

# Reconciling and Using Historical Person Registers as Linked Open Data in the AcademySampo Knowledge Graph

Petri Leskinen<sup>1</sup>[0000–0003–2327–6942] and Eero Hyvönen<sup>1,2</sup>[0000–0003–1695–5840]

<sup>1</sup> Semantic Computing Research Group (SeCo), Aalto University, Finland

<sup>2</sup> HELDIG – Helsinki Centre for Digital Humanities, University of Helsinki, Finland  
<http://seco.cs.aalto.fi>, <http://heldig.fi>, [firstname.lastname@aalto.fi](mailto:firstname.lastname@aalto.fi)

**Abstract.** This paper presents a method for extracting and reassembling a genealogical network automatically from a biographical register of historical people. The method is applied to a register that contains short textual biographical descriptions about 28 000 Finnish and Swedish university students in 1640–1899. The aim is to connect and disambiguate the relatives mentioned in the biographies in order to build a continuous, genealogical network, which can be used in Digital Humanities for data and network analysis of historical academic persons and their lives. An artificial neural network approach is presented for solving a supervised learning task to disambiguate relatives mentioned in the register descriptions using basic biographical information enhanced with an ontology of vocations and additional occasionally sparse genealogical information. Evaluation results of the record linkage are promising and provide novel insights into the problem of historical people register reconciliation. The outcome of the work has been used in practise as part of the AcademySampo semantic portal and linked open data service, a new member in the Sampo series of cultural heritage applications for Digital Humanities.

**Keywords:** Data Reconciling · Biographies · Linked Data · Digital Humanities.

## 1 Introduction

A key idea of Linked Data is to enrich datasets by integrating complementary local information sources in an interoperable way into a global knowledge graph. This involves harmonization of local data models used, as well as aligning the concepts and entities used in populating the local data models. The latter problem has been addressed traditionally in the field of *record linkage (RL)* [33,13,7], where the goal is to find matching data records between heterogeneous databases. For example, how to match person records in different registers, which may contain data about same persons, but where the data is represented using different metadata schemas and notational conventions? Using RL and data integration, richer global descriptions of persons can be created based on fusing local datasets. In addition, RL facilitates data enrichment by linking together local datasets that

use different vocabularies and identifiers for representing same resources, such as persons.

## 1.1 Research Problem Formulation

This paper concerns the problem of entity reconciliation and RL of persons in historical person registers. As a case study, university students and their relatives extracted automatically from the textual biographical descriptions of the Royal Academy of Turku and University of Helsinki are considered. The primary data contains some 28 000 short biographical descriptions of people in 1640–1899, covering virtually all university students in Finland during this time period. This data contains not only the 1) the explicit set of students recorded but also 2) the implicit set of persons mentioned in the short biography record texts of (1), such as relatives and prominent historical persons. The task is to construct a knowledge graph of all persons referred to in the data (1)–(2) in order to study the characteristics of the underlying academic network.

As solution approach, a probabilistic RL solution for linking person records is presented and tested with promising evaluation results. In our method, RL is based on the attributes of an actor, such as the name, life years, and vocations relating to her/his life. The key novel idea here is to enrich these attributes with genealogical information, i.e., information about the names and lifespans of actors’ relatives. Integrating local person registers into a single global *knowledge graph* (*KG*) facilitates biographical and prosopographical research based on enriched data. For this purpose, the aligned enriched person data has been used as a basis for a new semantic portal and data service, *AcademySampo – Finnish Academic People 1640–1899 on the Semantic Web*<sup>3</sup> [22].

## 1.2 Related Work

The RL field is presented [13,33,3]. Several nation-wide projects are underway on integrating person registries. For example, the Norwegian Historical Population Register (HPR) is pursuing to cover the country’s whole population in 1800–1964, based on combining church records and census data [29]. The Links project<sup>4</sup> in the Netherlands aims to reconstruct all nineteenth and early twentieth century families in the Netherlands based on civil certificates.

The problem of reconciling person records is evident in genealogical research. For example, in [23] Machine Learning has been applied to automatic construction of family trees from person records. Antolie et al. [2] present a case study of integrating Canadian World War I data from three sources: soldier records, casualty records, and census data. Here more traditional crafted RL processes

---

<sup>3</sup> The portal and its linked open data service, including a SPARQL endpoint, will be released on February 5, 2020. More information about AcademySampo can be found on the project homepage: <https://seco.cs.aalto.fi/projects/yo-matrikkelit/>

<sup>4</sup> Cf. the project homepage <https://iisg.amsterdam/en/hsn/projects/links> and research papers at <https://iisg.amsterdam/en/hsn/projects/links/links-publications>.

were used, and using the data in research is demonstrated. Also Cunningham [8] concerns military person data. Here World War I military service records have been integrated with a census data, and the integrated data is used for data analysis. In Ivie et al. [16] the RL process is enhanced with the available genealogical data, e.g. information about spouses and children, to achieve a higher accuracy. Also Pixton et al. [24] utilize the genealogical information and apply a neural network for RL.

Representing and analyzing biographical data has grown into a new research and application field, reported, e.g., in the Biographical Data in Digital World workshops BD2015 [4], BD2017 [11], and BD2019. In [20], analytic visualizations were created based on U.S. Legislator registry data, and the Six Degrees of Francis Bacon system<sup>5</sup> [31,19] utilizes data of the Oxford Dictionary of National Biography. Extracting Linked Data from texts has been studied in several works, such as [12]. In [10], language technology was applied for extracting entities and relations in RDF using Dutch biographies in the BiographyNet, as part of the larger NewsReader project [26].

Our own earlier works related to the topic include reconciling biographees and their relatives in the BiographySampo semantic portal [15,21]. Here genealogical statistics e.g. average ages of becoming a parent or getting married were extracted from the source data, and person’s life years are estimated according to that distribution. However, in this paper a neural network model is trained to learn similar rules from the data. References to World War II soldiers were reconciled for data linking in the in the WarSampo portal and knowledge graph [14,18].

This paper is structured as follows: We first present the primary data of our study and how it has been transformed into Linked Data. After this, the method of reconciling mentions of person in person registries is explained, and evaluation results in our case study are presented. In conclusion, contributions of the paper are discussed, and directions for further research are pointed out.

## 2 Knowledge Graph of Historical Academic Persons

This section presents the data to be used in our study: the Finnish university student registries “Ylioppilasmatrikkeli” containing short biographical descriptions about persons. A more detailed description about the data conversion is described in [22].

### 2.1 Primary Data Sources

The student registry datasets in our focus are based on original handwritten university enrollment documents. In an earlier project, the documents have been transliterated manually into textual form and extended with information from other sources about later life events of the biographees. It has been estimated that ten man years of manual work of archivists was needed to accomplish this.

---

<sup>5</sup> <http://www.sixdegreesoffrancisbacon.com>

Our work concerns two main parts of the student registry: the database covering the years 1640–1852<sup>6</sup> (*D1640*) available in Finnish and Swedish, and the registry of 1853–1899<sup>7</sup> (*D1853*) for the next years. The records contains short biographical descriptions of 28 000 students of the University of Helsinki<sup>8</sup>, originally the Royal Academy of Turku<sup>9</sup> in Finland. There are lots of mentions of relatives as well as of prominent related persons in the biographical descriptions. These student registries cover a significant part of the history of Finland and the Finnish university institution, since the University of Helsinki was the only university in the country during the time frame in focus. The data is widely used but genealogists and historians.

A key challenge in transforming this kind of data into Linked Data for data-analysis is how to reconcile mentions of people in the records and their biographical texts. For example, the data contains records of ten students with the same name *Johan Wegehus*. Furthermore, eight of them have a vocation related to clergy—more than half of the students who studied before the year 1780 worked as priests after their graduation.<sup>10</sup> In the textual descriptions of the students, there are 72 mentions of spouses or mothers with the name *Maria Johansdotter*. Furthermore, there are variations in how the names are written because the data has been collected from multiple sources by different archivists, when it was extended by additional information about the later lives of the students. For example, the name *Sofia Dorotea Cedercreutz* can also be written *Sophia Dorothea Cedercreutz*.

The data is divided into four parts: *D1640*: the students in 1640–1852; *R1640*: the relatives in *D1640*; *D1853*: the students of 1853–1899; *R1853*: the relatives of *R1853*. The aim is to link the people between these datasets. Therefore, the record linkage consists of the following partial tasks: 1) linkage from *R1640* to *D1640*, 2) linkage from *R1853* to *D1853*, 3) linkage from *R1853* to *D1640*, and 4) disambiguation of *R1640* and *R1853* data.

## 2.2 Extracting Information from Text

A comprehensive description about the data conversion as well as about the used data model is presented in an earlier article [22]. An extract of a registry entry for *Anders Israel Cajander*<sup>11</sup> is depicted in Fig. 1. The description starts with the date or year of enrollment, in this case *11.2.1830*. After that there is the full name and an unique database identifier followed by the place and time of birth ([Leppävirralla 24.2.1811](#)). Next there is a Finnish abbreviation *Vht* meaning parents; in the example case the father is *Zachris Johan Cajander* and

<sup>6</sup> <https://ylioppilasmatrikkeli.helsinki.fi>

<sup>7</sup> <https://ylioppilasmatrikkeli.helsinki.fi/1853-1899>

<sup>8</sup> [https://en.wikipedia.org/wiki/University\\_of\\_Helsinki](https://en.wikipedia.org/wiki/University_of_Helsinki)

<sup>9</sup> [https://en.wikipedia.org/wiki/Royal\\_Academy\\_of\\_Turku](https://en.wikipedia.org/wiki/Royal_Academy_of_Turku)

<sup>10</sup> This statistical result was obtained after we used the reconciled data in *AcademySampo* for data analysis.

<sup>11</sup> <https://ylioppilasmatrikkeli.helsinki.fi/henkilo.php?id=14689>

11.2.1830 **Anders Israel Cajander** 14689. \* [Leppävirralla 24.2.1811](#). Vht: Savon alisen kihlakunnan kruununvouti *Zachris Johan Cajander* (†1862) ja *Gustava Karolina Neiglick*. Kuopion triviaalikoulun oppilas 4.2.1822 – 22.6.1826 (betyg). Viipurin lukion oppilas 17.9.1827 – 1.7.1829. Ylioppilas Helsingissä 11.2.1830 (arvosana approbatur cum laude äänimäärällä 14). Viipurilaisen osakunnan jäsen 12.2.1830 *12/2 1830 \ Anders Israel Cajander \ 24/2 1811 \ KronoFogden Zachr. Joh. Cajander i Randasalmi \ Leppävirta \ [med betyg] fr. Gymn. i Wiborg \ Uttog betyg d. 12/10 1833 för att ingå vid Rättegångsverken*. Merkitty oikeustieteellisen tiedekunnan nimikirjaan 9.10.1832. Savokarjalaisen osakunnan perustajajäsen 1833 *Anders Israel Cajander*. Tuomarintutkinto 10.12.1833. Vaasan hovioikeuden auskultantti 24.12.1833. — Varatuomari 1837. Kihlakunnantuomarin arvonimi 1847. Äyräpään tuomiokunnan tuomari 1857, Jääsken tuomiokunnan 1870, Rannan tuomiokunnan 1877, ero 1891. Hovioikeudenasessorin arvonimi 1868. Laamannin arvonimi 1870. Valtiopäivämies 1872. †[Viipurissa 18.12.1901](#).

Pso: 1841 [Fredrika Emelie Schildt \(†1892\)](#).

Veli: Räisälän kappalainen [Gustaf Adolf Cajander 15376 \(yo 1835, †1882\)](#).

Veli: kirjailija [Zakarias Cajander 16147 \(yo 1843, †1895\)](#).

Lanko: lääninmetsänhoitajan apulainen [Berndt Vilhelm Kristoffer Schildt 14968 \(yo 1832, †1892\)](#).

**Fig. 1.** Partial extract from a register entry text for *Anders Israel Cajander*

the mother [Gustava Karolina Neiglick](#). After that there are two lists of events, one related to studies and academic career, and other describing the later career of the biographee. At the end of the first paragraph, a person’s death is marked with the symbol † and burial with ‡; the person in example died in Wyborg on December 18th, 1901 († [Viipurissa 18.12.1901](#)).

After the life time description, there are possible fields for relatives. In the example case, the spouse is mentioned first as [Pso: 1841 Fredrika Emelie Schildt](#) where *Pso* is a Finnish abbreviation for *puoliso* (spouse). There are three relatives who also have an entry in the register, i.e., two brothers ([Veli: Gustaf Adolf Cajander](#) and [Veli: Zakarias Cajander](#)) and a brother-in-law ([Lanko: Berndt Vilhelm Kristoffer Schildt](#)). The author of the *D1640* dataset Yrjö Kotivuori has manually added links from the description texts to the mentioned people also found in the register, like the three relatives in the example case. These links also contain linkage to the relatives in the *D1853* dataset.

### 2.3 Available Information

The previous person entry example was from the *D1640* data. However, the provided data in *D1853* differs in some aspects. For instance, *D1853* only mentions a person’s parents and spouses, never children or any other relatives, and the people are not interlinked. Abbreviations are used generally for, e.g., vocations, which was taken into consideration in the data conversion by using specific lists of abbreviations.

Table 1 shows an analysis of the known positive sample pairs in the both datasets. Here column *source* refers to the relative, and *target* to the corresponding student entry. The rows show how many of the example pairs of particular data field are available. One can notice that for the six uppermost properties, e.g., preferred label, gender, death, vocation, child, and spouse are available for both the source and target records. On the other hand, the data fields indicating the

place of death, year of birth, names of mother or father, as well as the alternative labels are usually not available. The column *common* indicates the number of cases where both the source and the target entries have the particular data field and *same* the number of entries where the source and the target values are equal.

This table clearly indicates which properties should be considered crucial in decision making. Notice that some attributes that are usually significant for a general case of RL, such as places of birth and death, are not chosen in this particular case study.

	Data 1640–1852				Data 1853–1899			
	source	target	common	same	source	target	common	same
preferable label	4285	4285	4285	3979	698	698	698	517
gender	4283	4284	4283	4283	698	698	698	688
year of death	4229	4208	4192	4141	135	352	134	130
vocation	4281	4270	4270	940	600	567	543	365
child	4285	4284	4284	3211	430	341	340	2
spouse	4285	4273	4273	-	698	687	687	2
place of death	2	3494	2	2	-	348	-	-
year of birth	-	2906	-	-	-	351	-	-
mother	-	3475	-	-	-	349	-	-
father	-	3478	-	-	-	348	-	-
alternative label	-	1761	-	-	30	165	29	22

**Table 1.** Available data fields in the training data

### 3 Method: Linking Person Records

This section describes the chosen formats for comparing two person registry entries. Generally, the input format for data comparison consists of numeric difference or similarity values between the data points of the two records, not the data of the records as it is. We first introduce the chosen input formats for data in different domains, e.g., for names and for vocations of the actors and the relatives. Finally, the architecture of the network model as well as the training setting are introduced.

#### 3.1 Person Names

Person names in the datasets consists of a preferred and possibly alternatives labels. Each label includes a family name and a sequence of given names. For the classifier input we considered four different variations of a label with a maximum

of three first given names, only 0.4% of people entries have more than three given names. The classifier input is in a matrix format where the entry elements are statistical values calculated from the dataset. Each family and given name gets a *rarity* value so that first the frequency of the appearances for each name is counted and the ranks are mapped into the numeric range [0.0, 1.0]: the most common names get a near-zero and the rarest values closer to 1.0 in order to distinguish the rare names.

Fig. 2 depicts an example of a name comparison matrix, in this case the family names of two person entries. The rows and columns mutually correspond to the data of two names that are compared. The uppermost row (0.000, 0.808, 0.983, ...) consists of the rarity values for the first, and likewise the leftmost column (0.000, 0.987, 0.991, 0.100, ...) for the second entry. The other values inside the matrix are Jaro-Winkler similarity values [32] between the name strings so that e.g., perfectly matching names get the value 1.0.

Hendriksson	Hendricius	Hendricius	Hendricius	Hendricius	
(rarity)	0.000	0.808	0.983	0.934	0.817
Hendriksson	0.987	0.733	0.717	0.967	0.859
Hendricius	0.991	0.600	0.842	0.813	0.933
Hendricius	0.100	0.970	0.735	0.737	0.660
Hendricius	0.100	0.000	0.000	0.000	0.000
	(rarity)	Hendriksson	Hendricius	Hendricius	Hendricius

**Fig. 2.** Example of a matrix for comparing family names

### 3.2 Vocations

The vocations are the titles extracted from the source data. These titles often consists of a place name and a related profession, e.g., *Bishop of Turku* or *Bishop of Porvoo*. To enrich the data the vocations are linked to the hierarchical AMMO [17] ontology of historical occupations. Statistical values are used here like with the name entries. A *rarity* value is calculated for each title following the same principle as with the titles. In addition to that, a value of co-occurrences between two titles is calculated.

Fig. 3 depicts an example of a vocation comparison matrix. The value in the leftmost upper corner (0.455) is the Jaccard index [28] between the two sets of vocations. Similarly to the name matrix, the rows and columns correspond to the vocation in two dataset entries with the rarity values on the uppermost row (0.909, 0.804, 0.249...) and leftmost columns. The rarity values are in a descending order so that the rarest vocations appear first on the lists. The other values filling the rest of the matrix are the co-occurrence values. In the data matrix, the co-occurrence value for a pair (*Law Reader*, *Mayor*) is 0.985, while the pair (*Court Attorney*, *Mayor*) has a value 0.250 indicating that this pair co-occurs in the data more frequently. The zero-valued elements on the right indicate that one of the title sets has less than the reserved seven data fields.

(rarity)	0.455	0.909	0.804	0.249	0.019	0.014	0.000	0.000
Law Reader	0.804	0.000	0.898	0.985	0.938	0.935	0.000	0.000
County Secretary	0.536	0.000	0.000	0.000	0.818	0.818	0.000	0.000
Mayor	0.249	0.966	0.985	0.250	0.250	0.250	0.000	0.000
Statesman	0.100	0.000	0.000	0.806	0.410	0.380	0.000	0.000
Public administration	0.081	0.000	0.000	0.960	0.398	0.358	0.000	0.000
Socio-administrative Work	0.029	0.000	0.000	0.938	0.290	0.269	0.000	0.000
Court Attorney	0.019	0.966	0.938	0.250	0.012	0.012	0.000	0.000
	(rarity)	Mayor of Oulu	Law Reader	Mayor	Court Attorney	Legal work		

**Fig. 3.** Matrix for comparing the vocations

### 3.3 Years of birth and death, gender

The difference in actor's and relatives' birth and death years and their genders were also input to the network. The years use a precision of one year due to the format used in source data: the birth and death of the actor is usually known with a precision of a day, while in the case of relatives only the precision of a year is used. The actual difference in years is mapped into a near-zero range by using the arctan function. Gender was indexed using value  $-1.0$  for female,  $1.0$  for male, and  $0.0$  for the rare cases where the gender was not known.



### 3.4 Relative information

The information of the relatives consists of details about the children and spouses of an actor, and basic information about her/his parents. The relative information uses the same matrix format as for the names, lifetime information, and vocations of the actor. It has reserved space for three children and three spouses, according to analyzing the data. In the data more than 99% have three or less spouses, and 95% three or less children.

### 3.5 Network model

The used network model is depicted in Fig. 4. It is a multi-input network based on the Keras functional API [6]. The network has eight inputs out of which six for the given and family names of the two actors, their spouses, and their children, one for the age comparison of actors and their relatives, and one for the actors' titles. The network as a probabilistic classifier and the output is  $\bar{y} \approx [0.0, 1.0]$  for matching entries and  $\bar{y} \approx [1.0, 0.0]$  for not matching pairs. For a binary decision these values are filtered by choosing the positive matches when the latter value exceeds a chosen threshold, e.g.,  $\lambda = 0.9$ .

Some inputs are in a matrix format, which are first flattened<sup>12</sup>, and after that run through a Dense<sup>13</sup> layer. Dropout layers with a ratio of 25% are used to prevent the overfit to the training data [27]. Different inputs of the same domain (e.g., names and years) are first concatenated<sup>14</sup> to one another. After a layer of Dense network the network concatenates into the final output.

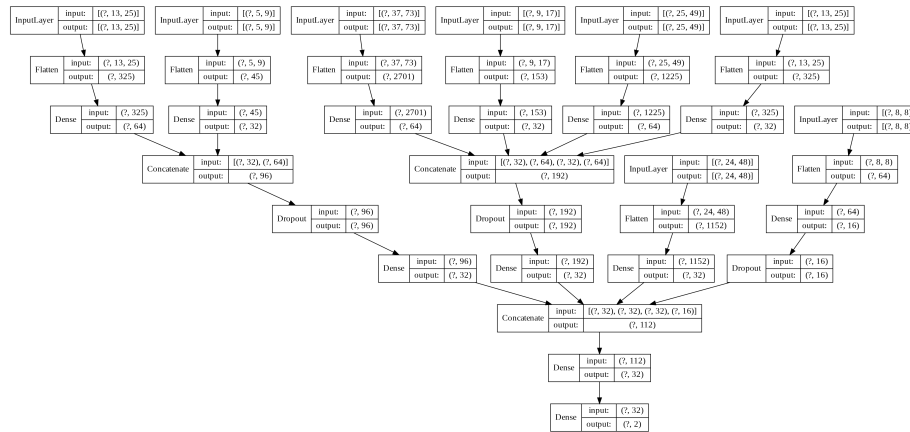


Fig. 4. Classifier model structure

<sup>12</sup> [https://keras.io/api/layers/reshaping\\_layers/flatten/](https://keras.io/api/layers/reshaping_layers/flatten/)

<sup>13</sup> [https://keras.io/api/layers/core\\_layers/dense/](https://keras.io/api/layers/core_layers/dense/)

<sup>14</sup> [https://keras.io/api/layers/merging\\_layers/concatenate/](https://keras.io/api/layers/merging_layers/concatenate/)

**Training Data.** The training data for the neural network could be input by either as an single data entry or in several smaller batches of data. We chose to feed the data in batches utilizing the Keras Data Generator Sequence [1] as described by A. Amidi and S. Amidi<sup>15</sup> due to the amount of data preprocessing from RDF format to numeric input.

Positive samples are created by reading the manually marked matches from the data. This linkage is many-to-one, so all the samples pointing to the same target can be chosen as training pairs pointing to each other. Finally, positive sample data is augmented with pairs where both the target and the source refer to the same resource.

The easiest way to gather the negative samples is to pick random pairs from the data. However, we chose to sample pairs that are likely to have some similar data values to improve the decision making. The dataset contains relations indicating, e.g., that two persons are siblings, cousins, or namesakes. Close relatives often have same similar characteristics, like family name or nearby years of birth.

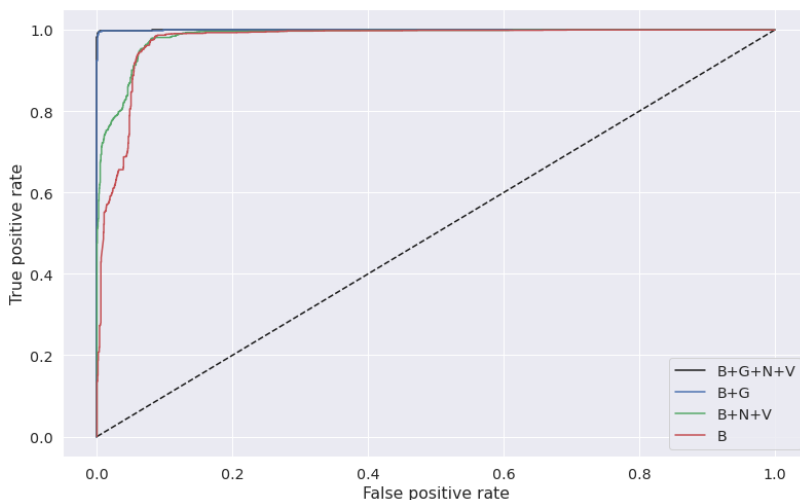
**Model Training.** For the training the data was split into separate sets for training, testing, and validation of sizes 70%, 15%, and 15%, respectively. The classes in the training data are imbalanced, e.g., the number of negative samples ( $N_n \approx 200000$ ) is significantly larger than the positive samples ( $N_p \approx 13000$ ). Therefore the positive samples were defined to have a larger weight than the negative ones [30,5]. The training was performed in Google Colab, and the training with 100 epochs using a GPU took 4242.2 seconds. Validation accuracy of more than 99.6% was achieved during the training.

## 4 Evaluation

The results were analyzed closely by the Receiver Operating Characteristic (ROC) curve (Fig. 5) and by taking look at the details of False Positive and False Negative classifications. To deal with the data imbalance, a validation set with equal amount of positive and negative sample was used. The classifier input was divided by four different types: basic biographical information (B), genealogical information (G), name frequencies (N), and vocation frequencies (V). To analyze how much each data entry contributes to the prediction, evaluation was performed for four times using the entire data (B+G+N+V), biographical and genealogical data (B+G), biographical data with name and vocation frequencies (B+N+V), and the plain biographical data (B). The threshold value  $\lambda$  for optimal performance was chosen from the ROC curve coordinates by the point closest to the upper left corner [9]. For the entire data (B+G+N+V) the threshold value was  $\lambda = 90.01\%$  and the resulting number of True Positives (TP) is 2035, True Negatives (TN) 2089, False Positives (FP) 0, and False Negatives (FN) 54 with measures precision of 100.00%, recall of 97.42%, F<sub>1</sub>-score of 98.69%, and accuracy of 98.71%.

---

<sup>15</sup> <https://stanford.edu/~shervine/blog/keras-how-to-generate-data-on-the-fly>



**Fig. 5.** ROC curve

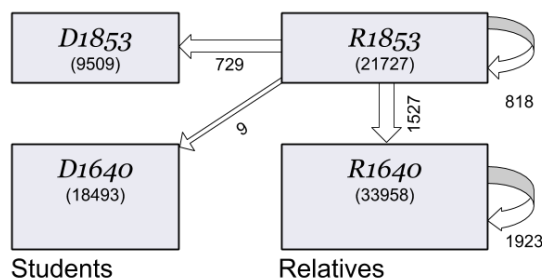
In the ROC visualization, the curve with basic and genealogical (B+G) almost emerges with the curve for the entire data (B+G+N+V). Also Table 2 shows how close these results are to one another. Furthermore, the validation results without the genealogical information (B, B+N+V) show lower accuracy.

Data Subset	TP	FP	FN	TN	Precision	Recall	F <sub>1</sub> -score	Accuracy	AUC	$\lambda$
B+G+N+V	2035	0	54	2089	100.00%	97.42%	98.69%	98.71%	99.98%	90.01%
B+G	2007	1	82	2088	99.95%	96.07%	97.97%	98.01%	99.97%	84.06%
B+N+V	2011	150	78	1939	93.06%	96.27%	94.64%	94.54%	98.47%	16.86%
B	587	12	1502	2077	98.00%	28.10%	43.68%	63.76%	97.48%	92.15%

**Table 2.** Validation results using different data subsets

**Full disambiguation** Record linkage with the real dataset was a many-to-one task, e.g. many records in the source set can be merged into one in the target data. When applying the model to the real dataset first blocking strategies [7] were applied to reduce the number of comparisons. For instance, candidate pairs of different gender or mismatching life years when known, could be omitted from candidate pairs. Likewise, candidates mentioned in a same register entry text e.g. siblings or different spouses could be omitted—same person is never mentioned twice in one text entry. Some preliminary disambiguation was performed already during the data conversion, e.g., aligning spouses of a person, if the names had a

high string similarity. The iterative process was run for several times because merging two person records furthermore can lead to finding more matches also among the relatives. To achieve a high precision and to minimize the number of false positive classification a high threshold values ( $\lambda \geq 0.9$ ) were used.



**Fig. 6.** Number of matches between the datasets

Fig. 6 depicts the number of records in each part of the dataset and the numbers of matches detected within them. The number in parenthesis inside the block is the number of records before RL. For example, 729 of the records in *R1853* were merged into *D1853*, 1527 into *R1640*, and 9 into *D1640*. The latter number is relatively small because this matching was a part of the existing manual linkage by the dataset author, so these results are links missing from manual linkage or errors in our data conversion process. Inside the *R1853* dataset, 818 and in *R1640* 1923 entries were matched, respectively.

## 5 Discussion

### 5.1 Contributions

Our work describe in this article shows that using genealogical information in RL is useful and can improve significantly the accuracy in person name reconciliation. This argument was tested and evaluated in detail in a case study using the AcademySampo datasets with promising results. We anticipate that similar results can be obtained in related use cases using other similar dataset. In the AcademySampo project, the genealogical information has been used also when linking the records with Wikidata for semantic data enrichment.

When analysing the resulting matched pairs some weak cases needing separate handling where found. Historically, patronymic family names, e.g., *Johansdotter* (*Daughter of Johan*) have been common for women. However, the chosen Jaro-Winkler similarity may not be optimal to always disambiguate between cases like *Jöransdotter* and *Johansdotter*. Likewise, the classifier made some false results with the vocation of a farmer. Farmer was a common vocation in the 17th–19th

century Finland, but yet rare in data records of academic people, for which reason we had put some excess weight on it in the classifying system.

This paper presented a method for reconciling person names mentioned in biographical texts of other persons. The method was applied to creating a semantic knowledge graph of persons to be used for studying and analysing academic networks of people. For this purpose, the AcademySampo portal has been created, but also the underlying open linked data service can be used for custom-made data-analyses using, e.g., YASGUI<sup>16</sup> [25] and SPARQL or Python scripting in Google Colab<sup>17</sup> or Jupyter<sup>18</sup> notebooks, and for developing new applications.

**Acknowledgements** Thanks to Yrjö Kotivuori and Veli-Matti Autio for their seminal work in creating the original Ylioppilasmatrikkeli databases used in our work, and for making the data openly available for our project, and to be published as a Linked Open Data service later on. We also wish to thank Mikko Koho and Jouni Tuominen for fruitful discussions. Our work has been funded by Aalto University and Helsinki Centre for Digital Humanities (HELDIG), and it is also part of the EU project InTaVia<sup>19</sup> (2020–2023) on integrating and studying intangible biographical data with tangible Cultural Heritage.

## References

1. Keras Documentation, Sequence. [https://www.tensorflow.org/api\\_docs/python/tf/keras/utils/Sequence](https://www.tensorflow.org/api_docs/python/tf/keras/utils/Sequence), accessed: 2020-12-10.
2. Antonie, L., Gadgil, H., Grewal, G., Inwood, K.: Historical Data Integration, a Study of WWI Canadian Soldiers. In: 2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW). pp. 186–193. IEEE (2016).
3. Barlaug, N., Gulla, J.A.: Neural networks for entity matching. arXiv preprint arXiv:2010.11075 (2020).
4. ter Braake, S., Anstke Fokkens, R.S., Declerck, T., Wandl-Vogt, E. (eds.): BD2015, Biographical Data in a Digital World 2015. CEUR Workshop Proceedings, Vol-1399 (2015), <http://ceur-ws.org/Vol-1272/>.
5. Brownlee, J.: Machine Learning Mastery: How to Develop a Cost-Sensitive Neural Network for Imbalanced Classification. <https://machinelearningmastery.com/cost-sensitive-neural-network-for-imbalanced-classification/>, accessed: 2020-12-10.
6. Chollet, F.: Keras, The Functional API. [https://keras.io/guides/functional\\_api/](https://keras.io/guides/functional_api/), accessed: 2020-12-10.
7. Christen, P.: Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection. Springer Science & Business Media (2012).
8. Cunningham, A.: After “it’s over over there”: Using record linkage to enable the reconstruction of World War I veterans’ demography from soldiers’ experiences to civilian populations. *Historical Methods: A Journal of Quantitative and Interdisciplinary History* 51, 1–27 (2018).

---

<sup>16</sup> <https://yasgui.triply.cc>

<sup>17</sup> <https://colab.research.google.com/notebooks/intro.ipynb>

<sup>18</sup> <https://jupyter.org>

<sup>19</sup> <https://seco.cs.aalto.fi/projects/intavia/>

9. Fawcett, T.: An introduction to ROC analysis. *Pattern recognition letters* 27(8), 861–874 (2006).
10. Fokkens, A., ter Braake, S., Ockeloen, N., Vossen, P., Legêne, S., Schreiber, G., de Boer, V.: *BiographyNet: Extracting Relations Between People and Events*. In: *Europa baut auf Biographien*. pp. 193–224. New Academic Press, Wien (2017).
11. Fokkens, A., ter Braake, S., Sluijter, R., Arthur, P., Wandl-Vogt, E. (eds.): *BD2017 Biographical Data in a Digital World 2015*. CEUR Workshop Proceedings, Vol-1399 (2017), <http://ceur-ws.org/Vol-2119/>.
12. Gangemi, A., Presutti, V., Recupero, D.R., Nuzzolese, A.G., Draicchio, F., Mongiovì, M.: *Semantic web machine reading with FRED*. *Semantic Web* 8, 873–893 (2017).
13. Gu, L., Baxter, R., Vickers, D., Rainsford, C.: *Record linkage: Current practice and future directions*. CSIRO Mathematical and Information Sciences (2003), cMIS Technical Report No. 03/83.
14. Heino, E., Tamper, M., Mäkelä, E., Leskinen, P., Ikkala, E., Tuominen, J., Koho, M., Hyvönen, E.: *Named entity linking in a complex domain: Case second world war history*. In: *Proceedings, Language, Technology and Knowledge (LDK 2017)*. pp. 120–133. Springer-Verlag (June 2017), [https://link.springer.com/chapter/10.1007/978-3-319-59888-8\\_10](https://link.springer.com/chapter/10.1007/978-3-319-59888-8_10).
15. Hyvönen, E., Leskinen, P., Tamper, M., Rantala, H., Ikkala, E., Tuominen, J., Keravuori, K.: *BiographySampo – publishing and enriching biographies on the semantic web for digital humanities research*. In: *Proceedings of the 16th Extended Semantic Web Conference*. Springer-Verlag (2019).
16. Ivie, S., Pixton, B., Giraud-Carrier, C.: *Metric-based data mining model for genealogical record linkage*. In: *2007 IEEE International Conference on Information Reuse and Integration*. pp. 538–543. IEEE (2007).
17. Koho, M., Gasbarra, L., Tuominen, J., Rantala, H., Jokipii, I., Hyvönen, E.: *AMMO Ontology of Finnish Historical Occupations*. In: *Proceedings of the The First International Workshop on Open Data and Ontologies for Cultural Heritage (ODOCH'19)*. vol. 2375, pp. 91–96. CEUR Workshop Proceedings (June 2019), <http://ceur-ws.org/Vol-2375/>, vol 2375.
18. Koho, M., Leskinen, P., Hyvönen, E.: *Integrating historical person registers as linked open data in the warsampo knowledge graph*. In: *Blomqvist, E., Groth, P., de Boer, V., Pellegrini, T., Alam, M., Käfer, T., Kieseberg, P., Kirrane, S., Meroño-Peñuela, A., Pandit, H.J. (eds.) Semantic Systems. In the Era of Knowledge Graphs. SEMANTiCS 2020. Lecture Notes in Computer Science*, vol. 12378, pp. 118–126. Springer, Cham (October 2020), [https://doi.org/10.1007/978-3-030-59833-4\\_8](https://doi.org/10.1007/978-3-030-59833-4_8).
19. Langmead, A., Otis, J., Warren, C., Weingart, S., Zilinski, L.: *Towards interoperable network ontologies for the digital humanities*. *Int. J. of Humanities and Arts Computing* 10(1), 22–35 (2016).
20. Larson, R.: *Bringing lives to light: Biography in context (2010)*, Final Project Report, University of Berkeley, [http://metadata.berkeley.edu/Biography\\_Final\\_Report.pdf](http://metadata.berkeley.edu/Biography_Final_Report.pdf).
21. Leskinen, P., Hyvönen, E.: *Extracting genealogical networks of linked data from biographical texts*. In: *Proceedings of the 16th Extended Semantic Web Conference (ESWC 2019)*. Springer-Verlag (2019).
22. Leskinen, P., Hyvönen, E.: *Linked open data service about historical finnish academic people in 1640–1899*. In: *DHN 2020 Digital Humanities in the Nordic Countries. Proceedings of the Digital Humanities in the Nordic Countries 5th Conference*. pp. 284–292. CEUR Workshop Proceedings, vol. 2612 (October 2020), <http://ceur-ws.org/Vol-2612/short14.pdf>.

23. Malmi, E., Gionis, A., Solin, A.: Computationally inferred genealogical networks uncover long-term trends in assortative mating. arXiv (2018), arXiv:1802.06055 [cs.SI].
24. Pixton, B., Giraud-Carrier, C.: Using Structured Neural Networks for Record Linkage. In: Proceedings of the Sixth Annual Workshop on Technology for Family History and Genealogical Research (2006).
25. Rietveld, L., Hoekstra, R.: The YASGUI family of SPARQL clients. *Semantic Web – Interoperability, Usability, Applicability* 8(3), 373–383 (2017). <https://doi.org/10.3233/SW-150197>.
26. Rospocher, M., van Erp, M., Vossen, P., Fokkens, A., Aldabe, I., Rigau, G., Soroa, A., Ploeger, T., Bogaard, T.: Building event-centric knowledge graphs from news. *Web Semantics: Science, Services and Agents on the World Wide Web* 37, 132–151 (2016).
27. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A Simple Way to Prevent Neural Networks From Overfitting. *The Journal of Machine Learning Research* 15(1), 1929–1958 (2014).
28. Tan, P.N., Steinbach, M., Kumar, V.: Introduction to data mining. 1st (2005).
29. Thorvaldsen, G., Andersen, T., Sommerseth, H.L.: Record linkage in the historical population register for Norway. In: Population reconstruction, pp. 155–171. Springer (2015).
30. Wang, S., Liu, W., Wu, J., Cao, L., Meng, Q., Kennedy, P.J.: Training deep neural networks on imbalanced data sets. In: 2016 international joint conference on neural networks (IJCNN). pp. 4368–4374. IEEE (2016).
31. Warren, C., Shore, D., Otis, J., Wang, L., Finegold, M., Shalizi, C.: Six degrees of Francis Bacon: A statistical method for reconstructing large historical social networks. *Digital Humanities Quarterly* 10(3) (2016).
32. Winkler, W.E.: String comparator metrics and enhanced decision rules in the fellegi-sunter model of record linkage (1990).
33. Winkler, W.E.: Overview of Record Linkage and Current Research Directions. Tech. rep., U.S. Census Bureau (2006).