

Automatic Annotation Service APPI: Named Entity Linking in Legal Domain

Minna Tamper^{1,2}[0000–0003–1695–5840], Arttu Oksanen^{1,3}[0000–0003–2327–6942],
Jouni Tuominen^{1,2}[0000–0003–4789–5676], Aki Hietanen⁴, and
Eero Hyvönen^{1,2}[0000–0003–1695–5840]

¹ Semantic Computing Research Group (SeCo), Aalto University, Finland

<http://seco.cs.aalto.fi>, firstname.lastname@aalto.fi

² HELDIG – Helsinki Centre for Digital Humanities, University of Helsinki, Finland

<http://heldig.fi>

³ Edita Publishing Ltd.

<http://www.editapublishing.fi>

⁴ Ministry of Justice, Finland

<http://oikeusministerio.fi>, firstname.lastname@om.fi

Abstract. Texts referencing court decisions, statutes, and EU directives can be difficult to understand without context. It can be time consuming and expensive to find related statutes or to learn about context specific terminology. As a solution, we utilized named entity linking tool for extracting information and tailored it into a service, APPI, that can automatically annotate legal documents to provide context to the readers. The service can identify and link named entities and references to legal texts to corresponding vocabularies and data sources by combining statistics- and rule-based named entity recognition with named entity linking. The results provide users with enhanced reading experience with contextual information and possibility to access related materials, such as statutes and court decisions.

Keywords: Automatic annotation service · legal texts · named entity linking · linked data

1 Introduction

The research hypothesis of this paper is that by annotating and linking legal texts to knowledge bases it is possible to assist readers to understand the text and context by offering information about legislation, context, and terminology. To understand and interpret legal texts correctly, it is often important to get acquainted with other related contextual material. The linking of texts through similarity or references can aid in finding information. To support end users in close reading and to enable linking of legal texts, we created a service called APPI⁵. It utilized a named entity linking tool, NELLI [10], for identifying domain specific information and to enable named entity linking of legal texts. As a result,

⁵ A demonstrator that is under development is available at <http://nlp.ldf.fi/appi/>.

the APPI service can identify and link named entities and references to legal texts to corresponding vocabularies and data sources by combining statistics- and rule-based named entity recognition (NER) with named entity linking (NEL). The end results of APPI can be edited in the application and they can be downloaded in JSON format.

2 Data

Semantic Finlex⁶ [7] is a web service that hosts the Finnish legislation and case law as Linked Open Data. Currently, the data published in Semantic Finlex includes consolidated statutes with version history (approx. 2500 statutes), the original statutes as published in the official journal (approx. 50000 statutes), the Judgments of the Supreme court (5500), and the Judgments of the Supreme administrative court (7500). In addition, the data contains keywords used by the Supreme Court and the Supreme Administrative Court to annotate the court judgments. The judgments are also linked to judges and personnel contributing to the case. The original statutes are also linked to the EU law and Finnish government bills. The service includes the legal texts in text, HTML, and XML formats. The documents are written in Finnish and Swedish.

3 Method

In order to automatically annotate the legal texts of Semantic Finlex, the NELLI tool [10] was utilized. NELLI is a combination of NER and NEL tools and disambiguates entities using a scoring scheme where the most popular named entity type, the longest string, and linked interpretation wins. Initially, NELLI was a command line tool that could be only used for annotating text documents. In order to annotate and provide context to legal texts, the tool was transformed into a restful API service. The number of input and output parameters was extended to support HTML, XML, and text formats, and the output format was changed to JSON that returns the annotated document in the original form and a list of entities. Also, new tools were added in order to recognize more named entity types; FinBERT⁷ [11], Regular expression-based named entity linker, Reksi⁸ (Finnish for rector), and Person Name Finder⁹. Reksi is a NEL tool that uses numerous regular expressions to identify named entities, such as registry numbers, references to statutes, and case law from the text and links them to corresponding knowledge bases. It utilizes the regularity of the forms of the entities in texts and formats them to find the matching entities from the target ontologies. The Person Name Finder service is a tool for identifying references to people by linking the names to a large Finnish person name ontology HENKO¹⁰. In addition,

⁶ <http://data.finlex.fi>

⁷ <http://turkunlp.org/FinBERT/>

⁸ <http://nlp.ldf.fi/reksi>

⁹ <http://nlp.ldf.fi/name-finder>

¹⁰ <http://light.onki.fi/henko/en/>

an existing tool, LINFER tool[10], was upgraded to identify more organizations from the texts. The service is currently only for the Finnish language documents but it is possible to configure NELLI for other languages.

4 Application

The APPI web application was built on top of the results of the NELLI service to visualize them and to provide context and recommendations to the legal texts by linking the given text to different ontologies and to other legal texts in the Semantic Finlex dataset. For this purpose, the application form for annotating consists of an input field, input format (e.g., text, XML) selection, toggles for selecting what tools to use in NELLI, and linking options. The linking options consist of ontologies and vocabularies located in a drop-down menu that have been configured in the ARPA tool [3]. The tool can form n-grams from the given text and linguistically manipulate it (e.g., lemmatize) to match it to the given ontology. Currently, the linking options have been set to link mentions in the text to common Finnish place names (YSO places¹¹), legal terminology (the consolidated vocabulary of Finnish legal terms (draft) [2], the Helsinki Term Bank for Arts and Sciences, DBpedia), and terms used by EU institutions (EuroVoc¹²), in addition to Semantic Finlex keywords, statutes, and case law. After configuring the application, the user can click the “Annotate” button, and it annotates the given input and retrieves recommendations (i.e., similar court decisions) using the Semantic Finlex case law finder [9]. The results are presented in Fig. 1.

Results

Legend: person, animal, mythical or fictional person, general location, address, political location (e.g. state), geographical location, buildings or structures, astronomical locations (e.g. planets, galaxies), organization, media organization, financial organizations, corporation and administration, date, time, product, event, units (e.g. grams, meters), money, registry numbers, social security numbers, statutes, case law, domain information, title, and vocation, and unknown entity ?

Työntekijän¹ palkkaatavia koskeva kanne² oli jätetty tutkimatta, koska siitä ei ollut nostettu työsopimuslain 13 luvun 9 §³:n 3 momentissa säädettyä kahden vuoden määräajassa⁴ työsuhteeseen päätymisestä vaan vasta 3.59 vuoden kuluttua siitä. Saman säännöksen mukaan palkkaatava vanhentuu kuitenkin pykälän 1 momentissa säädetyn tavoin ja siis viiden vuoden kuluttua eräntymispäivästä⁵, jos työntekijän⁶ saatavan perusteena olevia työehtosopimuksen⁷ määräyksiä⁸ on pidettävä ilmeisen tulkinnanvaraisina. Kanteessa⁹ oli kysymys¹⁰ työsuhteeseen¹¹ sovellettavasta työehtosopimuksesta¹² ja palkkaatavan perusteena olevana ilmeisen tulkinnanvaraisena määräyksistä¹³. Myös sellaista voitiin pitää palkkaatavan perusteena olevana ilmeisen tulkinnanvaraisena määräyksenä¹⁴. Kannetta¹⁵ ei olisi saanut jättää 3 momentin nojalla vanhentuneena tutkimatta. TSL 13 luku 9 §¹¹ 3 mom

Time
0. 3.5
Statutes, Court decisions
1. työsuopimuslaki 13 luku 9 §
11. TSL 13 luku 9 §

Related documents:
ECLI:FI:KKO:2011:85 ECLI:FI:KKO:2016:19 ECLI:FI:KKO:2001:T894 ECLI:FI:KKO:2001:29 ECLI:FI:KKO:1987:T1261 ECLI:FI:KKO:1980:I177

JSON Download JSON

Fig. 1. Results of annotating an abstract of a court decision.

The results are presented under the configuration interface accompanied by a legend that shows available named entity types and how they are shown in text.

¹¹ <https://finto.fi/ysopaikat/en/>

¹² <http://eurovoc.europa.eu>

Below the legend is the annotated text and on its right side a list of entities found in the text (by type). The linked entities are shown with links and by clicking them a popup appears and shows the description of the given entity. Occasionally, when there is more than one option for an entity, all of them are shown in the popup and the user can select the correct one. In case the application has not found a matching entity, the user can use an autocompletion search field in the popup to query for suitable entities and link the entity manually. Below the text, there is also a list of similar documents that have been retrieved for the input text. At the bottom of the page, the JSON response shown that can be viewed or downloaded by clicking the tab.

In this example (Fig. 1), APPI has identified a reference to time, statutes, and references to different contextual terms from an abstract of a court decision. The statutes and times are identified but not linked whereas the domain information entities have been linked. The linking options were set to link legal terminology (i.e., domain information) to the consolidated vocabulary of Finnish legal terms and to the Helsinki Term Bank for Arts and Sciences. The Reksi tool links statutes and case law to Semantic Finlex. However, currently the endpoint doesn't contain all the alternative names for the statutes and therefore the linking fails. Below the text, the application has retrieved six related court decisions. The user can click the links to read the related documents in Semantic Finlex.

5 Related Work and Discussion

The APPI service provides easy access to related legal texts and helps to understand the terminology. The inspiration for the application has been the contextual reader application CORE [4] that was created to link text into ontologies in real-time to provide related materials and context. This application was initially utilized in the Semantic Finlex portal [7], configured to use content-related ontologies to provide context for the user. However, the tool does not have a powerful disambiguation system like other named entity linking tools, e.g., DBpedia Spotlight¹³ [6] and Gate Cloud¹⁴ [5]. For this purpose NELLI was created and using it a contextual reader was created for the BiographySampo portal [10]. With NELLI, the entities are not extracted in real time but in a preprocessing phase that ensures robust semantic disambiguation similarly to [8,1]. The results are recorded in RDF format and visualized by the contextual reader by quering them from the BiographySampo endpoint.

The initial demo application, APPI, manages to identify, highlight, and link named entities from the text. The annotation accuracy using NELLI was approximately 80% [10] for people and places in biographical texts. The service has been upgraded and the results are promising but it still needs a formal evaluation, which will be carried out in the future. The recommendations and legal text references can be identified with varying accuracy partially due to lack of document metadata. The current version is still under development and more work

¹³ <https://www.dbpedia-spotlight.org/demo/>

¹⁴ <https://cloud.gate.ac.uk>

needs to be done so that it can be utilized to extract all references to legislative texts such as EU directives and link them to the CELLAR system¹⁵. The APPI demo presents how by annotating documents it is possible to cater information and related documents to provide context to the reader automatically.

Acknowledgments This work is part of the ANOPPI project¹⁶ funded by the the Ministry of Justice in Finland. CSC – IT Center for Science, Finland, provided us with computational resources.

References

1. Ferragina, P., Scaiella, U.: Tagme: on-the-fly annotation of short text fragments (by Wikipedia entities). In: Proceedings of the 19th ACM international conference on Information and knowledge management. pp. 1625–1628. ACM (2010)
2. Frosterus, M., Tuominen, J., Hyvönen, E.: Facilitating re-use of legal data in applications—Finnish law as a linked open data service. In: Proceedings of the 27th International Conference on Legal Knowledge and Information Systems (JURIX 2014). pp. 115–124. IOS Press (December 2014)
3. Mäkelä, E.: Combining a REST lexical analysis web service with SPARQL for mashup semantic annotation from text. In: Proceedings of the ESWC 2014 demonstration track. pp. 424–428. Springer-Verlag (2014)
4. Mäkelä, E., Lindquist, T., Hyvönen, E.: CORE – a contextual reader based on linked data. In: Proceedings of Digital Humanities 2016, Krakow, Poland (long papers). pp. 267–269 (2016)
5. Maynard, D., Roberts, I., Greenwood, M.A., Rout, D., Bontcheva, K.: A framework for real-time semantic social media analysis. *Journal of Web Semantics* **44**, 75–88 (2017)
6. Mendes, P.N., Jakob, M., García-Silva, A., Bizer, C.: DBpedia Spotlight: Shedding light on the web of documents. In: Proceedings of the 7th international conference on semantic systems. pp. 1–8. ACM (2011)
7. Oksanen, A., Tuominen, J., Mäkelä, E., Tamper, M., Hietanen, A., Hyvönen, E.: Semantic Finlex: Transforming, publishing, and using Finnish legislation and case law as linked open data on the web. In: Peruginelli, G., Faro, S. (eds.) *Knowledge of the Law in the Big Data Age*, *Frontiers in Artificial Intelligence and Applications*, vol. 317, pp. 212–228. IOS Press (2019)
8. Piccinno, F., Ferragina, P.: From TagME to WAT: A New Entity Annotator. In: Proceedings of the first international workshop on Entity recognition & disambiguation. pp. 55–62. ACM (2014)
9. Sarsa, S., Hyvönen, E.: Searching case law judgements by using other judgements as a query (2019), submitted article under evaluation
10. Tamper, M., Hyvönen, E., Leskinen, P.: Visualizing and analyzing networks of named entities in biographical dictionaries for digital humanities research. In: Proceedings of the 20th International Conference on Computational Linguistics and Intelligent Text Processing (CICling 2019). Springer-Verlag (2019), forthcoming
11. Virtanen, A., Kanerva, J., Ilo, R., Luoma, J., Luotolahti, J., Salakoski, T., Ginter, F., Pyysalo, S.: Multilingual is not enough: BERT for Finnish. *CoRR abs/1912.07076* (2019)

¹⁵ <https://data.europa.eu/euodp/data/dataset/sparql-cellar-of-the-publications-office>

¹⁶ <https://seco.cs.aalto.fi/projects/anoppi/en/>

JT: Viitteessä 11 (FinBERT) pitäisi näkyä julkaisu-alustana arXiv; mikä on CoRR?