

# A Method for Record Linkage with Sparse Historical Data

Giovanni Colavizza, Maud Ehrmann, Yannick Rochat

Digital Humanities Laboratory (DHLAB)  
Swiss Federal Institute of Technology in Lausanne (EPFL)  
CDH, INN 116, Station 14, Lausanne, Switzerland

{name.surname}@epfl.ch

## 1. Introduction

Massive digitization of archival material, coupled with automatic document processing techniques and data visualisation tools offers great opportunities for reconstructing and exploring the past. Unprecedented wealth of historical data (e.g. names of persons, places, transaction records) can indeed be gathered through the transcription and annotation of digitized documents and thereby foster large-scale studies of past societies. Yet, the transformation of hand-written documents into well-represented, structured and connected data is not straightforward and requires several processing steps. In this regard, a key issue is entity record linkage, a process aiming at linking different mentions in texts which refer to the same entity. Also known as entity disambiguation, record linkage is essential in that it allows to identify genuine individuals, to aggregate multi-source information about single entities, and to reconstruct networks across documents and document series.

In this paper we present an approach to automatically identify coreferential entity mentions of type *Person* in a data set derived from Venetian apprenticeship contracts from the early modern period (16th-18th c.). Taking advantage of a manually annotated sub-part of the document series, we compute distances between pairs of mentions, combining various similarity measures based on (sparse) context information and person attributes.

## 2. Task Definition

Major challenges when dealing with people-related data are homographic person names referring to different persons as well as the existence of name variants referring to the same person. These are well-known issues in the field of Natural Language Processing for which various approaches have been devised, first via mention clustering [Mann and Yarowsky 2003, Artiles et al. 2008], more recently via linking to a knowledge base [Ji and Grishman 2011, Shen et al. 2015].

In the context of historical data, dealing with person name ambiguity is all the more difficult since data is inherently sparse and uncertain (resulting in poor mention context) and since knowledge bases such as DBpedia [Lehmann et al. 2013] contain very little about past average laypersons (resulting in poor entity context). It is however an essential step prior to any historical data analysis [Bloothoof et al. 2015], which we address as part of the *Garzoni* project. This project aims at studying apprenticeship in early modern Venice by extracting information from archival material. Part of this material have been manually annotated, including mention links towards unique entities. Starting from a subset of the current data, we present a method for person record linkage, with the objective to complement its disambiguation coverage and to bootstrap a system to better automate entity disambiguation during annotation, in an active learning fashion.

count	whole period	1586-1600
# annotated contracts	11,525	2,687
# mentions	31,952	7,589
# entities	26,641	6,599
<b># entities with # mention &gt; 1</b>	1793	382
AVG mention per entity	1.09	1.08
AVG mention per entity with # mention > 1	2.44	2.38

**Table 1. Entity-Mention statistical profile for the whole vs. selected period.**

### 3. The *Accordi dei Garzoni*

The *Accordi dei Garzoni* is a document series from the State Archives of Venice which originates from the activity of the *Giustizia Vecchia* magistracy. This judicial authority was in charge of registering apprenticeship contracts in order to protect young people while they were trained and/or providing domestic services [Bellavitis 2006]. As a result of this regulation, information for much of apprenticeship arrangements got centralized, today reflected in a dense archival series.

The *Accordi* consists of about 55,000 contracts registered from 1575 until 1772. Each contract features an apprentice, a master and often a guarantor, sometimes two. A sample of 11,000 contracts have been manually annotated and the resulting data is stored in an RDF triple store. For each person mentioned in a contract, annotators created a *person mention* and, importantly, linked it to a *person entity*. They did so either by selecting an already existing entity in the database or by creating a new one. Given the difficulty of this task, only a limited number of entities have been disambiguated; the annotated dataset can therefore be considered as correct but not exhaustive.

The present work considers annotated documents from the period 1586-1600, for which statistics about contracts and entity/mention ratio are shown in Table 1. We use a subset of this dataset (bolded line in the table) as a *golden* set for our experiments.

### 4. Approach

Given a set of mentions, our objective is to estimate the likelihood that two mentions refer to the same entity. We represent each mention by a vector of features and compare them pairwise using various similarity measures. The list of selected features at mention and contract levels are presented in Table 2 and Table 3 respectively.

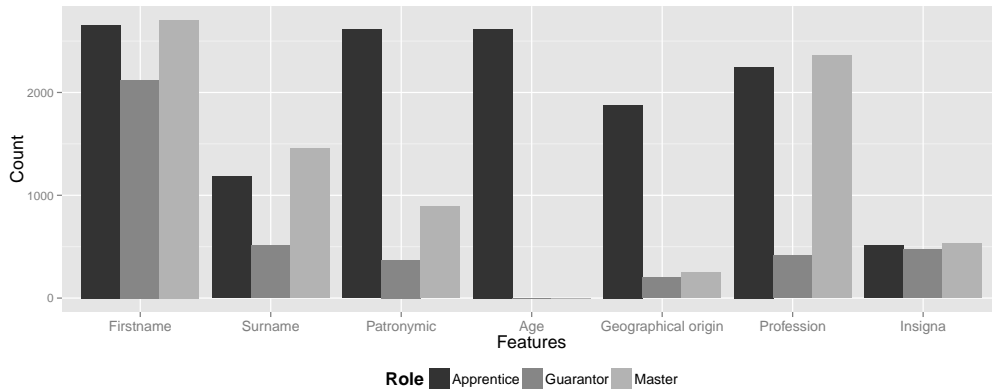
With respect to our dataset and features, several points should be emphasized. First, data sparsity: it is common for a mention to have just a few features. Second, features are not evenly sparse (cf. Figure 1) and do not contribute equally to a possible linkage. Core features such as *name*, *surname*, *patronymic*, *gender* and *profession* must strongly correspond in order to consider a link as reliable. On the other hand, rare features such as *workshop insigna* can be very informative when shared by two mentions and should also be valued by the linkage algorithm. Finally, features are dependent, particularly on the role of the person (e.g. age indicated only for apprentices).

Feature	Variable type
firstname	string
surname	string
patronymic	string
gender	categorical
age	integer
profession	categorical
geographical origins	string

**Table 2. Mention-level features**

Feature	Variable type
workshop toponym	string
workshop parish	string
workshop sestriere <sup>1</sup>	string
workshop insigna	string
contract year	integer
contract duration	string
master profession	categorical

**Table 3. Contract-level features**



**Figure 1. Distribution of features by role.**

We construct three matrices of size  $N \times N$ , where  $N$  is the number of mentions in the dataset. The first matrix  $\Phi$ , the *feature matrix*, stores similarity scores of mentions pairwise. Scores are computed using measures over features as follows:

- *year of contract*: the feature-score is measured via distance and diminishing returns. Each year of distance between 1 and 15 and between 15 and 30 decreases an initial score of 1 by 0.01 and 0.025 respectively, with a definitive cut-off after 30. For example, two contracts from 1590 and 1594 have a score of 0.96.
- *age*: similarly as per year, each year of distance of the difference between two ages decreases an initial score.
- *gender* and *profession*: the feature-scores are based on exact matches.
- *name*, *surname*, *patronymic* and *workshop toponym*: the feature-score is based on the Deverau-Levenshtein string metric[Cohen et al. 2003]. For example, *Polo* and *Pollo* have a similarity measure of 0.95.
- *geographical origins* and *insigna*: the feature-score is based on a token-based variant of the Jaro-Winkler metric. For example, *Friulano* and *del Friuli* have a similarity measure of 0.82.

The score of each pair is stored in  $\Phi$ : it is the L2 norm of the resulting feature-score vector.

<sup>1</sup>There are 6 *sestrieri* in Venice, i.e. groups of contiguous parishes.

The second matrix  $\Gamma$ , the *combination matrix*, stores values that indicate whether a pair of mentions shares similar feature combinations or not. To build such matrix, we leverage the *golden* set and identify combinations of features which produced a linkage on a role-by-role basis (e.g. master-master or guarantor-master). Features are considered activated when their feature-score is equal or above 0.84<sup>2</sup> and we filter out combinations occurring once. The score of a mention pair in  $\Gamma$  is 1.0 if the combination of activated features is valid for the given role pair; 0.5 if the role pair does not match but the combination is valid; 0.0 otherwise. This matrix accounts for feature dependencies and the different ways to name a person with respect to his/her role.

The third matrix  $\Delta$ , the *filtering matrix*, scores mention pairs according to the number of activated core features (1.0 if 3+ features – out of 5 – are activated, 0.0 otherwise<sup>3</sup>).

Given the three matrices, we normalize them and consider the following function to compute the similarity score of a mention pair  $p$ :

$$S(p) = \delta_p[\lambda\pi_p + (1 - \lambda)\gamma_p]$$

where  $\delta_p$  is a boolean parameter taken from  $\Delta$  which activates the filter over core features for pair  $p$ ;  $\pi_p$  is the feature score taken from  $\Phi$ ;  $\gamma_p$  is the combination score from  $\Gamma$ ; and  $\lambda$  is a parameter giving priority over vector features or combinations of features.  $\delta \in \{0, 1\}$  and  $0 \leq \lambda, \pi, \gamma \leq 1$ . This function allows us to adjust the different parameters: core vs sparse features ( $\delta$ ), feature scores ( $\pi$ ) and feature combinations ( $\gamma$ ).

## 5. Evaluation

We evaluate our approach in terms of coverage and precision. With respect to coverage, we verify our method over 100 thresholds from 0.99 to 0.0. For each threshold, we compare linkage curves as the percentage of links obtained over the total possible against the coverage of the *golden* set. Precision is based on manual annotation of 50 randomly selected linkages.

Both procedures are repeated with  $\lambda \in \{0.1, 0.5, 0.9\}$  and  $\delta$  activated or not, for a total of 6 configurations. The objective is to understand the individual contributions of the three components to our function.

## 6. Results and Discussion

Results for the first and second evaluation procedure are presented in Figure 2 and Table 4 (resp.). Highest precision (0.61 and 0.3 in Table 4) is obtained with a balance between feature combinations and feature scores ( $\lambda = 0.5$ ).  $\delta$  proves very useful for filtering the input space (from 28,7M possible pairs to 44,2K), and lowers the number of false positives, especially for links between apprentices (cf. line ‘w-o AA’ in Table 4). The combination of the two (filtered input space and equal weights) provides the best results, especially for masters and guarantors. Linkage curves can be explained similarly: low  $\lambda$  entails a step-like curve (three steps at 0.0, 0.5 and 1.0), while high  $\lambda$  creates a Gaussian over the disambiguation space.

<sup>2</sup>It has been shown in comparable settings that edit distance with cut-off at distance 3 (which for us is distance  $\geq 0.85$ ) provides good results [Kleanthi et al. 2015].

<sup>3</sup>Features are activated when their similarity is above 0.84.

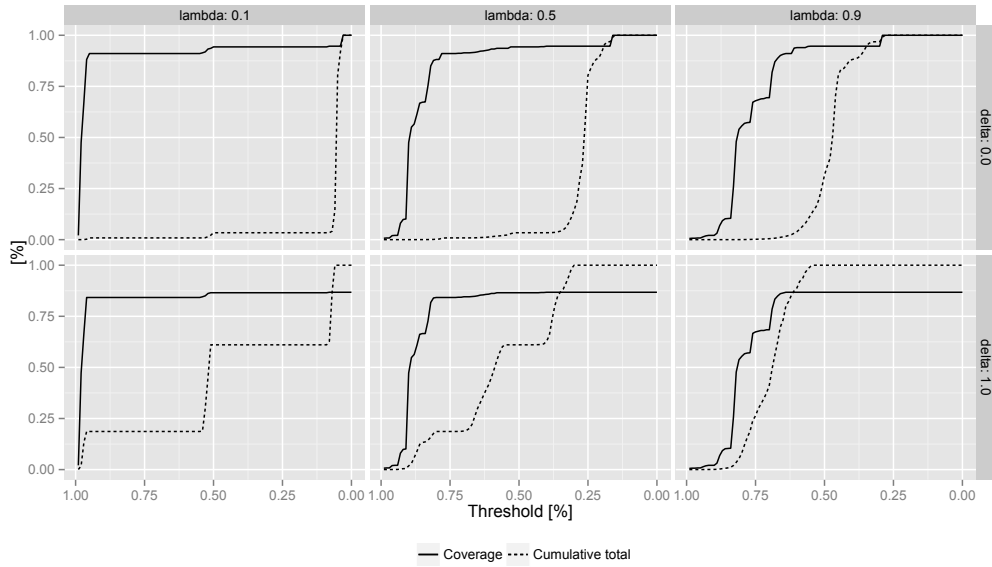


Figure 2. Linkage curves for the 6 parameter settings, over thresholds.

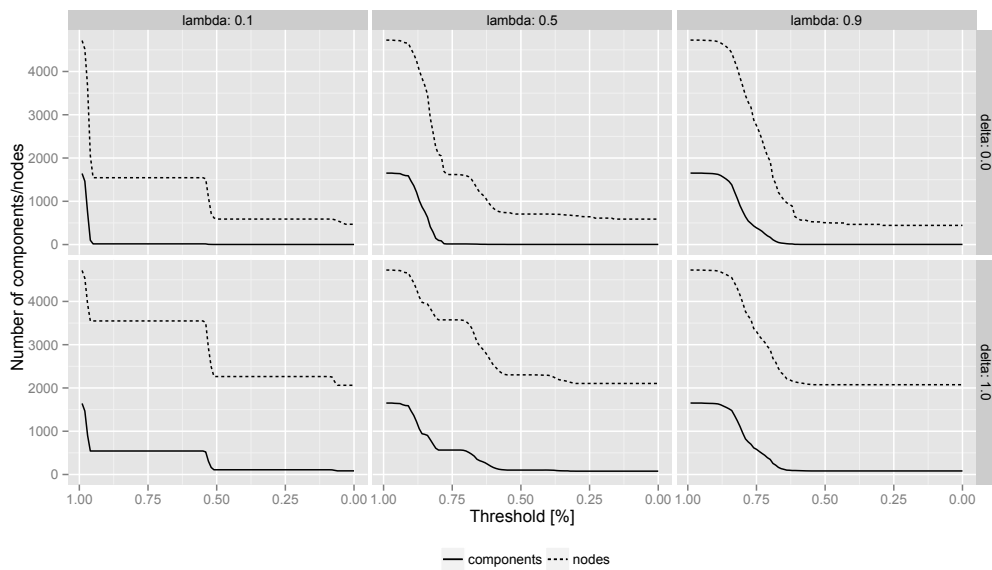


Figure 3. Graph properties for the 6 parameter settings, over thresholds.

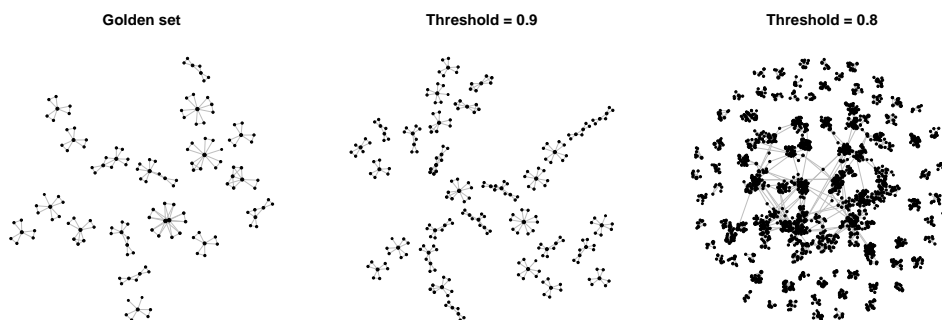
This confirms that a balanced approach might be the best solution in a setting where data is sparse (high  $\lambda$ ), the *golden* set is present but of limited coverage (low  $\lambda$ ), and some prior assumptions on the required features can be made ( $\delta$ ). As shown in Figure 3, the graphs with  $\lambda = 0.5$  and  $\delta = 1.0$  collapse more gradually, providing the widest effective linkage space to explore. Eventually, results also suggest to proceed in an active learning fashion, where the system learns iteratively with new data as part of the *golden* set.

Finally, in order to further motivate our work, Figure 4 shows the largest components of the deduced social network with and without automatic disambiguation. The

	$\delta$ active			$\delta$ not active		
	$\lambda = 0.1$	$\lambda = 0.5$	$\lambda = 0.9$	$\lambda = 0.1$	$\lambda = 0.5$	$\lambda = 0.9$
all	0.21	<b>0.3</b>	0.21	0.0	0.26	0.15
w-o A-A	0.22	<b>0.61</b>	0.22	0.0	0.48	0.67*

**Table 4. Precision with threshold  $\geq 0.9$  (\* means not-significant statistics).**

linkage method has the nice property of enlarging small components before gradually connecting them.



**Figure 4. Largest components of social networks from *golden set* (left-most) and from disambiguated datasets (center and right-most), with  $\lambda = 0.5$  and  $\delta = 1.0$ .**

## 7. Conclusion and Future Work

This paper presented a system to perform record linkage over mentions of persons from sparse historical data. It deals with different constraints such as data sparsity and limited prior knowledge. We plan to apply the system to different datasets and to integrate it into a transcription and annotation interface, in order to use it for live, aided record linkage.

## 8. Acknowledgments

The second author gratefully acknowledges the support of the FNS/ANR grant No. CR12I1L\_156272.

## References

- Artiles, J., Sekine, S., and Gonzalo, J. (2008). Web people search: results of the first evaluation and the plan for the second. In *Proceedings of the 17<sup>th</sup> international conference on World Wide Web*, pages 1071–1072. ACM.
- Bellavitis, A. (2006). Apprentissages masculins, apprentissages féminins à venise au XVIe siècle. *Histoire Urbaine*, pages 49–73.
- Bloothoof, G., Christen, P., Mandemakers, K., and Schraqgen, M., editors (2015). *Population Reconstruction*. Springer.
- Cohen, W., Ravikumar, P., and Fienberg, S. (2003). A comparison of string metrics for matching names and records. In *KDD workshop on data cleaning and object consolidation*, volume 3, pages 73–78.

- Ji, H. and Grishman, R. (2011). Knowledge base population: Successful approaches and challenges. In *Proceedings of the 49<sup>th</sup> Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 1148–1158. Association for Computational Linguistics.
- Kleanthi, G., van der Burgh, B., Meeng, M., and Knobbe, A. (2015). Record linkage in medieval and early modern text. In *Population Reconstruction*, pages 173–195. Springer.
- Lehmann, J., Isele, R., Jakob, M., Jentzsch, A., Kontokostas, D., Mendes, P. N., Hellmann, S., Morsey, M., van Kleef, P., Auer, S., and Bizer, C. (2013). DPedia - A Large-scale, Multilingual Knowledge Base Extracted from Wikipedia. *Semantic Web Journal*.
- Mann, G. S. and Yarowsky, D. (2003). Unsupervised personal name disambiguation. In *Proceedings of the 7<sup>th</sup> Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4, CONLL '03*, pages 33–40.
- Shen, W., Wang, J., and Han, J. (2015). Entity linking with a knowledge base: Issues, techniques, and solutions. *Knowledge and Data Engineering*, 27(2):443–460.