from different data sources on that unit.

Therefore, record linkage is the methodology of bringing together corresponding records from two or more files or finding duplicates within files [16]. In the first situation, the definition of record linkage in [9] is more precise "Record linkage is a solution to the problem of recognizing those records in two files which represent identical persons, objects, or events (said to be matched)"

The term record linkage originated in the public health area when files of individual patients were brought together using name, date-of-birth and other information [16].

One of the main motivations for the utilize of the record linkage method is the construction of the big data bases for answer to the new informative needs [8]. In order to better understand the problem, small practical example is now presented. Suppose the user wants to link two datasets of persons A and B, whose the variables Name, Address and Age are known.

Suppose that Table A contains the following values:

Table A: Data in the first dataset

Unit	Name	Address	Age
a1	John A Smith	16 Main Street	16
a2	Javier Martinez	49 E Applecross Road	33
a3	Gillian Jones	645 Reading Aev	22

Furthermore, suppose that Table B contains the following values:

Table B: Data in the second dataset

Unit	Name	Address	Age
b1	J H Smith	16 Main St	17
b2	Haveir Marteenez	49 Aplecross Raod	36
b3	Jilliam Brown	123 Norcross Blvd	43

The matching table  $A \times B$  contains two units referring probably to the same persons, that the method should individuate as matches: 'John A Smith' with 'J H Smith' and 'Javier Martinez' with 'Haveir Marteenez'.

Modern record linkage begins with the pioneering work of Newcombe et al. [14], who introduced odds ratio of frequencies and the decision rules for delineating matches and nonmatches. In recent years, advances have yielded computer system that incorporate sophisticated ideas from computer sciences, statistics and operational research [16].

Then, Fellegi and Sunter [9] introduced a mathematical foundation for record linkage. Their theory demonstrated the optimality of the decision rules used by Newcombe and introduced a variety of ways of estimating crucial matching probabilities (parameters) directly from the files being matches.

Formally, given two files A and B to be matched, each pair  $(a,b) \in \Gamma = A \times B$  has to be classified into *true match* or *true nonmatch*. The odds ratios of probabilities is:

$$R = \frac{Pr(\gamma \in \Gamma \mid M)}{Pr(\gamma \in \Gamma \mid U)}$$

where  $\gamma$  is an arbitrary agreement pattern in the comparison space  $\Gamma$ , M is the set of of true matches and U is the set of true nonmatches. Between these two sets, the intermediate set of the possible matches exists.

The decision rule reported below helps to classify the pairs:

- if R > Upper, then the pair (a, b) is a designated match,
- if  $Lower \leq R \leq Upper$ , then the pair (a, b) is a designated potential match,

• if R < Lower, then the pair (a, b) is a designated nonmatch.

The estimation of the thresholds Upper and Lower is not easy in an objective way; the choice is competence of the analyst. In the decision rule, three different sets were created: the designated matches, designated potential matches, designated nonmatches. They constitute the partition of the set of all the records in the space  $\Gamma$  in three subsets  $C_3$  (matches),  $C_2$  (potential matches) and  $C_1$  (nonmatches), whose intersections are empty sets.

The idea is to solve the record linkage problem as a multi-criteria based classification problem, whose a priori defined classes are the subsets of the partition.

Without going into too much details, in the next section a brief introduction to the method Electre Tri is presented.

#### 3 The multi-criteria method Electre Tri

In Multi Criteria Decision Aid, a finite set of objects (alternatives, actions, projects) is evaluated by a finite set of criteria, which measure their performances. A criterion is the real-valued function  $g_j: A \to \Re$ , such that  $g_j(a_k)$  indicates the performance of the alternative  $a_k$  on the criterion  $g_j$ . The comparison of any pair of alternatives  $a_i$  and  $a_k$  may be grounded to the comparison of the two values  $g_j(a_i)$  and  $g_j(a_k)$  [13].

In general, a criterion can be either of gain or cost type; gain means that the DM prefers the highest value, while cost means that the DM prefers the lowest value on the criterion.

Many types of criterion were studied in literature, such as true-criterion, pseudo-criterion, pre-criterion, semi-criterion and other types [13].

In the case of true-criterion, if the difference between two performances is positive, then the true-criterion structure implies that the alternatives are in the strict preference relation; while if the difference is equal to 0, then they are in indifference relation.

The Electre Tri is a pseudo-criterion-based method. This type of criterion takes into account that data can be affected by errors from uncertainty, imprecision and small differences or big can not imply the same binary relations. Small and big differences of performances have to imply different binary relations. To define "small" and "big", two values are considered, which are the preference and indifference thresholds.

In literature, grouping problems can be divided in clustering, classification and sorting problems, depending on the a *priori/posteriori* knowledge of classes. The sorting problem is a classification problem, dealt with multi-criteria approach, requiring to Decision Maker (DM) any *preference information*. So, the aim of an ordinal sorting problem consists in assigning each alternative in one of the ordered predefined categories.

Formally, given p predefined ordered categories  $C_1, C_2, \ldots, C_p$  and a finite set of n alternatives  $A = \{a_1, a_2, \ldots, a_n\}$ , evaluated on a finite set of m criteria  $G = \{g_1, g_2, \ldots, g_m\}$ , in the case all criteria are gain-type, the relations among

the categories are  $C_1 \prec C_2 \prec \ldots \prec C_p$ , such that  $b_h$  is the profile, upper limit of category  $C_h$  and lower limit of category  $C_{h+1}$ . In this way,  $C_1$  and  $C_p$  are the worst and the best categories respectively.

The Electre Tri method is based on outranking relations, indicated with S, which characterize how the alternatives are compared with the profiles. Because the assignment of an alternative to a specific category follows from the comparison, on all criteria, of its performances with the profiles ones.

The relation  $aSb_h$  validates or invalidates the assertion "a outranks  $b_h$ " whose meaning is "a is at least as good as  $b_h$ ", on the set G.

In the context of the Electre Tri method, the validation of outranking relation is made by the computation of four indices [12, 13]:

- 1. the partial concordance indices on each criterion;
- 2. the global concordance index on all the criteria;
- 3. the partial discordance indices on each criterion;
- 4. the credibility index on all the criteria.

For the computation of the partial concordance indices, it is necessary to know the profiles, preference and indifference thresholds values. In the case one of these parameters are not known, the index can not be computed. For the computation of the global concordance index is necessary to know the weights, representing the importance coefficients of the criteria. For the computation of the partial discordance indices, it is necessary to know the profiles, preference and veto thresholds values. And the credibility index corresponds to the global concordance index weakened by veto effects. If veto thresholds do not enter in the model, the credibility index is equal to the global concordance index. From the credibility index to the definition of an outranking relation, it is necessary to fix a cutting level lambda, which is the minimum credibility index value which permits to define the outranking relation. Finally, the assignment of an alternative to one category does not result from the outranking relation directly, but it is necessary to use one (or both) of the two proposed exploitation procedures. They are the pessimistic and the optimistic assignment procedures. These procedures analyze the way an alternative compares to the profiles so as to determine the category to which the alternative should be assigned.

One of the main difficulties is the elicitation of various parameters that in the Electre Tri are profiles, weights, thresholds (preference, indifference and veto) and cutting level lambda. Even if these parameters can be interpreted, it can be difficult to fix directly their values (*direct elicitation*) and to have a clear global understanding of the implications of these values in terms of the output [12].

In order to estimate indirectly the value of the parameters, De Leone and Minnetti [6] proposed new estimation methodology whose procedure is composed of two phases: the first dedicated to the profiles and thresholds estimations,

the second to the weights and cutting level estimations. The core of the procedure is the profiles' estimation, suggested with Linear Programming (LP) using training set.

Let p be the number of categories, m the number of criteria, the LP problem is the following:

min 
$$\sum_{j=1}^{m} \sum_{a_k \to C_h} \theta_j(a_k)$$
s.t. 
$$\theta_j(a_k) \ge g_j(a_k) - g_j(b_h) \quad \forall j = 1, \dots, m, \forall a_k \to C_{h,h \neq p}$$

$$\theta_j(a_k) \ge g_j(b_{h-1}) - g_j(a_k) \quad \forall j = 1, \dots, m, \forall a_k \to C_{h,h \neq 1}$$

$$g_j(b_h) \ge g_j(b_{h-1}) + \epsilon \quad \forall j = 1, \dots, m, \forall h = 2, \dots, p-1$$

$$\theta_j(a_k) \ge 0 \quad \forall j = 1, \dots, m, \forall a_k \to C_h$$

where  $\epsilon$  is a small positive value.

The problem (1) minimizes the sum of the classification errors  $\theta_j(a_k)$  on all criteria and on all the alternatives in the training set, when this alternative's performance lies out the belonged category. The first two constraints define the error  $\theta_j(a_k)$ .

### 4 Application to Real Data: a first experiment

As said in the previous section, the multi-criteria approach requires DM any preference information, including binary relations. Since it is possible to state binary relations between the subsets as  $C_3 \succ C_2 \succ C_1$ , the record linkage problem can be structured as ordinal sorting problem, that is, classification problem whose classes are ordered in the strict preference binary relations.

Moreover, the importance of using multi-criteria decision methods, with respect to the other classification methods, is in the possibility to assign weights to each criterion, not possible in all the classification problems, and to use the *preference information*, provided by DM, for estimating the classification model's parameters.

The proposed application wants to find a classification model (i.e. classifier or learner), assigning each record of the space  $\Gamma$  to one of the three categories  $C_1$ ,  $C_2$  and  $C_3$ , following the two phases procedure formulated by De Leone and Minnetti [6].

The input data, used in the application, were taken from Winkler from American Census (in SecondString file for approximate string matching techniques). Two data sets A and B are considered, containing 449 and 392 records respectively, and the true links are 327.

The variables (textual fields from synthetic census data) are the following:

- DS (labels of the data sets with A and B);
- IDENTIFIER;
- SURNAME;

- NAME;
- LASTCODE (middle name initial);
- NUMCODE (address street number);
- STREET (address street name).

In this short paper, results from preliminary experiments are reported, because the application is an ongoing research, due to its complexity.

Some variables contain missing values that cause difficulties in the analysis, making it more complicated. So in order to facilitate the analysis, the records with missing values are deleted.

There are a number of popular methods of estimating the learner's ability to generalize; the test set method was used here. In this experiment, the use of distance measure and the search of training set had played the most important roles; they had contributed to obtain good results of the classification model, found by Electre Tri [11].

The performance of the classification model, applied to the test set (83868 alternatives) was 99.09% when all the criteria have the same importance and the lambda parameter is  $\lambda = 0.50$ . If lambda increases, the performance increases, up to 99.81% when  $\lambda = 0.70$  and 99.89% when  $\lambda = 0.85$ .

In the case the importance coefficients of criteria were considered different, the performances of the models were substantially the same, varying the lambda parameter.

In the case of performance 99.09%, the classification errors were committed by the model on the false links; namely, the model saw almost all the true links. The opposite situation occurred in the case of performance 99.89%, when the model saw almost all the false links and misclassified the true links. To the DM the choice of the most interesting model, depending his preferences.

#### 5 Final Conclusion and Remarks

In this short paper, the Electre Tri machine learning technique was proposed for solving Record Linkage matching. It is the first time that multi-criteria decision technique is used to solve the record linkage problem.

The proposed application started with an initial experiment demonstrating that the application of the Electre Tri to record linkage shall provide good results in terms of classifier performances. This paper shows only the results of a preliminary experiment, which provided good results in terms of performances of the classification model. Also this experiment confirmed that record linkage is more sensitive to the quality of preprocessing and standardization that of matching, as said in [17].

As consequence, other measures of distance in the construction of the input data matrix, as well as, different schemes in the search of training set, will be used.

#### References

- [1] M. Bilenko and R. Mooney. Adaptive duplicate detection using learnable string similarity. *ACM SIGKDD*, pages 39–48, 2003.
- [2] P. Christen. Automatic record linkage using seeded nearest neighbour and support vector machine classification. *ACM SIGKDD*, pages 151–159, 2008.
- [3] P. Christen. Automatic training example selection for scalable unsupervised record linkage. *PAKDD*, *Springer LNAI*, 5012:511–518, 2008.
- [4] W.W. Cohen. The whirl approach to data integration. *IEEE Intelligent Systems*, 13,no.3:20–24, 1998.
- [5] W.W. Cohen, P. Ravikumar, and S.E. Fienberg. A secure protocol for computing string distance metrics. In *PSDM held at ICDM*, pages 40–46, 2004.
- [6] R. De Leone and V. Minnetti. The estimation of the parameters in multicriteria classification problem: The case of the electre tri method. In Donatella Vicari, Akinori Okada, Giancarlo Ragozini, and Claus Weihs, editors, Analysis and Modeling of Complex Data in Behavioral and Social Sciences, Studies in Classification, Data Analysis, and Knowledge Organization, pages 93–101. Springer International Publishing, 2014. ISBN: 978-3-319-06691-2.
- [7] M. Elfeky and A. Elmagarmid. Tailor: A record linkage toolbox. *IEEE ICDE*, pages 17–28, 2002.
- [8] I.P. Fellegi. Record linkage and public policy: a dynamic evolution. In W. Alvey and B. Jamerson, editors, *Proceedings of International Workshop and Exposition*. Record Linkage Techniques, 1997, March, 1997.
- [9] I.P. Fellegi and A.B. Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 64:1183–1210, 1969.
- [10] C. Kenneth and C. Poirier. Machine learning documentation initiative. Modernisation Committee on Production and Methods, Statistics Canada, Febrary 4, 2015.
- [11] V. Minnetti. On the parameters of the Electre Tri method: a proposal of a new two phases procedure. PhD thesis, Faculty of Information Engineering, Informatics and Statistics, Sapienza University of Rome, March 2015.
- [12] V. Mousseau and R. Slowinski. Inferring an electre tri model from assignment examples. *Journal of Global Optimization*, 12:157–174, March 1998.
- [13] V. Mousseau, R. Slowinski, and P. Zielniewicz. *ELECTRE TRI 2.0, a methodological guide and userś manual.* Document du LAMSADE no111, Universit Paris-Dauphine, 1999.

- [14] H.B. Newcombe, J.M. Kennedy, S.J. Axford, and A.P. James. Automatic linkage of vital records. *Science*, 130(3381):954–959, October 1959.
- [15] D.R. Wilson. Beyond probabilistic record linkage: using neural network and complex features to improve genealogical record linkage. In *Proceedings of International Joint Conference on Neural Networks*, San Jose, California, USA, 2011.
- [16] W.E. Winkler. The state of record linkage and current research problems. Technical report, U.S. Bureau of the Census, Statistics of Income Division, Internal Revenue Service Publication, R99/04, 1999.
- [17] W.E. Winkler. Matching and record linkage. WIREs Comput Stat, 6:313–325, 2014.