# Clustering-based Automatic Construction of Legal Entity Knowledge Base from Contracts

1st Fuqi Song
*Data Science, Hyperlex*

Paris, France
fsong@hyperlex.ai

2nd Éric de la Clergerie
*Data Science, Hyperlex*
*Alpage, INRIA*
Paris, France
Eric.De_La_Clergerie@inria.fr

*Abstract*—In contract analysis and contract automation, a Knowledge Base (*KB*) of legal entities is fundamental for performing tasks such as contract verification, contract generation and contract analytic. However, such a knowledge base does not always exist nor can be produced in a short time. In this paper, we propose a clustering-based approach to automatically generate a reliable knowledge base of legal entities from given contracts without any supplemental references. The proposed method is robust to different types of errors produced by pre-processing such as Optical Character Recognition (*OCR*) and Named Entity Recognition (*NER*), as well as editing errors such as typos. We evaluate our method on a dataset that consists of 800 real contracts with various qualities from 15 clients. Compared to the collected ground-truth data, our method is able to recall 84% of the knowledge.

*Index Terms*—legal entity extraction, ontology population, clustering, contract analysis, contract automation

## I. INTRODUCTION

In the life cycle of contract management, the clauses of party declaration play an important role. They declare the legal entities involved in an engagement and the legal responsibilities of each party. Many contract analysis and automation tasks are based on these clauses and legal entities such as:

- *Contract verification*: check inconsistencies and errors that are present in contracts;
- *Contract generation*: generate and complete automatically certain clauses and fields specified by contract templates;
- *Contract analytics*: analyze statistically the partnerships and provide advanced business insights for customers.

A reliable knowledge base of legal entities is critical to perform the above tasks. However such a knowledge base is not always publicly available, nor can the clients provide such data easily in a short time. In this paper, we propose an unsupervised clustering-based approach to automatically extract a such reliable knowledge base from contracts. Three preliminary document processing steps are required: OCR (*Optical Character Recognition*), NER (*Named Entity Recognition*) [1], [2] and entity aggregation [3], [4].

The rest of our paper is organized as follows: Section II explains the basic concepts relevant to legal entities and discusses the problems meet in legal entity extraction; Section III studies the related works in the domains of entity linking and ontology population; Section IV presents the proposed clustering-based approach; Section V demonstrates the evaluation of our method on a dataset that consists of 800 real contracts and Section VI draws some conclusions of this paper and extends the future works.

## II. PROBLEM STATEMENT

According to the Cambridge Dictionary[1], a legal entity is defined as *a company or organization that has legal rights and responsibilities*. Concretely, in a contract, a legal entity is characterized by a set of information. The following example is a typical way to describe a legal entity in a French contract. The corporate name is usually used to represent the legal entity and a few attributes (in bold) are associated with it. In this paper, we use the term *basic entity* to refer to these attributes and the term *legal entity* to refer to the structure formed by these basic entities. This paper focuses on the basic entities with the following roles: *corporate name, nature, capital, registration number, registration city, headquarter address and legal representative*.

> ..., **Hyperlex**, **société à actions simplifiées** au capital **2040,78 euros**, dont le siège social est sis **12 rue Anselme 93400 Saint-Ouen**, immatriculée au registre du commerce et des sociétés de **Bobigny** sous le numéro **832 146 237**, ...

For legal entity extraction, generally we apply three key processing steps as illustrated in Figure 1: 1) OCR for recognizing text from documents, 2) NER for extracting (basic) named entities from text and 3) entity aggregation for regrouping the basic entities that belong to the same legal entity. Assuming that these steps work perfectly without errors, we could get a clean legal entity from a contract. However, in practice, none of these steps works perfectly, for instance, OCR works poorly on scanned and old documents; NER performs variously on different contexts and different types of entities; not to mention editing errors such as wrong information sources and typos. Given a single contract, usually only partial basic entities could be extracted and often containing errors. In

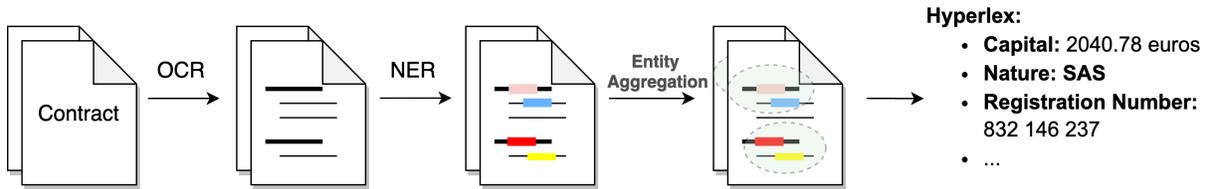[1] https://dictionary.cambridge.org/dictionary/english/legal-entity

Fig. 1. Pipeline for extracting a raw legal entity from a single contract

the following examples, two basic entities are processed with errors: *corporate name* and *legal representative*:

- "***John Doe is_representative_of Compamy AbcD***": appeared three times (typos in corporate name);
- "***Jean Doe is_representative_of Company ABC***": appeared once (wrong name of the legal representative);
- "***NA is_representative_of Company ABC***": appeared once (missing legal representative);
- "***John Doe is_representative_of Company ABC Ltd***": appeared twice (variation of corporate name),

where the ground-truth is: "***John Doe is_representative_of Company ABC***". In order to extract the correct legal entity from such noisy data, we propose two key operations:

1) we refer a group of aggregated basic entities as *raw legal entity* and use $g_{ij}$ to denote it. A raw legal entity may contain partial information with possible errors. The first operation is to aggregate these raw legal entities $\{g_{11}, g_{12}, \ldots, g_{21}, g_{22}, \ldots, g_{ij}, \ldots\}$ into groups $g_1 = \{g_{11}, g_{12}, \ldots\}$, $g_2 = \{g_{21}, g_{22}, \ldots\}$, $\ldots$, each group presenting potentially the same legal entity;

2) in each group $g_i$, find a representative value for each basic entity. In this way, we select the representatives for all types of basic entities and form the final legal entity for group $g_i$.

Figure 2 illustrates visually the two operations. Different sub-graphs on the left (a) represent the raw legal entities $g_{ij}$ extracted from a contract base. The graphs in middle (b) present the grouped raw legal entities and on the right (c) is the expected legal entities with complete correct attributes.

## III. RELATED WORKS

We categorize the related works into two domains: Named Entity Linking (NEL) and Ontology Population (OP). NEL assigns a unique identity to named entities as an entry in a structured knowledge base [5]. In our case, we group basic entities that belong to the same role and link them to the representative entity, i.e. unique identity. We don't link them to an external KB such as Wikipedia used in [5]. Many works have been published about NEL, Shen [6] gave a comprehensive review of different techniques and applications of NEL. Raiman [7] stands as the current state-of-the-art (SOTA) in Cross-lingual Entity Linking for WikiDisamb30 and TAC KBP 2010 datasets. They constructed a deep type system and applied as a constraint of a neural network to keep the symbolic structure in the outputs.

An ontology serves as an *explicit specification of a conceptualization* [8] and is generally used to express knowledge. Ontology population is *the process of inserting concept and relation instances into an existing ontology* [9]. We can regard a legal entity as an ontology and its basic entities as the properties. The purpose is to populate different raw legal entities into an empty ontology with a common structure. Ontology population is well studied in the domain of biology and medicine due to their needs of conceptualization for large sets of complex concepts. Petasis [9] discussed different research issues addressed in ontology population, and they also compared BOEMIE [10] to other ontology extraction tools. Fareh [11] presented a similarity-based approach for ontology population based on measurements of three aspects: terminological, structural and semantic.

Regarding relevant works in legal industry, Natural Language Processing (*NLP*) and machine learning techniques are frequently adopted for: 1) extracting information from unstructured data, such as named entities [12] from legal texts and legal document types [13], and then 2) establishing links with existing KB. The European project MIREL [14] studied and applied many NLP and ML techniques for extracting information from legal texts, such as privacy agreements, and connecting these extracted information to two knowledge bases: DAPRECO [15] and PrOnto [16].

Most of the above-mentioned works explored establishing links between pieces of information (e.g. named entities from text) and entries in an existing knowledge base. However, often, such a knowledge base doesn't always exist nor easily accessible. In contract analysis, the clients do not often maintain their own knowledge base of legal entities, and the access to database of legal entities is not always free (e.g. in France). This is the motivation of our work on this paper. In the next section, we explain how our approach constructs such a KB in a autonomous way.

## IV. AFFINITY PROPAGATION CLUSTERING

Affinity Propagation (**AP**) [17] clustering is an approach based on sending (real-valued) messages between pairs of examples until a high-quality set of exemplars and corresponding clusters gradually emerge. The messages represent the suitability for one example to be the representative of the other one. Updating occurs iteratively until convergence, at which point the final exemplars are chosen, and then the final clustering is obtained. In this way, exemplars are chosen by
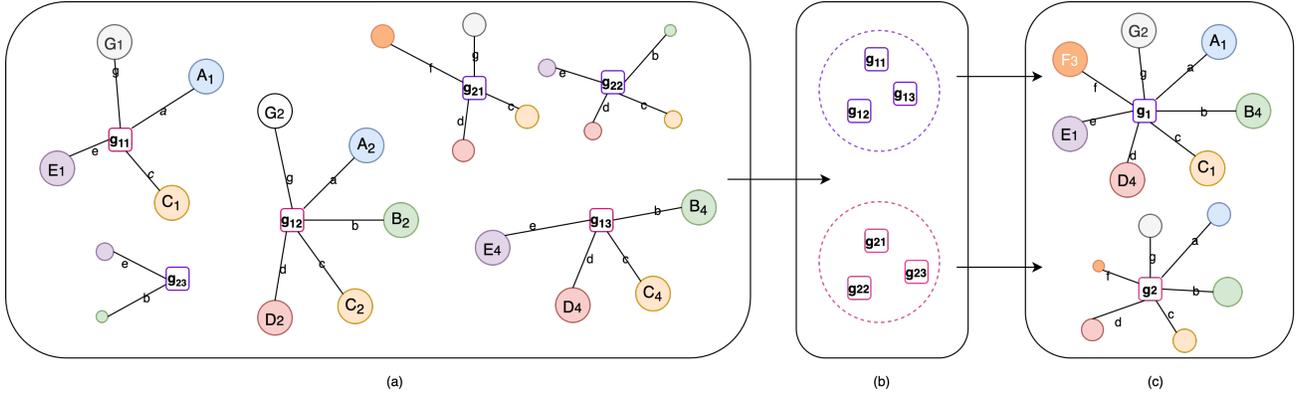
Fig. 2. Schema of general problem: the nodes with same color represent the basic entities with the same role (denoted also by the edge name), (a) raw legal entities with partial information and variations of basic entities, (b) grouped raw legal entities that present potentially the same legal entity and (c) the expected legal entity with complete information and representative value for each basic entity.

examples if they are similar enough to many other examples and chosen by many examples to be their representative.

To apply AP clustering, we need to compute a similarity matrix that indicates the degree of similarity between each pair of elements. In our approach, we compute respectively the similarity matrix for basic entities and legal entities.

For basic entities, we use a hybrid function basic_entity_sim that combines char-level (Sequence Matcher [18]) and token-level (Jaccard Index [19]) string metrics to capture maximum information between two basic entities $e_1$ and $e_2$.

$$\text{basic\_entity\_sim}(e_1, e_2) = \max(\text{sequence\_matcher}(e_1, e_2),$$
$$\text{jaccard\_index}(e_1, e_2)) \quad (1)$$

As illustrated in Figure 2, a potential legal entity is denoted as a group of basic entities $e$ (e.g. $A_1$) with role $r$ (e.g. $a$). For legal entity-level similarity legal_entity_sim, we: 1) calculate the similarity of pairwise basic entities ($e_1$, $e_2$) that share the same role $r$ in groups $g_1$ and $g_2$; 2) assign weight $w_r$ to each role, which is assigned empirically according to the importance of the roles, e.g. the corporate name has a greater $w_r$ since it is a critical for distinguishing two legal entities, and 3) compute similarity between $g_1$ and $g_2$ using weighed sum.

$$\text{legal\_entity\_sim}(g_1, g_2) = \Sigma_{(r,e_1) \in g_1, (r,e_2) \in g_2, w_r \in W} w_r$$
$$\cdot \text{basic\_entity\_sim}(e_1, e_2) \quad (2)$$

Demonstrated in Figure 3, first we cluster the raw legal entities extracted from contract into clusters of groups using similarity matrix obtained by legal_entity_sim($g_1, g_2$). Then we aggregate the basic entities with the same roles in each cluster to find their exemplars using similarity matrix calculated by basic_entity_sim($e_1, e_2$). In Figure 3, we show an example of the role of *corporate name*. In the complete pipeline, we apply the same operations for all the required roles of entities mentioned in Section II. Additionally, in practice, legal entities generated with less than three raw legal entities are removed since we regard them as unreliable.

## V. EVALUATIONS AND DISCUSSIONS

We use a dataset[2] consisting of 800 contracts from 15 clients with ground-truth data collected on Infogreffe[3] for French companies and on companies' websites for those of the other countries. Each contract declares several legal entities and is pre-processed by the pipeline described in Figure 1.

In total, 572 reference legal entities are generated, most of them (90%) have a cluster size smaller than three. Cluster size refers to the number of raw legal entities used in forming the final legal entity. This is due to the phenomenon that a client signed contracts with many different customers, thus this client himself appears much more frequently than the others. After eliminating those small clusters, we obtained 51 legal entities that are referred by the party declaration clauses in more than 50% of the contracts.

We use two metrics to evaluate the performance of our method: 1) accuracy of key information extraction *%key* and 2) accuracy of general information extraction *%all*. In a legal entity, the corporate name and the registration number are the two most important pieces of information. The first metric mentioned above evaluates the capability for capturing key information. The second metric evaluates the general extraction ability for all required information. Formally, for a generated legal entity:

- *%key* = 1 if both corporate name and registration number are correctly generated, otherwise 0
- *%all* = *the number of correctly generated basic entities* divided by *the number of all expected basic entities*.

Table I lists the average accuracy of the 51 generated legal entities. Overall, we obtained 80% accuracy for key information and 84% for all information in the legal entities that we produced fully automatically. These scores decrease as the size of clusters decreases: when the size is greater than 15, we generate 100% key information and 91% of all information, while 68% and 78% respectively when the size
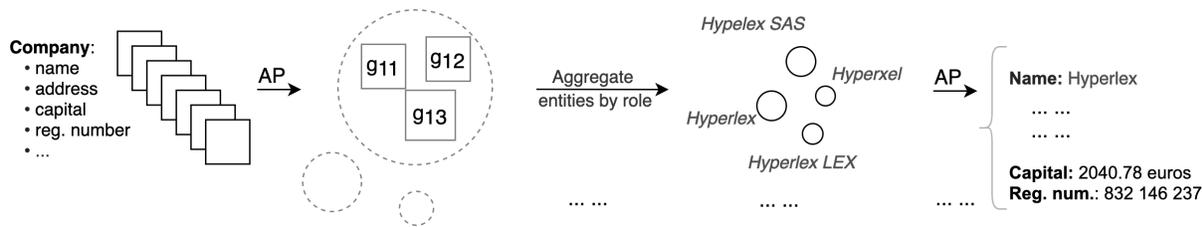
Fig. 3. Pipeline using Affinity Propagation to generate final legal entity from raw legal entities

TABLE I
ACCURACY *%key* AND *%all* BY DIFFERENT RANGES OF CLUSTER SIZES

|  | all | size > 15 | 7 < size ≤ 15 | 3 < size ≤ 7 |
|---|---|---|---|---|
| **%key** | 80% | 100% | 87% | 68% |
| **%all** | 84% | 91% | 87% | 78% |
| **Nb. of samples** | 51 | 11 | 15 | 25 |

is smaller than 7. More examples benefits obviously to get better results. We investigated the legal entities with cluster size smaller than 3, the quality depends greatly on the pre-processing steps, meaning if we improve the accuracy of NER and entity aggregation, the performance of our method can be further improved.

By analyzing other errors, we observe that 40% of mis-matched entities are due to the mismatch of headquarter addresses. This is due to the fact that many clients often use the address of their offices in their contacts instead of the legally registered one. Another 35% mismatched entities are the capitals of companies that change regularly. The rest 25% errors relate to the preliminary document processing steps and insufficient samples.

## VI. CONCLUDING REMARKS

We presented a clustering-based automatic approach for legal entity base construction without relying on external knowledge base. The extracted knowledge base contributes considerably to different tasks in contract analysis and con-tract automation. Currently, our approach performs better on relatively static information with more data samples, and relies on a good pipeline for contract preliminary processing.

Although our approach was tested on French contracts, it is language independent and can be applied to other languages with an adequate language specific pre-processing (OCR and NER). It is also applicable more broadly to other types of knowledge construction using the same principles. Based on our evaluation results, we plan on focusing on two main aspects in future work:

- quantify the reliability of generated legal entities using supplementary information provided by OCR and NER in the pre-processing phase in order to increase the generation capacity as the majority of the generated legal entities are currently discarded due to weak reliability;
- manage evolving information such as addresses, capitals and legal representative using approaches based on statis-tics and versioning in order to make the knowledge base more reliable.

## REFERENCES

[1] D. Nadeau and S. Sekine, "A survey of named entity recognition and classification," *Lingvisticae InvestigationesLingvisticæ Investigationes-Lingvisticæ Investigationes. International Journal of Linguistics and Language Resources*, vol. 30, 2007.
[2] V. Yadav and S. Bethard, "A survey on recent advances in named entity recognition from deep learning models," *arXiv preprint arXiv:1910.11470*, 2019, unpublished.
[3] P. Pons and M. Latapy, "Computing communities in large networks using random walks," *Journal of Graph Algorithms and Applications*, vol. 10, 2006.
[4] S. Fortunato, "Community detection in graphs," *Physics Reports*, vol. 486, 2010.
[5] B. Hachey, W. Radford, J. Nothman, M. Honnibal, and J. R. Curran, "Evaluating entity linking with wikipedia," vol. 194, 2013.
[6] W. Shen, J. Wang, and J. Han, "Entity linking with a knowledge base: Issues, techniques, and solutions," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, 2015.
[7] J. Raiman and O. Raiman, "Deeptype: Multilingual entity linking by neural type system evolution," 2018.
[8] T. R. Gruber, "A translation approach to portable ontology specifica-tions," *Knowledge acquisition*, vol. 5, pp. 199–220, 1993.
[9] G. Petasis, V. Karkaletsis, G. Paliouras, A. Krithara, and E. Zavitsanos, "Ontology population and enrichment: State of the art," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6050, 2011.
[10] P. Fragkou, G. Petasis, A. Theodorakos, V. Karkaletsis, and C. D. Spyropoulos, "Boemie ontology-based text annotation tool," 2008.
[11] M. Fareh, O. Boussaid, R. Chalal, M. Mezzi, and K. Nadji, "Merging ontology by semantic enrichment and combining similarity measures," *International Journal of Metadata, Semantics and Ontologies*, vol. 8, 2013.
[12] M. Bruckschen, C. Northfleet, D. D. Silva, P. Bridi, R. L. Granada, R. Vieira, P. Rao, and T. Sander, "Named entity recognition in the legal domain for ontology population," *Workshop Programme SPLeT*, 2010.
[13] F. Wei, H. Qin, S. Ye, and H. Zhao, "Empirical study of deep learning for text classification in legal document review," 2019.
[14] L. Robaldo, L. D. Caro, L. A. Alemany, M. Palmirani, and S. Vil-lata, "Ontology population: connecting legal text to ontology concepts and instances." https://www.mirelproject.eu/publications/D2.4.pdf, 2018. Union European Project MIREL: MIning and REasoning with Legal texts.
[15] L. Robaldo, C. Bartolini, M. Palmirani, A. Rossi, M. Martoni, and G. Lenzini, "Formalizing gdpr provisions in reified i/o logic: The dapreco knowledge base," *Journal of Logic, Language and Information*, 2019.
[16] M. Palmirani, M. Martoni, A. Rossi, C. Bartolini, and L. Robaldo, "Pronto: Privacy ontology for legal reasoning," vol. 11032 LNCS, 2018.
[17] B. J. Frey and D. Dueck, "Clustering by passing messages between data points," *Science*, vol. 315, 2007.
[18] "Sequence matcher python library." https://docs.python.org/2/library/difflib.html#difflib.SequenceMatcher. Accessed: 2020-10-22.
[19] P. Jaccard, "Distribution de la flore alpine dans le bassin des dranses et dans quelques régions voisines," *Bull Soc Vaudoise Sci Nat*, vol. 37, pp. 241–272, 1901.