# Repopulating Paris: massive extraction of 4 Million addresses from city directories between 1839 and 1922.

*Isabella di Lenardo (isabella.dilenardo@epfl.ch), DHLAB-EPFL, Switzerland and Raphaël Barman (raphael.barman@epfl.ch), Institut National d'Histoire de l'Art and Albane Descombes (albane.descombes@epfl.ch), Institut National d'Histoire de l'Art and Frédéric Kaplan (frederic.kaplan@epfl.ch), DHLAB-EPFL, Switzerland*

## 1. Introduction

In 1839, in Paris, the Maison Didot bought the Bottin company. Sébastien Bottin trained as a statistician was the initiator of a high impact yearly publication, called "Almanachs" containing the listing of residents, businesses and institutions, arranged geographically, alphabetically and by activity typologies (Fig. 1). These regular publications encountered a great success. In 1820, the Parisian Bottin Almanach contained more than 50 000 addresses and until the end of the 20th century the word "Bottin" was the colloquial term to designate a city directory in France. The publication of the "Didot-Bottin" continued at an annual rhythm, mapping the evolution of the active population of Paris and other cities in France.
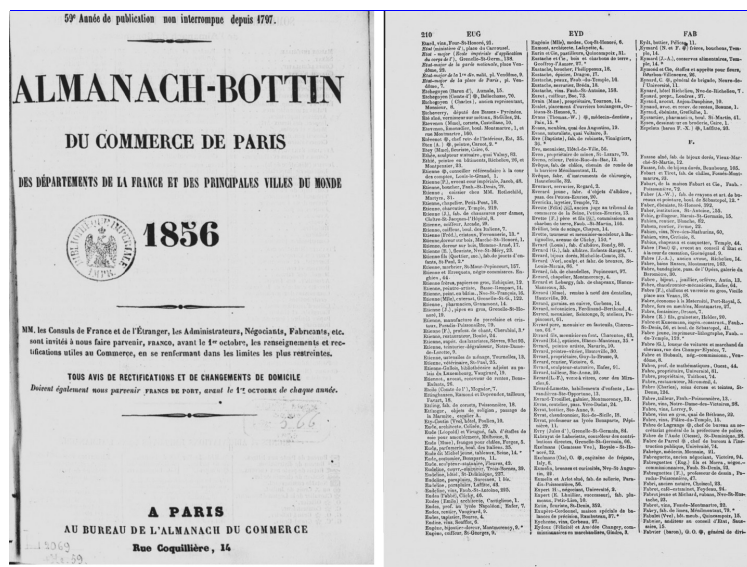


Figure 1: Example of the cover page and alphabetical list for the year 1856

The relevance of automatically mining city directories for historical reconstruction has already been argued by several authors (e.g Osborne, N., Hamilton, G. and Macdonald, S. 2014 or Berenbaum, D. et al. (2016). This article reports on the extraction and analysis of the data contained in "Didot-Bottin" covering the period 1839-1922 for Paris, digitized by the Bibliotheque nationale de France. We process more than 27 500 pages to create a database of 4,2 Million entries linking addresses, person mention and activities. The quality of the document analysis process is assessed diachronically and a conservative strategy was chosen in order to populate the database with only information of high confidence. An initial analysis of the data is presented, reporting on the overall statistics of the distribution of professions in Paris and their evolution during more than 80 years, as well a general overview of the diversity of family names through time. Seven case studies corresponding to different streets are briefly compared, showing how information in city directories capture statistically the dynamics of segmentation of the city into functionality differentiated neighborhoods.
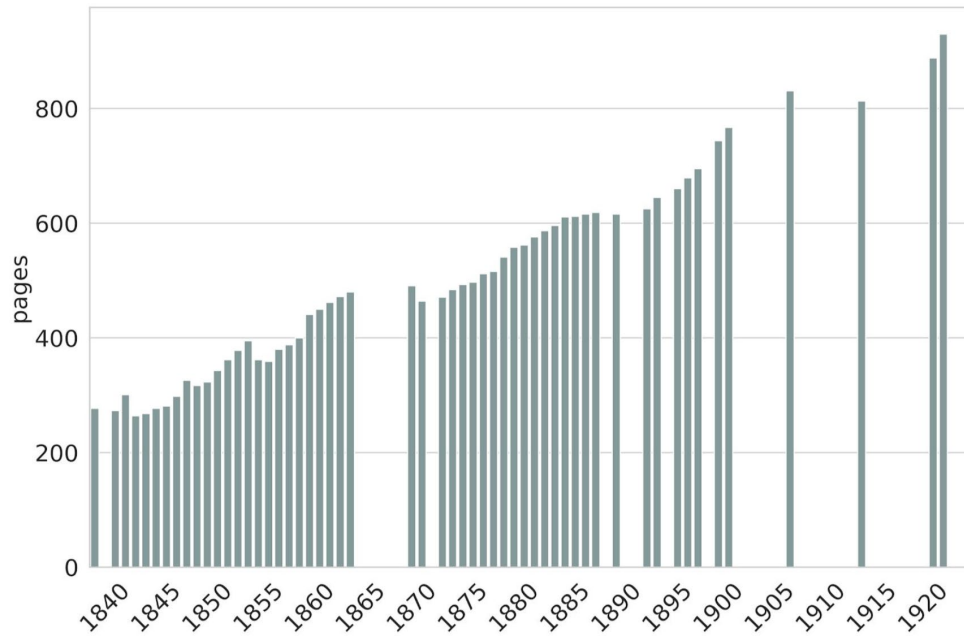
Figure 2: Number of pages per city directory year by year. 29 years are missing in the digital collection used for this study.

## 2. Method

The document analysed in this article were digitized by Bibliotheque nationale de France and made available online through the Gallica portal. The dataset corresponds to three different published series which are homogeneous in their structure and aims (Bing 1897): Annuaire général du commerce (1839-1856); Annuaire-Didot-Bottin (1857-1907); Annuaire du commerce Didot-Bottin (1914-1922). The documents were associated with an ALTO description containing a structural decomposition of each page into text blocks and lines, associated with a transcription obtained by an Optical Character Recognition (OCR) process. We designed a parsing process converting each line/entry into a record in a database documenting the name, the activity, the place and when relevant the street number (Fig. 3). Only a subset of the entries was successfully parsed (4,2M over 5,6M) and included in the database.



Figure 3: Structure of an entry in the directories

A general discussion on quality assessment methods of the OCR for the BnF digital collection can be found in ( Salah - Moreux - Ragot - Paquet 2015) . In order to assess specifically the OCR quality of the inserted data, 14 pages were randomly picked for the years 1839, 1848, 1856, 1857, 1881, 1907, 1921 and manually controlled.
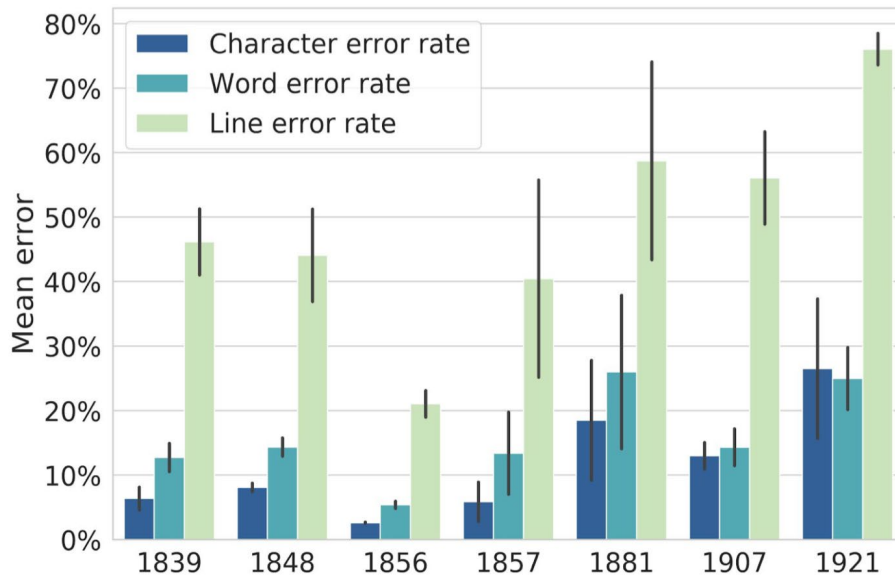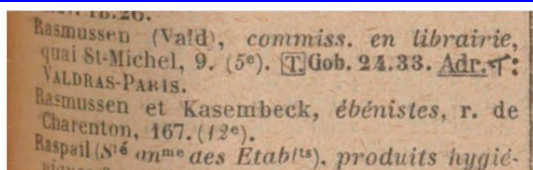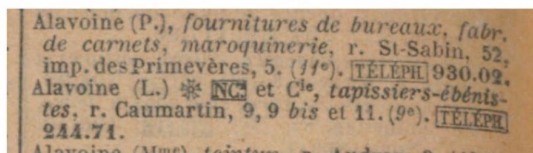
Figure 4: Character error rate, Word error rate and Line error rate for 7 years of the corpus (14 random pages analysed).
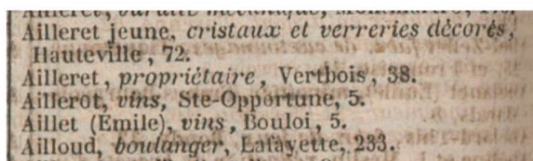
On well parsed entries, the mean of the character error is of 2.6% with a standard deviation of 0.1%. The error per line is of 21% with a standard deviation of 2.9% (Fig. 4). The overall transcription quality tends to decrease as sources are more recent. This is essentially due to three factors: (1) the insertion at the beginning of the 20th century of special symbols used for saving textual space but unparsed by most OCR system, (2) the increasing thickness of the volumes leading to the curvature issue during the scanning process making line detection and word identification more difficult (3) the use of continuously thinner paper sheets leading to problems of transparency between the verso and the recto of a page (Fig. 5).



Figure 5: Examples of problems of scanning (1) symbols, (2) page straightening (3) transparency and corresponding OCR

The entries of the database for seven cases studies were realigned on the Vasserot cadastre digitized and analysed during the ALPAGE project (Noizet-Bove-Costa 2013) and available online. The Vasserot cadastre is giving a full director of addresses in Paris for the period 1810-1837. Paris Open Data covers the structural change due to the Haussmann period and the evolution of the 20th century. Using these two sources, 89,2% of the addresses were successfully realigned for the seven case studies considered.

## 3. Results

In total 4.2M person mentions were extracted for the period 1839 to 1922. This database could be the starting point of numerous studies, we are only giving here a broad illustration of the content of this dataset and discussing their potential relevance for future research. For instance, the diversity of family names, an indirect proxy of the social effects of urbanization, clearly increases during the 19th century and then stabilizes at the beginning of the 20th century (Fig. 6).
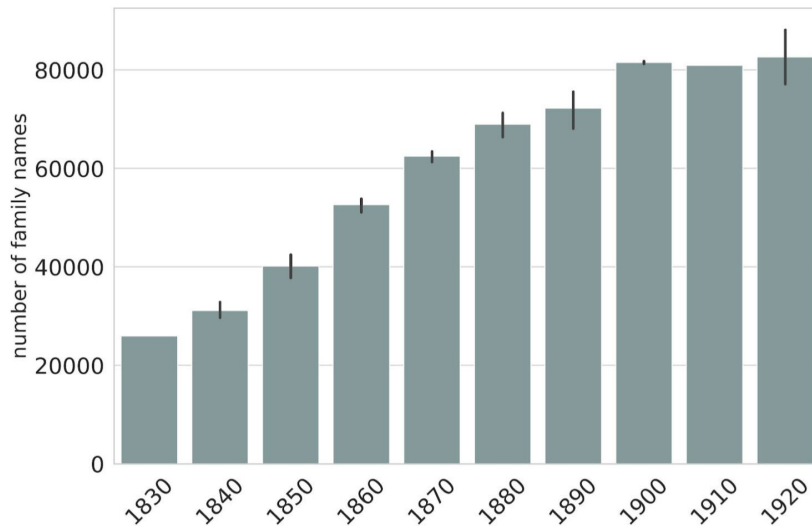
Figure 6: Number of unique family names by year during the entire period. Name diversity keeps increasing during the 19th century and then stabilizes.
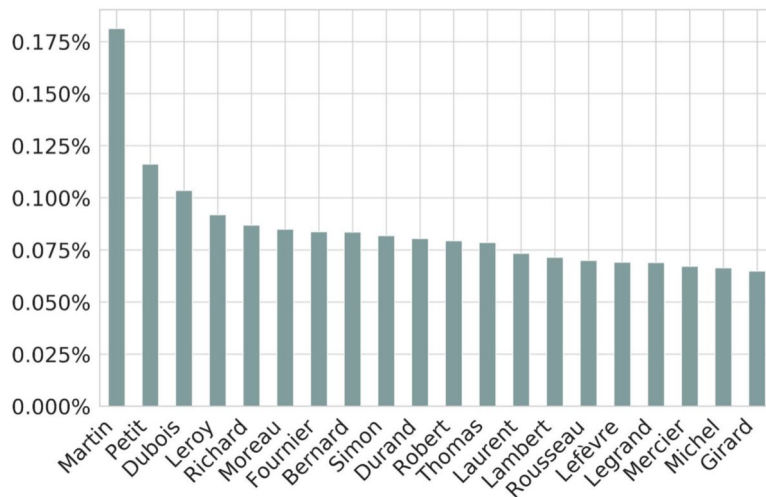


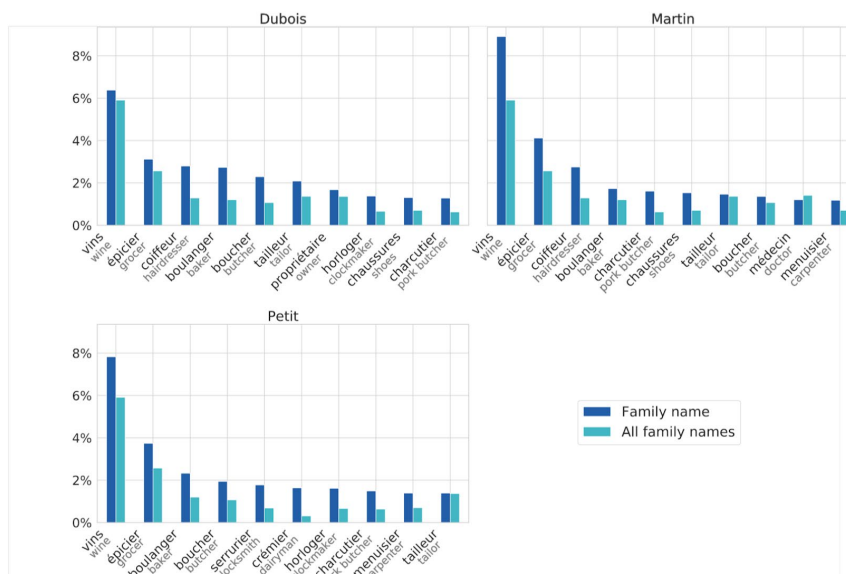Figure 7: Family names frequency on the period 1839-1922



Figure 8: Frequency of trade type by family names "Dubois" ; "Martin" ; "Petit" on the period 1839-1922

If one compares the three most frequent family names "Martin", "Petit", "Dubois" (Fig. 7) with the typologies of trading businesses, the sale of wine appears to be, by far, the most commonly practiced activity (Fig. 8). If you encounter a Martin, a Petit or a Dubois at the end of the 19th century, there is a certain chance, he'll be a wine seller. This result corresponds to a general trend, at the urban scale. In the figures 9 and 10 the dominance of the wine business is 2.5 times more important than the one of grocery stores, the second activity on the list. As confirmed by other studies (Pinol-Garden, 2009), the relative proportion of the wine activity keeps increasing during the 19th century. This is in line with other figures like the wine consumption (one million hectolitre in 1800, 5 million in 1920, about 150-200 litres per person/per year) and the construction of the Bercy wine hall in 1869 (42 hectares on the Seine bank) and their extension in 1910 (Gallet 1939; Thillay-Reynald 2004). Grocery professions, tailor, hairdressers, bakers are equally represented, alongside liberal professions like doctors. Architects, cabinet makers and carpenter are also well represented, a sign of an important building activity on a city scale.
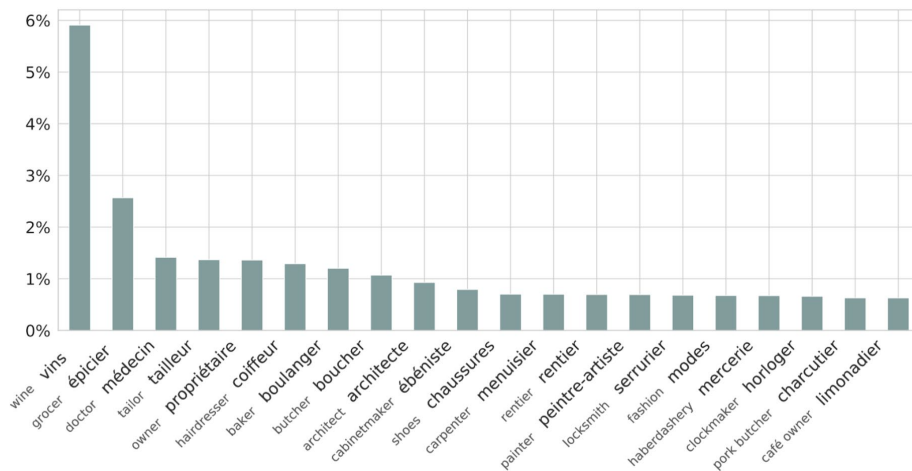


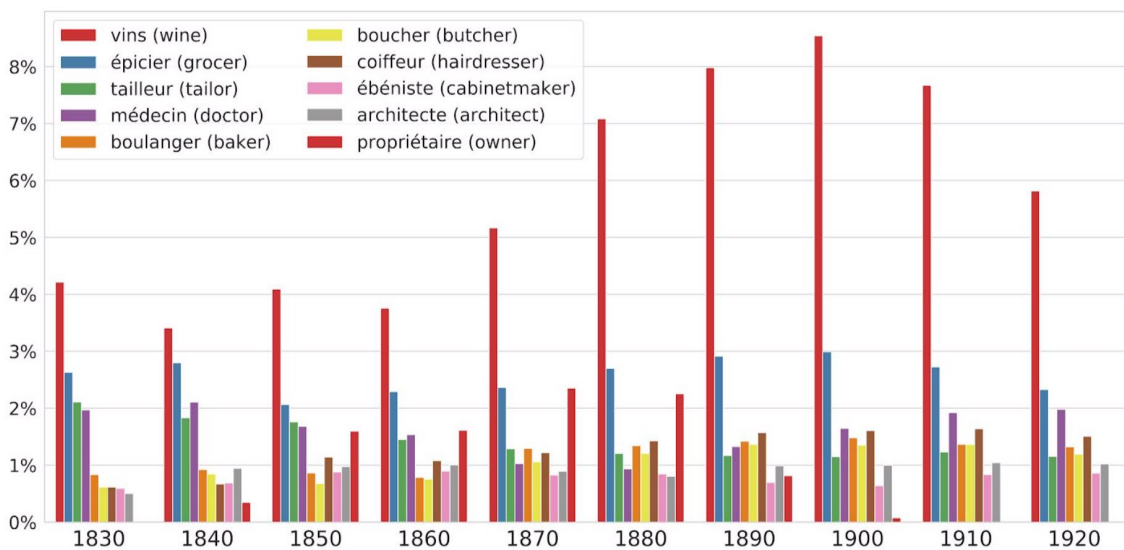Figure 9: most common trades over the whole period



Figure 10: most common trades by decade

The seven chosen case studies correspond to a specific analysis crossing historical information and urban space, focusing on very ancient streets (Rue Montmartre-Rue Saint Denis-Rue Saint Antoine-Rue Saint Jacques) organising a variety of activities over a deep urban "palimpsest" (Corboz 2001) and others, much more recent, with a strong characterization by trade businesses (Rue Richelieu-Rue, Saint Anne) or residential vocation (Boulevard Saint Michel) (Fig. 11). Fig 12 compares the situation of seven cases studies showing the diversification of activities by neighborhood. The ancient roman streets (Rue Saint-Denis, Rue Saint-Antoine, Rue Saint-Jacques) are characterized by the equally distributed presence of services for the population, reflecting a residential activity and a mix of activities acquired over time. Only the mention of the 'bookbinders' on the Rue Saint Jacques, highlights the presence of the university district strongly linked to book production. Rue Montmartre, mentioned since the 13th century, includes many more activities than other cases, without clear differentiators, except for the presence of the tailors. Indeed, this axis is the interface to the 'Richelieu District' considered as a place of production of fabrics and clothes. Rue Richelieu and Rue Saint Anne have a concentration of activities linked with fashion

("tailleurs", "couturière", "bottiers"). This denotes the original vocation not residential but specific to these trades. The distribution of activities on the Boulevard Saint-Michel constructed after the transformation from Baron Haussmann is significantly different from all the others with higher level of owners and liberal professions (doctors, architects).



Fig. 11: The seven case studies: *Rue Montmartre; Rue Saint-Denis; Rue Saint-Antoine; Rue Saint-Jacques; Rue Richelieu; Rue Saint Anne; Boulevard Saint-Michel*
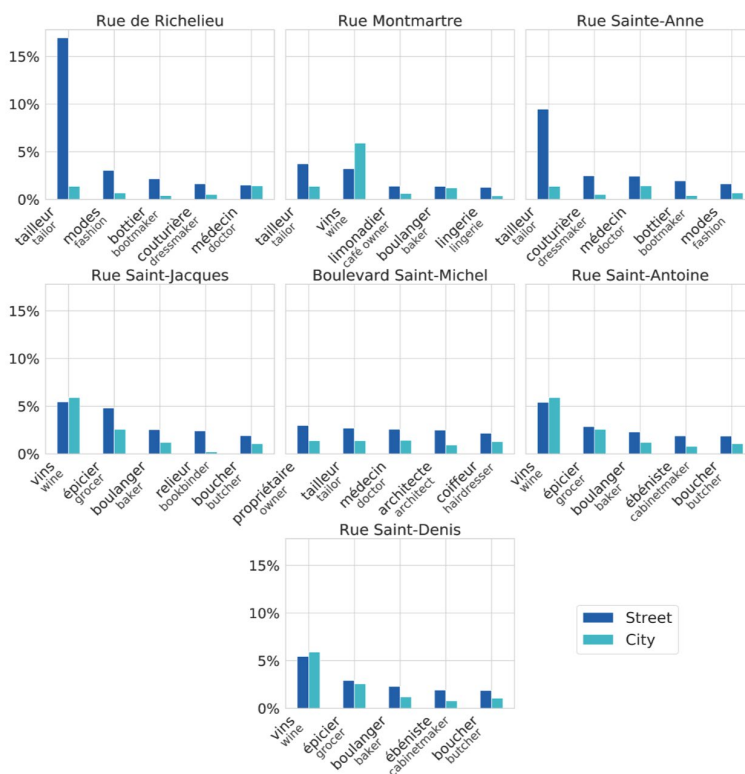
Fig. 12: Street trades (1839-1922)

These preliminary results on Paris hope to demonstrate the potential of city directories to conduct large-scale urban analysis at different level of granularity. The automatic extraction process designed for this article permits to envision to easily conduct similar studies on the population of many other important cities in the world. Provided that the quality of the extraction process is monitored, such kind of massive datasets will open new avenues to study the transformations of the urban structure at different geographical and temporal scales during the ongoing industrialization and other significant societal transformations of the 19th century, connecting these large datasets from the past with the ones of the present.

## Appendix A

Bibliography

1. Ben Salah, A., Moreux, J.-P., Ragot, N. and Paquet, T. (2015). OCR performance prediction using cross-OCR alignment. 13th International Conference on Document Analysis and Recognition (ICDAR), pp. 556Ð60

2. Berenbaum, D., Deighan, D., Marlow, T., Lee, A., Frickel, S. and Howison, M. (2016). Mining Spatio-temporal Data on Industrialization from Historical Registries. Journal of Environmental Informatics. d oi:10.3808/jei.201700381 https://arxiv.org/abs/1612.00992

3. Bing, A.-D. A. (1897). Les annuaires parisiens, de Montaigne à Didot, 1500 à 1900. Le Havre.

4. Corboz, A. (2001), Le territoire comme palimpseste et autre essais, Payot, 2001 Gallet, P. (1939). LÕapprovisionnement en vin de Paris. Annales de géographie , vol. 48, no. 274. Paris, pp. 359-368.

5. Noizet, H., Bove, B. and Costa, L. (2013). Paris de parcelles en pixels . A nalyse géomatique de l'espace parisien médiéval et moderne . Paris : Presses universitaires de Vincennes-Comité d'histoire de la Ville de Paris.

6. Osborne, N., Hamilton, G. and Macdonald, S. (2014). Historical Post Office Directory Parser (POD Parser) Software From the Addressing History Project. Journal of Open Research Software , vol. 2, no.1, p. 23.

7. Pinol, J.-L. and Garden, M. (2009). Atlas des Parisiens. De la révolution à nos jours. Paris.

8. Thillay, A. and Reynald, A. (2004). Le grand marché. L'approvisionnement alimentaire de Paris sous l'Ancien Régime. Histoire, économie et société , vol. 23, no. 2 : La société, la guerre, la paix, 1911-1946. pp. 307-308.