

A method and a tool for geocoding and record linkage

Omar CHARIF¹
Hichem OMRANI¹
Olivier KLEIN¹
Marc SCHNEIDER¹
Philippe TRIGANO²

*CEPS/INSTEAD, Luxembourg¹
Heudiasyc Laboratory, Technology University of Compiègne, France²*

CEPS/INSTEAD Working Papers are intended to make research findings available and stimulate comments and discussion. They have been approved for circulation but are to be considered preliminary. They have not been edited and have not been subject to any peer review.

The views expressed in this paper are those of the author(s) and do not necessarily reflect views of CEPS/INSTEAD. Errors and omissions are the sole responsibility of the author(s).

A method and a tool for geocoding and record linkage

Omar Charif^{*‡}, Hichem Omrani[‡], Olivier Klein[‡], Marc Schneider[‡] and Philippe Trigano^{*}

^{*} *Heudiasyc laboratory*

Technology University of Compiègne, France

Email: omar.charif@hds.utc.fr

[‡] *GEODE department, CEPS/INSTEAD, Luxembourg*

Email: hichem.omrani@ceps.lu

Abstract—For many years, researchers have presented the geocoding of postal addresses as a challenge. Several research works have been devoted to achieve the geocoding process. This paper presents theoretical and technical aspects for geolocalization, geocoding, and record linkage. It shows possibilities and limitations of existing methods and commercial software identifying areas for further research. In particular, we present a methodology and a computing tool allowing the correction and the geo-coding of mailing addresses. The paper presents two main steps of the methodology. The first preliminary step is addresses correction (addresses matching), while the second carries geocoding of identified addresses. Additionally, we present some results from the processing of real data sets. Finally, in the discussion, areas for further research are identified.

Keywords—addresses correction, geocodage, matching, data management, record linkage.

I. INTRODUCTION

A lot of research works have been devoted to Geocoding of postal addresses. The interest in this topic is supported by the need to transform postal addresses into geographical coordinates which are essential for various domains of scientific and social research. The benefits of the address geocoding precision are numerous. Geocoding can be used for a wide range of applications such as market segmentation, demographics, geo-spatial distribution of plants, sales territories, taxes, elections. Geocoding is also a very important tool to target certain demographics characteristic. The results of geo-coding have provided fundamental components for wide variety of research works in many fields (e.g. health [4], crime analysis [8], political science [6], computer science [5], etc.). The geocoding operation plays for example an important role for marketing in companies; it helps to cluster peoples with specific characteristics that might be interested in their products.

Many research centre and companies have developed free and commercial geocoder. A big number of these softwares use the linear interpolation method to calculate the spatial coordinates. This method estimates the coordinates of an address using the coordinates of bordering addresses of the street where the address is located. Many research papers [1], [5] have described the error in localization produced by using the linear interpolation method. It was mentioned in [1] that the error in localization can reach 3 kilometres (the distance between the true position and the estimated localization). In addition to the localization

error, a big number of the developed tools do not take into account misspelled and abbreviation errors which are made while writing the postal addresses.

These tools are not able to deal with miswritten addresses such as miswritten road name, city name, etc. After studying some of the existing solutions for geocoding, we decided to develop our own geocoder. The developed tool is able to detect and correct errors as well as to deliver the precise coordinates for input addresses. The structure of the paper is as follows. First, we present a brief overview about geocoding. Second, we describe the developed methods. Subsequently, the Results of processing administrative files are summarised. Finally, we conclude the paper and show some areas for future development.

II. OVERVIEW

Most existing works in the field of geocoding are developed based on the structure shown in Fig. 1. The geocoding process is divided into three main steps:

- 1) Structuring and normalizing: it consist to clean and normalize the input address.
- 2) Record linkage: it allows finding a match of the inputted address in the reference database.
- 3) Geo-coding: it calculates coordinates of the indentified address.

Existing research works usually differ with respect to the methods which are used on each step of the geocoding process. Fig. 2 summarizes methods currently available for each step.

- Structuring and normalizing step: this step is required for cleaning and structuring the input address. The most difficult part of this step is the normalization where each different part of a postal address (postal code, address, road name, etc.) must be identified from a completely input address. Fig. 2 presents details about different methods already used in this step.
- Record linkage phase: It allows comparing names and address information across to pairs of files (or data sets) to find out if they are describing the same entity. It is during this step that errors in writing an address will be detected (methods are shown in fig 2)
- Geocoding: the final step of the process is to calculate the spatial coordinates. This step finds the coordinates while considering the desired scale (see methods shown in fig. 2)

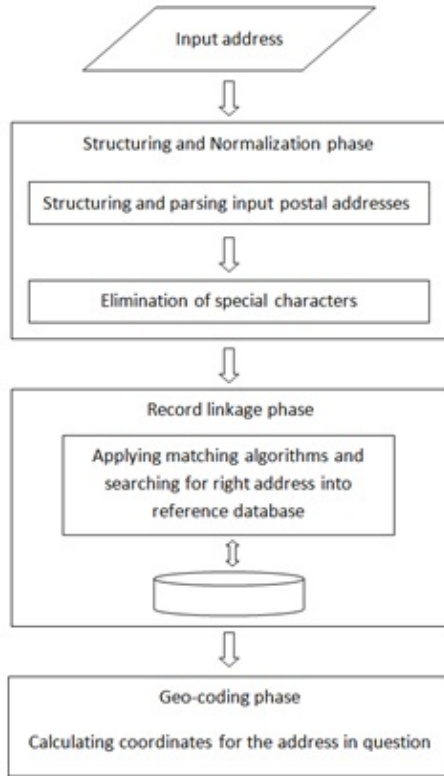


Figure 1. General structure for geocoding process

III. METHODOLOGY

We have developed a general purpose tool for geocoding while taking into account the particularities of our case of study (Luxembourg). Consequently, some parts of our methodological choice have been influenced by both characteristics of Luxembourgish postal address system and the type of files that we were willing to process. Below are the presentation and the justification of our methodological choices:

A. Step 1: Structuring and normalizing

This tool has been developed to process administrative files, where the postal addresses were divided into fields (road name, postal code, municipality, etc.). Thus, we didn't need to perform a complicated normalization technique that parse the input address into fields. In some cases, we used substitution based normalization in order to distinguish two parts of an address that were coded under the same variable (e.g. L-3123 where the letter "L" represents country Luxembourg and 3123 stands for postal code). This method relies on the type of the fields to identify them (e.g. postal code is usually a number and country is a string). It divides the input address into tokens by using "space", "comma", etc as a separator. It will then associate tokens to fields that have the same type. On the other hand, the fact that Luxembourg is a multi language country (i.e. Luxembourgish, French and German) has brought up the need of cleaning (eliminating special character) and standardization step.

B. Step 2: Record linkage

Besides geocoding, we also developed a tool able to correct mistakes produced while inputting data. The decision in this step was very important for the success of the work presented in this paper. Thus, the biggest work lies in the effort to find an algorithm able to detect and correct mistakes while matching the inputted address with addresses in the reference database. Although the first two choices (Match-Merge and deterministic, see fig. 2) were very simple they were not able to deal with complicated misspelled and mistakes. According to Dey [3] the string comparison methods have shown higher reliability than probabilistic methods.

Following the results presented in table I which obtained by applying different String similarity metric methods on one road written in a different way ("AVENUE J.F.KENNEDY" and "AVENUE J-F KENNEDY"), we have noticed that the "Jaro", "Jaro winkler", "Levenshtein", "Mongo Elkan" and "soundex" were the best in detecting misspelling errors.

On the other hand, the results (shown in table II) of comparing two roads with very similar name demonstrate that "Levenshtein" method is more reliable than "Jaro", "Jaro Winkler", "Mongo Elkan" and "Soundex". Thus we decided to combine two techniques of string comparison "Livenshetein distance" [7] and "vectorial" approach (e.g. Q_Grams algorithm [2]). The "Levenshtein distance" calculates the number of operation (i.e. add, remove, substitution) which is needed for passing from one string to another, which helps to detect and correct the misspelled errors. Yet this method is not able to detect abbreviation based errors. This type of errors requires the intervention of the "vectorial technique" which consists of dividing the compared names into tokens or words. Fractions of each matched name will then be compared (by comparing the two words using "Levenshtein distance" or just by comparing first letters of the two words) while a percentage of similarity is calculated. These choices were made by considering the processing time and the reliability of the similarity metric results (tables I and II).

The matching procedure begins with verifying the existing of the couple postal code and road name by querying the reference database. If the answer of this query is null we always assume that the postal code is correct. The reason for this is because errors are most likely to be committed while inputting text data. We then create a list of roads which are associated with the input postal code. An algorithm is then executed to match the input address with captured road list. If the matching did not succeed then the same procedure is repeated but with a list of roads associated to the input Municipality. In order to accelerate the processing time, we have created a knowledge database which helps to memorize the variants of names writing (errors already detected). This knowledge database becomes richer as we run a new file process.

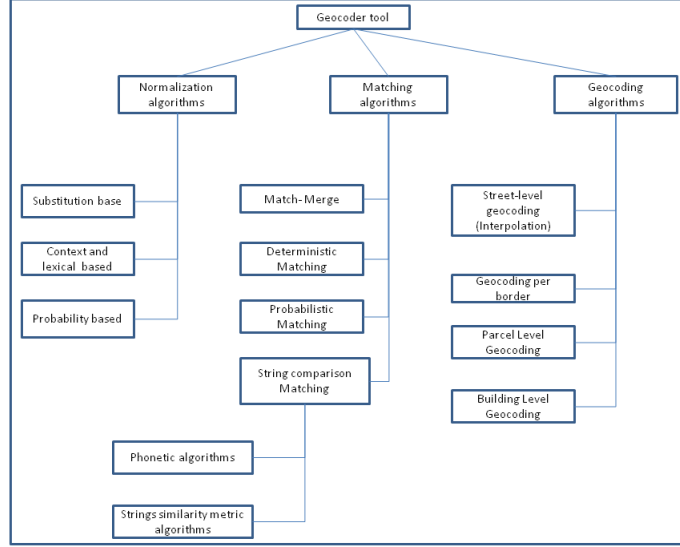


Figure 2. Methods used in geocoding process

	Jaro	Jaro Winkler	Levenshtein	Mongo Elkan	QGrams	Jaccard	Soundex
Similarity index	0.925	0.970	0.888	1.0	0.75	0.25	1
Processing time	0.013	0.014	0.058	0.86	0.043	0.0008	0.020

Table I

SIMILARITY CALCULATION RESULTS FOR MATCHING "AVENUE J.F.KENNEDY" AND "AVENUE J-F KENNEDY"

	Jaro	Jaro Winkler	Levenshtein	Mongo Elkan	QGrams	Jaccard	Soundex
Similarity index	0.893	0.957	0.75	0.888	0.571	0.5	1
Processing time	0.009	0.010	0.043	1.238	0.032	0.001	0.017

Table II

SIMILARITY CALCULATION RESULTS FOR MATCHING "RUEDESARDENNES" AND "RUEDESJARDINS"

	Total record	Geocoding percentage	Missing data percentage	Building geocoding percentage	Geocoding nearest Neighbor percentage	Geocoding by road barycentre
Data set 1	19409	98.20%	0.015%	88.83%	10.4%	0.45%
Data set 2	35594	97.51%	0.073%	86.32%	12.31%	1.37%
Data set 3	38672	81.98%	16.97%	84.69%	13.43%	1.88%
Data set 4	457339	95.43%	0.024%	97.87%	1.68%	0.45%

Table III

RESULTS OF PROCESSING FOUR DATA SETS FROM DIFFERENT ADMINISTRATIVE SOURCE

C. Step 3: Geo-coding

It has shown in [9] that the quality of geocoding has a big influence on the result of analyses which use it. According to [1], the error in localization which is produced using parcel geocoding method is significantly smaller than the error in localization which is produced by using street geocoding method. These two facts and our disposal of a database which contains coordinates for buildings have encouraged us to use the building localization method. In case the reference data base does not contain the input building, we associate to the input address the coordinates presented in equations 1. We have combined the building geocoding and some kind of linear interpolation to calculate these coordinates. We called this method "geo-coding by nearest-neighbour".

$$\begin{cases} X = X_{n_n} + (n_i - n_n) \times \overline{distance} \times \cos(\Theta) \\ Y = Y_{n_n} + (n_i - n_n) \times \overline{distance} \times \sin(\Theta) \end{cases} \quad (1)$$

Where: n_i : the input building number

n_n : the nearest neighbours building number

X_{n_n} : $X_{nearest_neighbour}$: longitude of the nearest neighbour (in term of address building number) form the same side

Y_{n_n} : $Y_{nearest_neighbour}$: latitude of the nearest neighbour (in term of address building number) form the same side

$$\overline{distance} = \frac{\sum_{i=2}^n \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}}{(n-1)}$$

where n represents the number of building exist on the side of inputted address.

Algorithm 1 Calculate angle

```

if ( $Y_{first\_address} < Y_{last\_address}$ ) then
  if ( $X_{first\_address} < X_{last\_address}$ ) then
     $\Theta \leftarrow \arctan \frac{(Y_{last\_address} - Y_{first\_address})}{(X_{last\_address} - X_{first\_address})}$ 
  else
     $\Theta \leftarrow 180 - \arctan \frac{(Y_{last\_address} - Y_{first\_address})}{(X_{last\_address} - X_{first\_address})}$ 
  end if
else
  if ( $(X_{first\_address} < X_{last\_address})$ ) then
     $\Theta \leftarrow 360 - \arctan \frac{(Y_{last\_address} - Y_{first\_address})}{(X_{last\_address} - X_{first\_address})}$ 
  else
     $\Theta \leftarrow 180 + \arctan \frac{(Y_{last\_address} - Y_{first\_address})}{(X_{last\_address} - X_{first\_address})}$ 
  end if
end if

```

with:

- $X_{first_address}$: longitude of the building which has the smallest address number from the same side as the inputted address
- $Y_{first_address}$: latitude of the building which has the smallest address number from the same side as the inputted address

IV. RESULTS

We present in table III, the result of processing of four data sets from different administrative sources. The first three are results for geocoding data sets containing addresses of six test municipalities in Luxembourg. The fourth is the result of process a data set that contains addresses from all over Luxembourg. The developed tool contains interactive, friendly user interfaces which facilitate the setup of settings needed for data sets processing as shown in figs 3 and 4.

V. CONCLUSION

In this paper, we have presented two new methods. The first one is for record linkage and the second is for coding. These two methods have given good results with a more than 95% percentage of success. We have implemented and developed the computation tool using Java programming language. In the future, a normalization of the input address module however must be added to this tool.

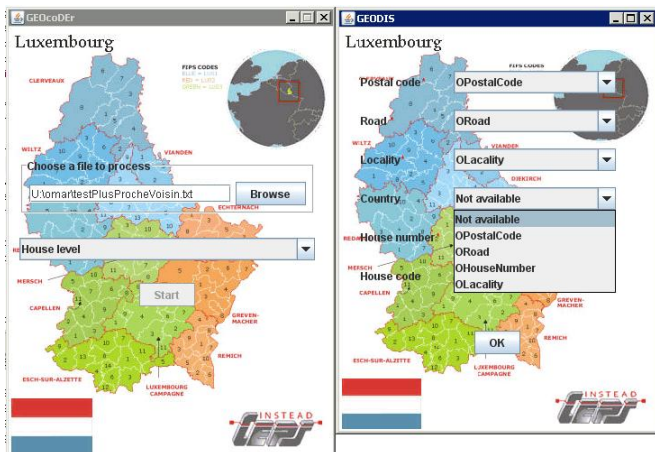


Figure 3. Select file and setup processing settings

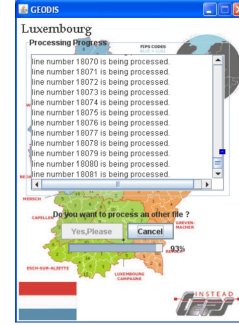


Figure 4. Processing progress user interface

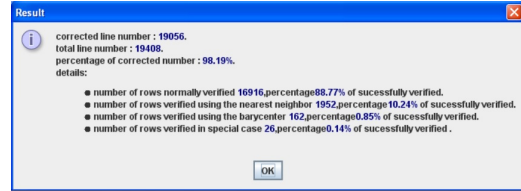


Figure 5. Result user interface

ACKNOWLEDGMENT

The work reported in this paper was partially supported by CEPS/INSTEAD research centre and GEODIS project. The authors would like to thank Dr. Philippe Gerber for his comments and two anonymous reviewers for their useful and constructive comments on the earlier version of this paper. Any errors in this paper are the responsibility of the authors.

REFERENCES

- [1] M.R. Cayo and T.O. Talbot. Positional error in automated geocoding of residential addresses. *International journal of health geographics*, 2(1):10, 2003.
- [2] T. Churches and P. Christen. Some methods for blindfolded record linkage. *BMC Medical Informatics and Decision Making*, 4(1):9, 2004.
- [3] D. Dey et al. A distance-based approach to entity reconciliation in heterogeneous databases. *IEEE Transactions on Knowledge and Data Engineering*, 14(3):567 – 582, 2002.
- [4] G. Rushton et al. Geocoding in cancer research : a review. *American Journal of Preventive Medicine*, 30(2):16–24, 2006.
- [5] R. Bakshi et al. Exploiting online sources to accurately geocode addresses. In *ACM-Gis*, page 194203, 2004.
- [6] M. Haspel and H.G. Knotts. Location, location, location: Precinct placement and the costs of voting. *The journal of politics*, Cambridge University Press, 67:560–573, 2005.
- [7] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics Doklady*, 10:707, 1966.
- [8] J.H. Ratcliffe. Geocoding crime and a first estimate of a minimum acceptable hit rate. *International Journal of Geographical information Science*, 18(1):6172, 2004.
- [9] P.A. Zandbergen. Influence of geocoding quality on environmental exposure assessment of children living near high traffic roads. *BMC Public Health*, 7(37):1–13, 2006.

CEPS
I N S T E A D

B.P. 48
L-4501 Differdange
Tél.: +352 58.58.55-801
www.ceps.lu