

Title	構造化ドキュメントの類似性をモデル化するためのベクトル表現構築
Author(s)	Tran, Duc Vu
Citation	
Issue Date	2019-09
Type	Thesis or Dissertation
Text version	ETD
URL	<a href="http://hdl.handle.net/10119/16169">http://hdl.handle.net/10119/16169</a>
Rights	
Description	Supervisor:Nguyen Minh Le, 先端科学技術研究科, 博士

氏名	Tran Duc Vu		
学位の種類	博士(情報科学)		
学位記番号	博情第 421 号		
学位授与年月日	令和元年 9 月 24 日		
論文題目	Building vector representation for modeling similarity of struted documents		
論文審査委員	主査 Nguyen Minh Le	JAIST	Assoc. Prof
	Satoshi Tojo	JAIST	Professor
	Kiyoaki Shirai	JAIST	Assoc. Prof.
	Shinobu Hasegawa	JAIST	Assoc. Prof.
	Akria Shimazu.	JAIST	(Emeritus) Professor
	Ken Satoh	NII	Professor

### 論文の内容の要旨

Representing texts into vector space is revolutionary to Natural Language Processing, which brings the ability to apply deep learning, the very popular and very powerful machine learning technique, on texts which was previously infeasible. Remarkable works have been done on this topic, namely, *word2vec*, *GloVe*, *doc2vec*, *deepwalk*, and dependency-based word embeddings. Theirs models represent texts into vector space which enabling the computability or compatibility of texts with well-known deep learning models which were previously only applicable for digital data such as images, speeches, etc. *word2vec/GloVe* models context distribution of each word via the concept of surround context when a word is used in various text sequences. *dependency-based* word embedding is another work that builds the surround context by traversing through the dependency tree of a sentence, hence, takes care about positionally distant dependencies.

While it is convenient to have word vectors, it is usually not straightforward to compose a document vector from its word vectors. Based on specific tasks, document vectors are learned with certain algorithms or deep learning architecture specialized for the said tasks. *doc2vec* leverages this problem by introducing document-context presence into each word-context and learning the vector representations altogether. However, the implementation does not cover internal structures of the document. Besides, *deepwalk* is another work on context-based vector representation by learning node vectors of a given graph. Similar to *dependency-based* word embedding, *deepwalk* focuses on building the surround contexts of each node by performing random walks through the node.

Document structures can contain relationships including (but not limited to) hierarchy (sections, paragraphs, sentences), discourse (relationships between text-pairs such as agreement, contradiction, or equivalence), and cross-references, though, the previous works only cover a part or none of structural properties of documents.

We aim to build document embedding frameworks that can capture the dependencies within a document in multiple levels of hierarchy: words, sentences, and so on. We develop several methods for capturing those dependencies including context expansion on document hierarchy, *pq*-gram on dependency trees, rhetorical structure, and multi-level contextual features for encoded summarization.

We applied our methods successfully on tasks related to sentence pair modeling and information retrieval.

**Keywords: Deep Learning, Representation Learning, Information Retrieval, Legal Domain, Case Law.**

#### 論文審査の結果の要旨

The thesis explores a promising and challenging problem of structured document similarity and its application in legal text processing. This problem is challenging because legal documents are extremely long and complicated. To deal with this problem, the research proposes a novel model that can represent a structured document into vector representation using deep learning models. The ultimate goal aims at building document embedding frameworks that can capture the dependencies within a document in multiple levels of hierarchy: words, sentences, and so on. The thesis developed several methods for capturing those dependencies, including context expansion on document hierarchy, including *pq*-gram on dependency trees, rhetorical structure, and multi-level contextual features for encoded summarization. To learn the document representation, the research presents a deep learning model that can represent documents via latent phrases scoring. This method is sufficient to learn the phrase representation and its application to case-law retrieval and legal document summarization.

As a result, the proposed method performed successfully on tasks related to sentence pair modeling and information retrieval. The proposed system participated in the two years (2018 and 2019) for COLIEE competitions and won the best price among many competition teams. The results of the thesis also published in the top conference of Artificial Intelligence in Law 2019, which is the main forum in the AI and Law filed.

In conclusion, the candidate shows an excellent dissertation, and we approve awarding a doctoral degree to Mr. Tran Duc Vu.