

Classification of Contract-Amendment Relationships

Fuqi Song^a

^aData Science, Hyperlex, 13 Rue de la Grange Batelière, 75009 Paris, France

Abstract

In Contract Life-cycle Management (*CLM*), managing and tracking the master agreements and their associated amendments is essential, in order to be kept informed with different due dates and obligations. An automatic solution can facilitate the daily jobs and improve the efficiency of legal practitioners. This paper proposes an approach based on machine learning (ML) and Natural Language Processing (NLP) to detect the amendment relationship between two documents. The algorithm takes two PDF documents preprocessed by OCR (*Optical Character Recognition*) and NER (*Named Entity Recognition*) as input, and then it builds the features of each document pair and classifies the relationship. Different configurations are experimented on a dataset consisting of 1124 pairs of contract-amendment documents in English and French. The best result obtained a F1-score of 91%, which outperformed 23% compared to a heuristic-based baseline.

Keywords

amendment detection, document linking, NLP, relationship classification, contract life-cycle management(CLM)

1. Introduction and Problem Statement

In Contract Life-cycle Management(*CLM*), the contracts and other documents are not isolated elements. There exists links among them, the most common and important one being the contract-amendment relationship between a master agreement (*MA*) and an amendment. Tracking and handling such links is essential in different CLM tasks so as to lower the potential legal risks and be up to date relatively to the evolution of a contract through its amendments. Conventionally, the task is performed manually or semi-automatically within a digital solution, which is time-consuming and error-prone. A fully automatic solution is expected to overcome these drawbacks and to facilitate the CLM process.


This article therefore proposes a method for automatically detecting linked documents based on machine learning algorithms and NLP techniques. The problem can be formulated as a binary classification problem that takes two documents as input and classifies the relationship between them. A key problem is to identify a good feature set for the classification algorithms. I apply a ML-driven preprocessing pipeline, including principally OCR and NER. The pipeline outputs the recognized document content and named entities, such as corporate names and contract numbers. Depending on the quality and content of documents, the extracted text and named entities might contain errors and be inaccurate. Taking these factors into account and trying to be robust, this article proposes a similarity and cross reference-based approach

RELATED - Relations in the Legal Domain Workshop, in conjunction with ICAIL 2021, June 25, 2021, São Paulo, Brazil

✉ fsong@hyperlex.ai (F. Song)



© 2021 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

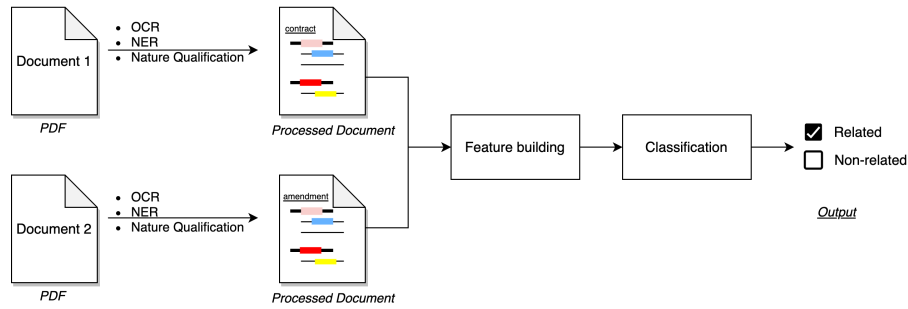


Figure 1: General pipeline of contract-amendment relationship classification

for extracting features from a pair of documents. The approach is robust to different errors introduced during preprocessing and allows taking into account multiple uncertain factors to classify the relationship.

The general schema of the entire process is illustrated in Figure 1: the pipeline takes a pair of PDF documents as input and perform three main steps to detect whether or not the two documents are related, namely, preprocessing, feature extraction, and classification. The paper focuses on the definition of the features and the classification algorithms. The preprocessing step that applies OCR (*Optical Character Recognition*), NER (*Named Entity Recognition* [1, 2]) and entity aggregation [3, 4] will not be elaborated.

The rest of the paper is organized as follows: Section 2 analyzes the key features that can distinguish the related contract-amendment documents from nonrelated ones and explains how the features are represented. Section 3 presents the dataset and the baseline algorithm used to evaluate the approach. Section 4 experiments different configurations to classify the relationships and analyzes the benchmarking results. Section 5 discusses two typical application scenarios using the contract-amendment classification as a key component. Section 6 draws some conclusions and extends the future works.

2. Feature Building

2.1. Analysis

In the quantification of the contract-amendment relationship of two documents, the following key pieces of information allowing to distinguish the related documents from nonrelated ones are identified:

- Document name: In general, the document name provides a lot of indices that help to deduce the relationship between a pair of documents. Indeed, often the document name follows certain patterns (which vary for different persons and different organizations), for instance, *Contract No. X12345.pdf* and *Contract No. X12345 Amendment 1.pdf*;
- Legal parties: Very frequently, the legal parties engaged in the contract are the same for the master agreement and its amendments, with the roles of legal parties in a contract explained in [5];

- Document body: The document body of the master contract and amendments tend to be similar in general and are semantically related. An amendment recalls the key information in the master contract and specifies the modifications in relation to the master contract;
- References: The indices that are referred explicitly in two documents to establish the relationship. The typical ones are dates and contract numbers, for instance, “... *Contract N°X12345 signed on 14 May 2003* ...” is a typical way used in an amendment to address the relationship with the master agreement.

Once the features are identified, the next question is how to represent these features with numerical values. One of the key issues is that the extracted information is not 100% accurate, for instance, the extracted dates or legal parties might be inaccurate or missing. Therefore, this paper proposes to build the features based on the distance between two pieces of information in two documents. Section 2.2 explains the representation of a single document and Section 2.3 illustrates the feature representation of a document pair that will be used to classify the relationship.

2.2. Document Representation

A preprocessed document is formally denoted by:

$$doc = (\text{name}, \text{text}, \text{legal_parties}, \text{keywords}, \text{nature})$$

wherein

- *name* denotes the file name given by users;
- *text* denotes the plain text extracted by OCR;
- *legal_parties* lists all distinct corporate names extracted by NER from the clause of declaration of parties;
- *keywords* includes named entities extracted by NER that could be used as cross references, and more specifically *dates* and *contract identifiers* in this paper. It is however important to note that the same entity may play distinct roles in the master agreement and in the amendment. For instance, a named entity *signature date* in master agreement that is used as a reference in amendment can be with simple type *date*;
- *nature* represents the type of a document in three categories: *contract*, *amendment* or *other*. *nature* is determined during the preprocessing, thanks to a text classification algorithm (with F1-score about 90%). This information is used to filter the document pairs to classify. More precisely, the classification of relationships is only performed between a pair of documents where one is a *contract* and the other is an *amendment*.

2.3. Feature Representation

The feature associated with a pair of documents (doc_1, doc_2) is denoted as $\mathcal{F} = (f_1, f_2, f_3, f_4)$ wherein:

- f_1 (Document name): f_1 represents the similarity between the document names, the string metric is a string and token-based compound metric described in [5];

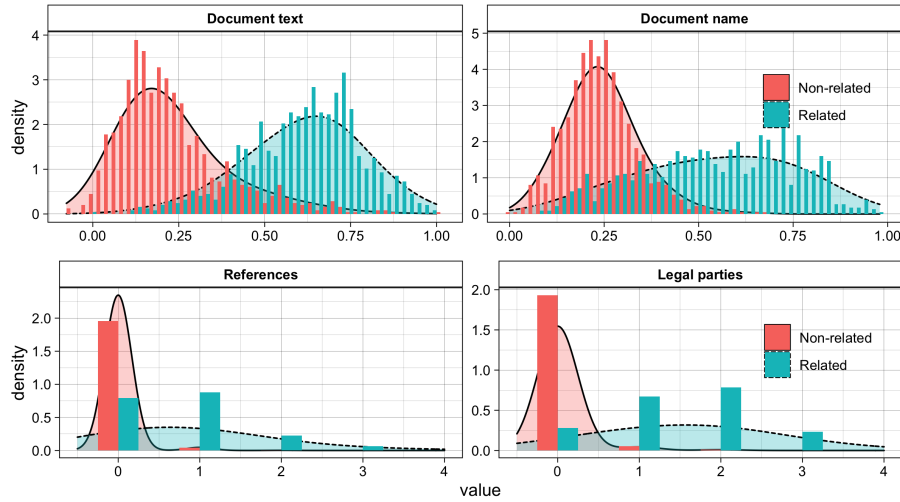


Figure 2: Distribution of feature values in relation to link types

- f_2 (Document text): To compute the similarity of two texts, the first step is to embed the text to numerical vectors and then compute the cosine similarity [6]. In this article, two embedding methods are tested: TF-IDF and FastText, which is elaborated in Section 4;
- f_3 (Legal parties): The absolute number of shared legal parties. For each corporate name in $doc_1.legal_parties$, we compute the string similarity with each corporate name in $doc_2.legal_parties$. When the similarity is greater than the defined threshold (0.85), we increment the number of shared legal parties;
- f_4 (References): The absolute number of shared keywords, computed following the same principle as f_3 .

Features f_1 and f_2 are real numbers ranging $[0, 1]$ whereas features f_3 and f_4 are discrete numbers $(0, 1, 2, \dots)$.

3. Dataset and Baseline

The dataset¹ consists of 1124 pairs of documents in relationship *contract-amendment* of real contracts from different companies. The dataset has been annotated manually by legal experts by presenting them the pairs of potentially linked documents. The dataset includes different types of contracts with different levels of qualities. 617 pairs of documents are in French and 507 pairs in English in order to test a robust multilingual approach. As for the negative samples, 1124 document pairs have been sampled randomly from the contract base and verified manually by legal experts. The measurements are *precision*, *recall* and *F1-score (macro)* [7].

The dataset is preprocessed using the pipeline illustrated in Figure 1 which outputs the processed documents in the format defined in Section 2.3. The correlation between the selected

¹Due to confidentiality reasons, the dataset is not publicly accessible.

features and the relationships to classify is analyzed in Figure 2. The figure shows the histogram and density of each feature in relation to the link types. TF-IDF is used as embedding for computing the text similarity. A clear pattern can be observed between the related and nonrelated datasets on the four features, more precisely the values being generally higher for related pairs. However no feature is discriminating enough to separate related pairs from nonrelated ones.

To the best of the author’s knowledge, due to the specificity of the research problem, few works have been published on the topic of classification of contract-amendment relationships. To evaluate the proposed approach, a heuristic-based baseline is used without the application of ML techniques, namely, using only the document name and the extracted text. The rules are as follows: if the similarities of document name and text (with TF-IDF) between two documents are both greater than 0.5 (as observed in Figure 2), the two documents are considered as related, otherwise nonrelated.

4. Classification

Different configurations are experimented to evaluate the impacts of these variables on the classification and to try to find the best configuration for the final model.

- Text embedding: **TF-IDF** [8] and **FastText** [9] for evaluating the impacts on the document content feature. FastText is a non-contextual word embedding taking the subwords into account, which is potentially more robust to OCR errors. The pretrained English and French word vectors ² are used in the experiments;
- Transformations: The strategies to transform the feature values: 1) **None**: no transformation, 2) **Binary**: the real values are mapped to binary 0 or 1 relatively to a threshold, and 3) **Decimal**: the real values are mapped proportionally to an integer between 0 and 10;
- Classification algorithms: 1) Random Forest (**RF**), 2) Linear SGD classifier (**Linear**), and 3) Multi-layer perceptron (**MLP**).

All configurations use the same split of the dataset with 60% for the training set, 20% for the validation set, and 20% for the test set. Table 1 lists the results of all combinations of different variables. The first row lists the scores of the non ML baseline. The best ML configuration is the combination of *RF*, *TF-IDF*, and *None* transformation, which obtained an F1-score of 90.9%, with a strong gain of 23% over the baseline.

In Figure 3, to get a better understanding of the impact of each variable, the average F1-score for each variable is computed, aggregating on the other variables. On the one hand, it is observed that the classifier *RF* performs better than *MLP* and *Linear* while no transformation on feature values performs better than the other two strategies. On the other hand, it appears that no significant difference between the two text embedding representations FastText and TF-IDF. Additionally there is no significant differences observed between French and English documents, which is not surprising since the features selected are generic and language-dependent.

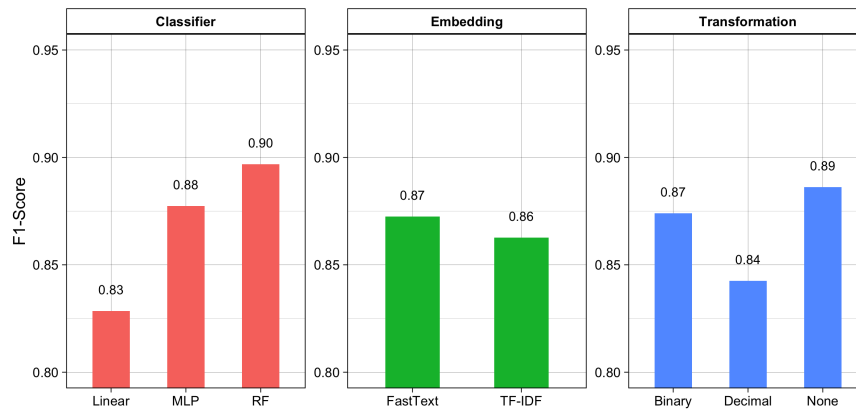
In terms of errors, about 35% of them arise from some required piece of information being not correctly extracted (*partial, missing or incorrect*), for instance, a missed contract number

²<https://fasttext.cc/docs/en/crawl-vectors.html>

Table 1

Benchmarking results of different configurations on test set

Classifier	Embedding	Transformation	Precision (%)	Recall (%)	F1-score (%)
Baseline	TF-IDF	None	77.5	64.7	67.6
RF	FastText	Decimal	90.9	87.7	89.2
RF	FastText	Binary	90.3	88.5	89.4
RF	FastText	None	89.7	88.5	89.1
RF	TF-IDF	Decimal	90.8	89.0	89.8
RF	TF-IDF	Binary	91.3	88.3	89.7
RF	TF-IDF	None	90.4	91.4	90.9
MLP	FastText	Decimal	89.5	85.0	87.0
MLP	FastText	Binary	89.5	87.0	88.2
MLP	FastText	None	88.8	88.6	88.7
MLP	TF-IDF	Decimal	89.1	86.1	87.5
MLP	TF-IDF	Binary	89.1	84.4	86.4
MLP	TF-IDF	None	89.2	88.1	88.6
Linear	FastText	Decimal	83.2	82.5	82.9
Linear	FastText	Binary	87.9	81.9	84.4
Linear	FastText	None	87.1	85.5	86.3
Linear	TF-IDF	Decimal	84.1	65.5	69.1
Linear	TF-IDF	Binary	86.6	86.0	86.3
Linear	TF-IDF	None	88.0	88.2	88.1

**Figure 3:** Average F1-score by aggregating different variables on test set

will lead to an incomplete feature representation. About 25% of errors are due to confusion between amendments and other document types (such as appendixes) that exhibit naming patterns similar to amendments. 20% of errors result from the trained model being not able to capture the specific conventions followed by each organization, for instance, the contract referencing system. For the remaining 20%, no reason could be clearly identified.

5. Applications

Identifying the amendment relationship between a pair of document is a key step for real life CLM automation. For instance, the following two concrete scenarios are identified and implemented:

- **Linked documents suggestion:** When the user uploads a new document, the software can suggest a list of potentially linked documents to the user. The user will need to simply verify and validate the suggested documents instead of searching or selecting manually the linked documents;
- **Automatic sorting:** For new users, when they upload their contracts (generally a big volume) for the first time, this function could help them to structure their contract database by making explicit links between the master agreements and their amendments.

In practice, some settings and thresholds on the prediction probability may differ according to the above-mentioned scenarios. In the first application, we wish to favor the *recall* so as to suggest all possible linked documents: a relatively low threshold of probability will be sufficient. Additionally, a parameter *top_x* to limit the number of suggestions can be set, namely, to keep *x* top best predictions. Whereas in the second case, we prefer to set a higher threshold to guarantee a high *precision*, in order to ensure that the automatically sorted documents are correct.

6. Conclusions and Future Works

This paper addresses the problem of contract-amendment classification in CLM and shows some promising results. A distance and cross reference-based approach is proposed to build the features of a pair of documents and several configurations are evaluated to classify the relationships. The best configuration outperforms 23% in terms of F1-score compared to the baseline, which is a heuristic-based method without the application of machine learning techniques. The obtained results can be applied to different application scenarios in the CLM automation and bring real benefits to the final users.

Based on the error analysis performed in Section 4, the following aspects will be studied in future to improve the approach:

- **Reinforce preprocessing:** As 34% of errors are related to the fact that the required information is not well extracted, the reinforcement of OCR and NER could improve these issues. Furthermore, these improvements will contribute to other tasks in the whole CLM pipeline;
- **Train model by user:** To capture the user preferences, training on the dataset of each user would help to make the model more customized and accurate;
- **Improve cross-reference detection:** The current method checks the number of some shared keywords as a feature to detect cross document references. However, this can be improved with Named Entity Linking (*NEL*) [10], particularly a graph-based approach [11];

- Explore textual features: The current distance-based features could be enriched by adding textual features using such as Doc2Vec and BERT;
- Fine-tune the settings: The thresholds for computing features f_3 and f_4 are chosen empirically, which can be fine-tuned in order to find the optimal configuration.

Acknowledgment

I thank Dr. Éric de la Clergerie, researcher at INRIA (team Alpage³) and the members of Data Science team at Hyperlex⁴ for discussions and comments that improved this manuscript.

References

- [1] D. Nadeau, S. Sekine, A survey of named entity recognition and classification, *Linguisticae InvestigationesLingvisticæ InvestigationesLingvisticæ Investigationes. International Journal of Linguistics and Language Resources* 30 (2007). doi:10.1075/li.30.1.03nad.
- [2] V. Yadav, 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, M. Latapy, Computing communities in large networks using random walks, *Journal of Graph Algorithms and Applications* 10 (2006). doi:10.7155/jgaa.00124.
- [4] S. Fortunato, Community detection in graphs, *Physics Reports* 486 (2010). doi:10.1016/j.physrep.2009.11.002.
- [5] F. Song, É. de la Clergerie, Clustering-based automatic construction of legal entity knowledge base from contracts, in: *2020 IEEE International Conference on Big Data (Big Data)*, 2020, pp. 2149–2152. doi:10.1109/BigData50022.2020.9378166.
- [6] H. Schütze, C. D. Manning, P. Raghavan, *Introduction to information retrieval*, volume 39, Cambridge University Press Cambridge, 2008.
- [7] D. M. Powers, Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation, *arXiv preprint arXiv:2010.16061* (2020).
- [8] A. Rajaraman, J. D. Ullman, *Data Mining*, Cambridge University Press, 2011, p. 1–17. doi:10.1017/CB09781139058452.002.
- [9] P. Bojanowski, E. Grave, A. Joulin, T. Mikolov, Enriching word vectors with subword information, *Transactions of the Association for Computational Linguistics* 5 (2017) 135–146.
- [10] B. Hachey, W. Radford, J. Nothman, M. Honnibal, J. R. Curran, Evaluating entity linking with wikipedia, volume 194, 2013. doi:10.1016/j.artint.2012.04.005.
- [11] B. Hachey, W. Radford, J. R. Curran, Graph-based named entity linking with wikipedia, in: *International conference on web information systems engineering*, Springer, 2011, pp. 213–226.

³<https://www.rocq.inria.fr/alpage-wiki/tiki-index.php?page=accueil>

⁴<https://hyperlex.ai/>