

## Performance of Spanish Encoding Functions during Record Linkage

María del Pilar Angeles, Noemi Bailón-Miguel

Facultad de Ingeniería

Universidad Nacional Autónoma de México

Ciudad de México, México

e-mail: pilarang@unam.mx, mimibailon@hotmail.com

**Abstract**—Nowadays, many businesses suffer from duplicate records. For instance, information about the same provider, customer or product appears in multiple systems and in multiple formats across the company and simply does not tally from system to system. This situation seriously prevents managers to make well informed decisions. In the case of low data quality written in Spanish language, the identification and correction of problems such as spelling errors with English language based coding techniques is not suitable. In this paper, we have implemented, modified, and utilized three Spanish phonetic coding functions in our prototype called Universal Evaluation System of Data Quality (SEUCAD). A Spanish phonetic coding based on Soundex algorithm, a Spanish Metaphone coding, and a Modified version of the latter were utilized to detect duplicate text strings in the presence of spelling errors in Spanish. The results were satisfactory, the Spanish phonetic algorithm performed well most of the time, demonstrating opportunities for an improved performance of Spanish encoding during the record linkage process.

**Keywords**— data mining; data matching; de-duplication; record linkage.

### I. INTRODUCTION

The existence of duplicate records has strong implications on the use and scope of data. Low data quality affects decision making. For instance, the financial industry has faced several frauds caused by duplicate data. All financial institutions are interested in decreasing the already existing number of duplicates and implementing a more efficiently data handling in order to avoid future duplicate data. In the case of duplicate medical records, when the system is unable to find a reliable patient record the risk of wrong medical treatment, or over-immunization, is present, along with the corresponding cost of unnecessary immunizations, or the risk of adverse effects on patients, etc. Therefore, there has been a significant research in the area of data quality and data matching during the last decade.

We have developed a prototype called Universal Evaluation System of Data Quality (SEUCAD) [1] on the basis of the Freely Available Record Linkage System (FEBRL) [2].

We have previously compared, added and improved a number of data matching methods. Our prototype allows end users to assess density, coverage, completeness [1][3], and performs a complete data matching process in order to identify duplicate records. However, most of the coding

algorithms are based on the English language, few approaches are oriented to the Spanish language. Consequently, the encoding algorithms are not efficient enough to detect common errors and misspellings in the process of data matching for the improvement of quality of data. This problem impacts all the industry projects that are related for instance, to data mining, data science, business intelligence, and big data for companies where data are written in the Spanish language.

Our research has been lately focused on the implementation and enhancement of Spanish encoding functions in order to improve the performance of the encoding phase during entity resolution when data have been written in Spanish language.

Within our SEUCAD prototype, the Phonex, Soundex, and Modified Spanish phonetic functions have been previously compared, and our findings published in [3].

The Spanish phonetic coding was proposed in [4], which is an extended Soundex coding, where Spanish characters have been added. Besides, we have modified the Spanish Phonetic Algorithm so the encryption code is resizable, and all white spaces are removed during encoding. The previous comparison showed that the modified version of the Spanish Phonetic Algorithm had a better performance in terms of precision.

The present document shows the implementation of two more Spanish encoding functions: the Spanish Metaphone algorithm [5][6], and a second version of such an algorithm, which applies the same code to similar sounds derived from very common misspellings.

The record linkage outcomes for these three coding functions have been evaluated under a number of different scenarios, where the true match status of record pairs was known. We have obtained precision, recall, and f-measure because they are suitable measures to assess data matching quality.

The present paper is organized as follows: The next section briefly explains the data matching process and how it has been implemented within SEUCAD. Section III explains the phonetic encoding functions proposed from previous research, the enhancements we have implemented on some of them, along with their role within the process of data matching. Section IV presents the experiments carried out, and analyses the results. Finally, the last section concludes the main topics achieved regarding the performance of the encoding functions and the future work to be done.

## II. RELATED WORK

The data matching process is mainly concerned to the record comparison among databases in order to determine if a pair of records corresponds to the same entity or not [7]. It is also called record linkage or de-duplication. In general terms, this process consists on the following tasks:

a) A standardization process [7], which refers to the conversion of input data from multiple databases into a format that allows correct and efficient record correspondence between two data sources.

b) Phonetic encoding is a type of algorithm that converts a string into a code that represents the pronunciation of that string. Encoding the phonetic sound of names avoids most problems of misspellings or alternate spellings, a very common problem on low quality of data sources.

c) The indexing process aims to reduce those pairs of records that are unlikely to correspond to the same real world entity and retaining those records that probably would correspond in the same block for comparison; consequently, reducing the number of record comparisons. The record similarity depends on their data types because they can be phonetically, numerically or textually similar. Some of the methods implemented within our prototype SEUCAD are for instance, Soundex [9], Phonex [2], Phonix [2], NYSIIS [10], and Double metaphone [5].

d) Field and record comparison methods provide degrees of similarity and define thresholds depending on their semantics or data types. In the prototype, the algorithms Qgram, Jaro - Winkler Distance [11][12], Longest common substring comparison are already implemented.

e) The classification of pairs of records grouped and compared during previous steps is mainly based on the similarity values that were already obtained, since it is assumed that the more similar two records are, there is more probability that these records belong to the same entity of the real world. The records are classified into matches, not matches or possible matches.

The SEUCAD prototype was aimed to the development of algorithms that reduce the quadratic complexity of the naive process of pair-wise comparing each record from one database with all records in the other database, and how to accurately classify the compared record pairs into matches and non-matches considering attributes dependency.

Nowadays, SEUCAD is able to measure, assess and help during the analysis of data quality process [1] under a number of open and licensed database management system (DBMS), such as Oracle DB, MySQL, IBM DB2, SAP-Sybase Adaptive Server Enterprise, SAP-Sybase IQ, and EnterpriseDB PostgreSQL.

The SEUCAD application extracts the database schema directly from the data dictionary and measures the intrinsic quality of the data through the following indicators: coverage, density, completeness [13]. Since these measures are intrinsically computed through SQL queries, the assessed granularity levels are at database, table and column where applicable as we have done in previous research [14]. Furthermore, the prototype implements a specific framework for the detection, classification and fusion (cleaning) of

duplicate records within a number of databases (data matching and de-duplication) with no regard of the type of data source.

During the implementation of some data matching algorithms, we have realized that the coding functions mainly used were on the basis of English language. Such algorithms were not suitable for Spanish written data already stored in our databases. Therefore, we were focused on the implementation and experimentation of Spanish encoding functions in order to improve the performance of the encoding phase during entity resolution.

We have implemented and enhanced two Spanish encoding functions in order to improve the performance of the encoding phase during entity resolution when data has been written in Spanish language, and the corresponding results are shown in the present work.

The aim of the following section is to briefly explain the phonetic encoding functions that we have implemented and enhanced in order to quantify and compare their performance during the record linkage process.

## III. PHONETIC ENCODING PROPOSALS TO COMPARE

### A. Phonetic coding functions

Phonetic encoding is a type of algorithm that converts a string (generally assumed to correspond to a name) into a code that represents the pronunciation of that string. Encoding the phonetic sound of names avoids most problems of misspellings or alternate spellings, a very common problem on low quality of data sources.

### B. Spanish phonetic

The Spanish phonetic coding function compared in the present document is a variation of the Soundex algorithm. Soundex is a phonetic encoding algorithm developed by Robert Russell and Margaret Odell in [9], and patented in 1918 and 1922. It converts a word in a code [15]. The Soundex code is to replace the consonants of a word by a number; if necessary zeros are added to the end of the code to form a 4-digit code. Soundex choose the classification of characters based on the place of articulation of the English language.

The limitations of the Soundex algorithm have been extensively documented and have resulted in several improvements, but none oriented to the Spanish language. Furthermore, the dependence of the initial letter, the grouping articulation point of the English language, and the four characters coding limit are not efficient to detect common misspellings in the Spanish language.

The Spanish phonetic coding was proposed in [4], it is an extended Soundex coding, where Spanish characters have been added. In general terms, the algorithm is as follows:

1. The string is converted to uppercase with no consideration of punctuation signs.
2. The symbols "A, E, I, O, U, H, W" are eliminated from the original word.
3. Assign numbers to the remaining letters according to Table 1.

TABLE I. SPANISH CODING

Characters	Digit
P	0
B, V	1
F, H	2
T, D	3
S, Z, C, X	4
Y, LL, L	5
N, Ñ, M	6
Q, K	7
G, J	8
R, RR	9

We have modified the Spanish Phonetic Algorithm [3] so the encryption code is resizable, and all white spaces are removed during encoding. This model allows us to analyze a larger number of cases where we can have misspellings. The modified Spanish phonetic algorithm is called as soundex\_sp in our SEUCAD prototype.

C. The Spanish Metaphone Algorithm

The Metaphone is a phonetic algorithm for indexing words by their English sounds when pronounced, it was proposed by Lawrence Philips in 1990 [5]. The English Double-Metaphone algorithm was implemented by Andrew Collins in 2007 who claims no rights to this work. The Metaphone port adapted to the Spanish Language is authored by Alejandro Mosquera in [6]; we have implemented this function and called as Esp\_metaphone in our SEUCAD prototype. Some of the changes applied in order to adjust to the Spanish language are shown in Table II, which considers typical cases of the Spanish language with letters such as á, é, í, ó, ú, ll, ñ, h.

TABLE II. SPANISH METAPHONE

Char	Replacement
á	A
ch	X
C	S
é	E
í	I
ó	O
ú	U
ñ	NY
ü	U
b	V
Z	S
ll	Y

D. Modified Spanish Metaphone coding function

In Spanish language, there are words such as “oscuro”, “oscurio” or “combate”, “convate” that should share the same code because even they are written different, their sound is similar and the misspelling is common. The second version of Esp\_metaphone contains the following enhancements:

The Royal Academy of the Spanish Language reviewed words that originally were written with “ps” as “psicología”, and introduced some changes, because "the truth is that in

Castilian the initial sound ps is quite violent, so the ordinary, both in Spain and in America, it is simply pronounced as “sicología”. Moreover, our language, differing French or English, is not greatly concerned with preserve the etymological spelling; He prefers the phonetic spelling and therefore tends to write as it is pronounced” [16]. Words that begin with “ps” can be written and pronounced as “s”, and are called silent letters; for example, words psicólogo and sicólogo. We have added some cases to the Spanish Metaphone algorithm in order to consider these possible variations in Spanish written words and to assign the same code in both cases. Therefore, in case there is a word that starts with “ps”, it will be replaced by “s”. A special case with silent letter is presented with words like “oscuro” and “obscuro”, where both words have the same meaning so that the use of both is correct. In this case both its meaning and pronunciation is usually the same. Then, in case there is a word that starts with “bs”, it shall be replaced by “s”. One case of a common misspelling in Spanish language is given with words like ”tambien” and “tanbien” were the latter is orthographically wrong, but phonetically is very similar to the former, and in case of typos, the letter “n” is close to letter “m” in a keyboard. Thus, we have decided to replace "mb" by "nb" and assign the same code. We have decided to replace "mp" by "np" and assign the same code in case of words such as “tampoco” and “tanpoco”. The words that begin with “s” followed by a consonant are replaced by 'es' such as “scalera” and “escalera”. Later all the letters “s” are replaced by “z”. Table III shows the additions contained in the Spanish Metaphone version 2.

TABLE III. MODIFIED SPANISH METAPHONE

Char	Replacement
mb	nb
mp	np
bs	s
ps	z

Table IV shows coding from Metaphone and Metaphone\_v2, the former is not able to apply the same code to words “psiquiatra”, “siquiatra”, “oscuro”, “obscuro”, “combate”, “convate”, “conbate”. All these words have the same meaning and in order to identify duplicates they should have the same code.

TABLE IV. SPANISH METAPHONE AND SPANISH METAPHONE V2 CODING

Word	Metaphone	Metaphone_v2
Cerilla	ZRY	ZRY
Empearar	EMPRR	ENPRR
Embotellar	EMVTYR	ENVTYR
Xochimilco	XXMLK	XXMLK
Psiquiatra	PSKTR	ZKTR
siquiatra	SKTR	ZKTR
Oscuro	OVSKR	OZKR
Oscuro	OSKR	OZKR
Combate	KMBT	KNVT
Convate	KNVT	KNVT
Conbate	KNBT	KNVT
Comportar	KMPRTR	KNPRTR

Conportar	KNPRTR	KNPRTR
Zapato	ZPT	ZPT
Sapato	SPT	ZPT
Escalera	ESKLR	EZKLR
scalera	ESKLR	EZKLR

In the case of code generated by Metaphone\_v2 the code is the same, although there are not identical texts because of spelling mistakes but same meaning.

The three coding functions we have explained in this section are meant to increase the similarity between the words written and the sound they represent in Spanish language in order to avoid common Spanish misspelling and errors and enhance the performance of the following steps during the data matching process. For instance, the level of similarity obtained among two words should be increased even in the case of a word was written as “siquiatra” rather than “psiquiatra”.

The following section is concerned with the set of experiments we have carried out in order to identify how the coding functions we have implemented within SEUCAD can help to data matching with data written in Spanish language

#### IV. EXPERIMENTS

We have developed and executed a set of experiments for the record linkage process through four scenarios, each scenario containing a different data-source. In this Section we will explain a) how the quality of the data matching process will be computed ;b) the configuration of all the record linkage process; c) the characteristics of each data-source according to each scenario; and c) the analysis of the outcomes from each scenario.

These experiments are aimed to identify for each data-set which encoding function has the best performance. The performance of the record linkage process is measured in terms of how many of the classified matches correspond to true real-world entities, while matching completeness is concerned with how many of the real-world entities that appear in both databases were correctly matched [7][8]. Each of the record pair corresponds to one of the following categories: True positives (TP). These are the record pairs that have been classified as matches and are true matches. These are the pairs where both records refer to the same entity. False positives (FP). These are the record pairs that have been classified as matches, but they are not true matches. The two records in these pairs refer to two different entities. The classifier has made a wrong decision with these record pairs. These pairs are also known as false matches.

True negative (TN). These are the record pairs that have been classified as non-matches, and they are true non-matches. The two records in pairs in this category do refer to two different real-world entities. False negatives FN). These are the record pairs that have been classified as non-matches, but they are actually true matches. The two records in these pairs refer to the same entity. The classifier has made a wrong decision with these record pairs. These pairs are also known as false non-matches. Precision calculates the

proportion of how many of the classified matches (TP + FP) have been correctly classified as true matches (TP). It thus measures how precise a classifier is in classifying true matches [9]. Precision is calculated as  $TP/(TP+FP)$ . Recall measures how many of the actual true matching record pairs have been correctly classified as matches [9]. It is calculated as:  $recall= TP/(TP+FN)$ . F-measure is a measure that combines precision and recall is the harmonic mean of precision and recall. Thus, is calculated as  $2TP/(2TP+FP+FN)$ .

An ideal outcome of a data matching project is to correctly classify as many of the true matches as true positives, while keeping both the number of false positives and false negatives small. Based on the number of true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN), different quality measures can be calculated. However, most classification techniques require one or several parameters that can be modified and depending upon the values of such parameters, a classifier will have a different performance leading a different numbers of false positives and negatives.

For each data source, the number of total records, the number of duplicated records, the maximum number of duplicated record for an original record, the maximum number of changed fields per item, and the maximum number of record modifications were considered as independent variables. The dependent variables will be the amount of matches, non-matches and possible matches. The quality of the data matching process will be obtained from precision, f-measure, and all the metrics we have already mentioned in this Section. The control variables (also known as constant variables) will correspond to the indexing, comparison and classification steps within the data matching process because the experiments are aimed to identify which coding function will perform the best. All the data sources presented a uniform probability distribution for duplicates. Fig. 1 shows the structure and sample source data utilized for experimentation.

nombre	apellido_pat	apellido_mat	calle
santiago		gonzalez	calle de san gumersindo
david	hernandez	cruz	calle de arnedillo
jessica	perez	martinez	calle de barbara de braganza
martha	sanchez	lopez	calle de jordi sole tura
patricia	garcia	aviles	calle de santa maria reina
alfonso	garcia	hernandez	calle del iridio
adriana	vazquez	gonzalez	calle de jose espelius
tania	mendez	lopez	plaza de arguelles
vicente		reyes	calle de infiesto
angelica	hernandez	brito	calle de los hermanos carpi
maria elena	perez	ramirez	calle de los bascones
isaac	martinez	gutierrez	calle de julia garcia bouton
berenice	ramirez	reyes	calle de elvira barrios
alejandro	alonso	flores	calle de la anunciacion
enrique	cordero	ramirez	calle del gladiolo

Figure 1. Sample of data source

The configuration of indexing, comparison and classification for all scenarios has been the same and repeated for each encoding function (Esp-Metaphone, Esp\_metaphone\_v2 and Soundex\_sp). Such configuration is presented as follows:

1. Indexing: Fields that form the record require to be encoded and indexed in order to avoid a large number of comparisons between records whose fields are not even similar. During the coding phase, we have executed for each experiment one of the coding functions: esp-metaphone, esp\_metaphone\_v2 or soundex\_sp. In order to execute the indexing step, we have chosen "Blocking index" as indexing method based on fields: "nombre", "apellido paterno", "apellido materno", "calle". Fig. 2 shows the configuration utilized for indexing and encoding methods.

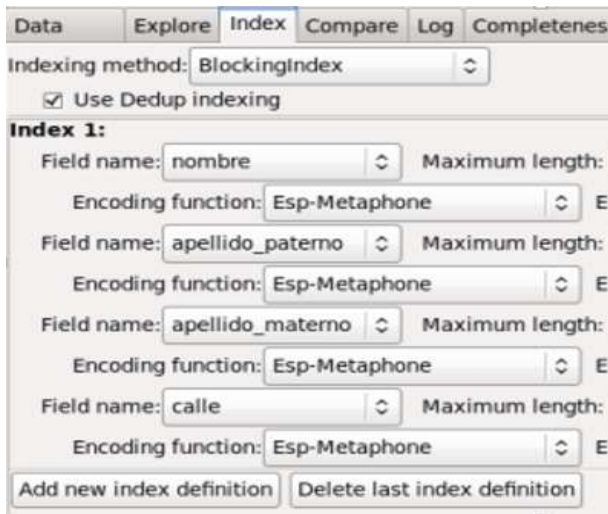


Figure 2. Indexing and encoding configuration

2. Comparison: Once records have been ordered and grouped in terms of the previous fields specified. Each encoded field will be compared. In order to obtain quality measures during the comparison step, we have chosen an exact function "Str-Exact", which requires an exact match on strings compared. This function will be used with the fields named as "nombre", "apellido paterno", "apellido materno", "calle". Fig. 3 shows the comparison specification for the experiments.

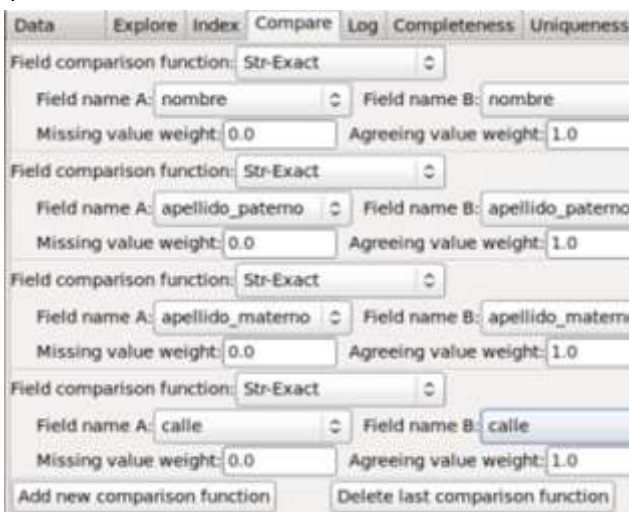


Figure 3. Comparison by String Exact method

3. Classification: In the case of pairs of record classification, we have selected the Optimal Threshold method, with a

minimized false method of Positives and negatives, and a bin width of 40 for the range of values to be considered for the output graphic. Fig. 4 shows the classification configuration for the experiments.

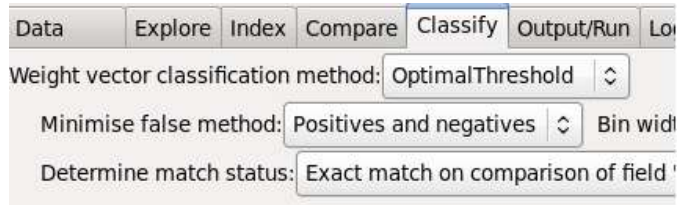


Figure 4. Classification by Optimal Threshold

The following scenarios are presented in order to show the performance of the encoding functions. The corresponding tables show the values of true positives, false positives, precision computed as  $TP/(TP+FP)$ , and F-measure computed as  $2TP/(2TP+FP+FN)$ . The value of false negatives was 0 in all scenarios and encoding functions.

A. Scenario 1

The first file was generated with a total length of 1000 records, 100 duplicated records, one duplicated record for an original record as maximum, one change field per item as maximum, one maximum record modification, with a uniform probability distribution for duplicates. The quality metrics obtained for each encoding method are presented in Table V.

TABLE V QUALITY METRICS FOR SCENARIO I

Encode Method	Total Classif.	TP	FP	Precision	F-measure
Metaphone_sp	68	65	3	0.95588	0.977443
Metaphone_v2	69	66	3	0.95652	0.977777
Soundex_sp	76	73	3	0.96052	0.979865

According to the outcomes obtained from the first scenario, we can observe that in the case of the Modified Spanish coding function (soundex\_sp), there were 76 record pairs classified, with 73 duplicated record pairs as true positives and 3 record pairs as false positives. Therefore, this method was more precise with 96% than the rest of the functions.

B. Scenario II

The second data source contained a total length of 5000 records, 500 duplicated records, one duplicated record for an original record as maximum, one change field per item as maximum, one maximum registry modification, with a uniform probability distribution for duplicates. The quality metrics obtained for each encoding method are presented in Table VI.

TABLE VI QUALITY METRICS FOR SCENARIO II

Encode Method	Total Classif.	TP	FP	Precision	F-measure
Metaphone_sp	320	319	1	0.9968	0.9984
Metaphone_v2	341	340	1	0.99706	0.99853
soundex_sp	353	352	1	0.99716	0.99581



From Table VI we can observe that the Modified Spanish function classified 353 record pairs, with 352 duplicated record pairs as true positives and 1 record pair mistakenly classified as true match, corresponding then as one false positive. Therefore, this method was 99.7% precise, with more records classified than the Metaphone\_sp and Methaphone\_v2 with 320 and 341 records classified respectively.

C. Scenario III

The third data source contained a total length of 10000 records, 5000 duplicated records, one duplicated record for an original record as maximum, one change field per item as maximum, one maximum registry modifications, with a uniform probability distribution for duplicates.

The process of record linkage under this scenario showed that the Modified Spanish coding function classified 3622 record pairs out of a total of 5000 potentially to detect, with 3620 duplicated record pairs as true positives and 2 record pairs mistakenly classified as true match. Therefore, this method was 99.94% precise. The Metaphone\_sp and Methaphone\_v2 phonetic functions obtained less records classified and more false positives than Spanish soundex function. The quality metrics obtained for each encoding method are presented in Table VII.

TABLE VII QUALITY METRICS FOR SCENARIO III

Encode Method	Total Classif.	TP	FP	Precision	F-measure
Metaphone_sp	3333	3324	9	0.997299	0.9986
Metaphone_v2	3489	3480	9	0.99742	0.9987
Soundex_sp	3622	3620	2	0.99944	0.9997

D. Scenario IV

The fourth file has a total length of 1000 records, 100 duplicated records, one duplicated record for an original record as maximum, two changed fields per item as maximum, three maximum registry modifications, with a uniform probability distribution for duplicates.

The Modified Spanish coding function, allowed that 964 record pairs could be classified, the total number of duplicates was actually 2500 records. However, this method did not present any false positive. The rest of the phonetic algorithms were 99% precise with two false positives, but the number of classified records was lower than those with Soundex\_sp. The outcomes obtained for each encoding method under scenario IV are presented in Table VIII.

TABLE VIII QUALITY METRICS FOR SCENARIO IV

Encode Method	Total Classif.	TP	FP	Precision	F-measure
Metaphone_sp	812	810	2	0.997536	0.99876
Metaphone_v2	884	882	2	0.99773	0.99886
Soundex_sp	964	964	0	1	1

E. Analysis of Outcomes

According to the outcomes shown in previous section, we can observe that the Modified Spanish Phonetic algorithm was always more precise than the rest of the algorithms. Therefore, the Modified Spanish-Phonetic

algorithm allows a higher proportion of how many of the classified matches (TP+FP) have been correctly classified as true matches. The Spanish phonetic algorithm allows a greater number of similarities than the remaining algorithms in all cases, because is more effective codifying Spanish words. The Spanish phonetic algorithm achieved a slightly higher f-measure than the two versions of the Spanish Metaphone algorithm.

The graphics presented in this section, have been generated according to the variation of the coding function in order to observe the behavior of the algorithms. The precision obtained from each encode method for all the scenarios have been compared, graphed and shown in Fig. 5 shows the trend of the contribution of each encoding method to the precision of the classification. As we can observe, the Spanish coding function was above the Metaphone base coding algorithms.

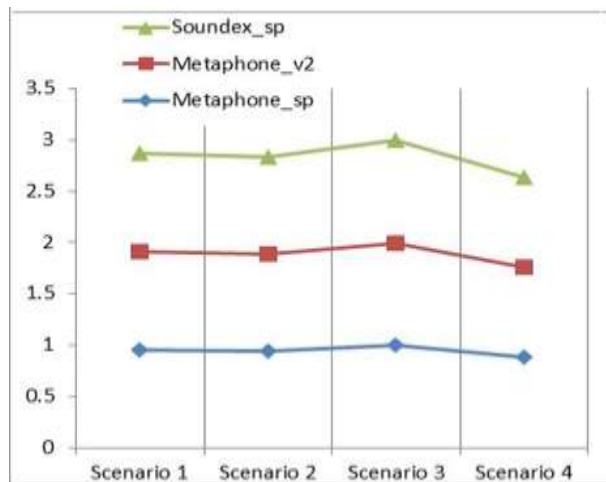


Figure 5. Precision of each encode function

Fig. 6 shows the trend of the contribution of each encoding method to the completeness of the classification. In other words, the proportion of record pairs classified against the entire number of duplicates per scenario.

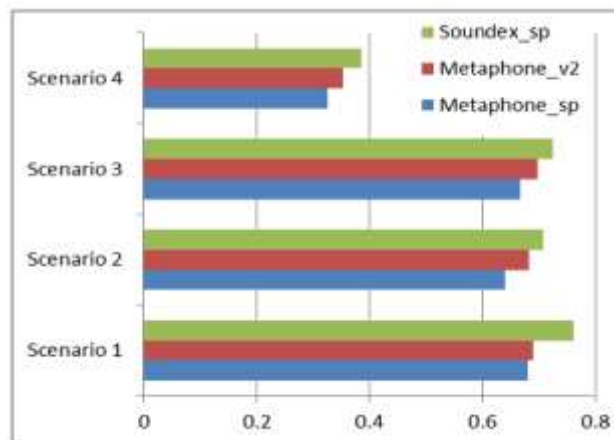


Figure 6. Completeness for each coding function per scenario

According to the outcomes shown in previous section, we can observe that the Modified Spanish Phonetic algorithm was always more precise than the two versions of Metaphone. Therefore, the Modified Spanish-Phonetic algorithm allows a higher proportion of true matches. The Spanish phonetic algorithm allows a total similarity greater than the remaining algorithms in all cases, because is more effective codifying Spanish words.

The Spanish phonetic algorithm achieved a slightly higher f-measure than the rest.

As we can observe from Fig. 6, the Spanish phonetic algorithm obtained a larger number of pairs of records classified than the rest of the phonetic algorithms.

## V. CONCLUSION

The problem of detection and classification of duplicate records the integration of disparate data sources affects business competitiveness. A number of encoding, comparison and classification methods have been utilized until now, but there still some work to do in terms of effectiveness and performance.

The present work has evaluated the record linkage outcomes under a number of different scenarios, where the true match status of record pairs was known. We have obtained precision, recall, and f-measure because they are suitable measures to assess data matching quality.

The Modified Spanish Soundex function presented a better performance than the rest of the phonetic functions during most of the experiments. However, it takes the longest execution time with a difference of some milliseconds. It is important to be aware that the performance of a de-duplication system or technique is dependent on the type and the characteristics of the involved data sets, having good domain knowledge is relevant in order to achieve good matching or deduplication results.

We have previously concluded in [3] that the Modified Spanish Phonetic algorithm was always more precise and complete than Soundex y Phonex. Under a new set of experiments we have carried out against a Spanish version of the Metaphone algorithm and an enhanced version of the Spanish Metaphone, the Modified Spanish Phonetic algorithm still having the best performance in terms of precision in the majority of the cases we have experimented during the present research. However, the precision presented for the three Spanish coding functions varies slightly as we have utilized a String exact comparison function, the experimentation with different comparison functions that provide different levels of similarity might give more information regarding encoding effectiveness. We will also focus on performance in terms of massive data processing and its corresponding response time, these elements might give us a better criteria in order to identify the best encoding function.

The proposed framework may also be developed and extended to other languages as part of future work.

## ACKNOWLEDGMENT

This work is being supported by a grant from Research Projects and Technology Innovation Support Program (Programa de Apoyo a Proyectos de Investigación e Innovación Tecnológica, PAPIIT, UNAM Project 1N114413 named Universal Evaluation System Data Quality (Sistema Evaluador Universal de Calidad de Datos).

## REFERENCES

- [1] P. Angeles, et al., "Universal evaluation system data quality," DBKDA 2014 : The Sixth International Conference on Advances in Databases, Knowledge, and Data Applications, vol. 32, pp. 13–19, 2014.
- [2] P.Christen, "Febrl a freely available record linkage system with a graphical user interface," Second Australasian Workshop on Health Data and Knowledge Management (HDKM 2008), vol. 80, pp. 17–25, 2008.
- [3] P. Angeles, J. García-Ugalde, A. Espino-Gamez, and J. Gil-Moncada, "Comparison of a Modified Spanish Soundex, and Phonex Coding function during data matching process", International Conference on Informatics, Electronic and Vision, ICIEV, Kytakyushu, Fukuoka Japan, ISBN:978-1-4673 6901-5, DOI:10.1109/ICIEV.2015.7334028, IEEE, pp.1-6,2015.
- [4] F. M. I. Amon, J. Echeverria, "Algoritmo fonetico para deteccion de cadenas de texto duplicadas en el idioma espanol," Ingenierias Universidad de Medellin, vol. 11, no. 20, pp. 120– 138, 2012.
- [5] L. Philips, "The double metaphone search algorithm," C/C++ Users J, vol. 18, no. 6, pp. 38–43, 2000.
- [6] A. Mosquera, E. Lloret, and P. Moreda, "Towards Facilitating the Accessibility of Web 2.0 Texts through Text Normalisation", Proceedings of the LREC workshop: Natural Language Processing for Improving Textual Accessibility, pp. 9–14, 2012.
- [7] P. Christen, Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution and Duplicate Detection. Springer Data-Centric Systems and Applications, 2012.
- [8] T. Churches, P. Christen, K. Lim, and J. X. Zhu, "Preparation of name and address data for record linkage using hidden markov models." BMC Medical Informatics and Decision Making, vol. 2, no. 1, p. 9, 2000.
- [9] M. Odell, R. Russell, "The soundex coding system," American Patent 1 261 167, 1918.
- [10] C. L. Borgman, S. L. Siegfried, "Gettys synonymetm and its cousins: A survey of applications of personal name-matching algorithms," Journal of the American Society for Information Science, vol. 43, no. 7, pp. 459–476, 1992.
- [11] M. A. Jaro, "Advances in record-linkage methodology applied to matching the 1985 census of Tampa, Florida," Journal of the American Statistical Association, vol. 84, pp. 414–420, 1989.
- [11] W. Winkler, "String comparator metrics and enhanced decision rules in the Fellegi-Sunter model of record linkage," Proceedings of the Section on Survey Research Methods, American Statistical Association., pp. 354–359,1990.
- [12] F. Naumann, J. Freytag, and U. Lesser, "Completeness of Integrated Information Sources", Workshop on Data Quality in Cooperative Information Systems (DQCIS2004), Cambridge, Mass., pp.583–615, 2004.
- [13] P. Angeles, F. Garcia-Ugalde, "Assessing data quality of integrated data by quality aggregation of its ancestors", Computación y Sistemas, Centro de Investigación en Computación, Instituto Politécnico Nacional (IPN), vol. 13 No. 3, ISSN 1405-5546, pp. 331–334, 2010.