

# 令和元年度 情報工学コース卒業研究報告要旨

松原 研究室	氏 名	角 掛 正 弥
卒業研究題目	学術論文で引用された Web上の研究データの同定と分類	

オープンサイエンスは、論文や研究データの参照や利活用を促進するための活動である。この促進において重要な役割を担うのが論文における論文や研究データの引用である。このうち、論文の引用については、統一的な規定がなされており、著者や題目の取得、種類の判別を機械的に行える。一方、研究データの引用については、統一的な規定がない。

近年、データ中心科学が広まり、論文で研究データを引用するケースが増えている。論文で引用された研究データを機械的に収集し、リポジトリとして整備できれば、研究データの有効活用につながる。図1に研究データリポジトリのメタデータの例を示す。現在、研究データの多くがWeb上で利用可能であり、それらはURLで引用されている。しかし、論文に記載されたURLの全てが研究データであるとは限らない。

本論文では、メタデータにおける研究データの「種類」の自動獲得を目指し、学術論文で引用されたWeb上の研究データの同定と分類について述べる。本研究では、URLを以下の3つに分類する分類タスクとして実現した。

- tool: コード、ソフトウェア、ツールキットなど 例 <https://pytorch.org/>
- data: データ資源や知識のソースなど 例 <http://qwone.com/~jason/20Newsgroups/>
- other: 研究データを指し示さないサイト 例 <http://arxiv.org/abs/1301.3781>

論文におけるURLの引用目的がわかれば、上記のいずれかに適切に分類できる。このため本手法では、URLの引用文脈を分散表現として獲得し、分類の入力素性として用いる。引用文脈の例を図2に示す。また、URLのリンク先の内容を判断する際、引用文脈だけでなく、URLを構成する文字列を考慮することが考えられる。例えば、“<http://trec.nist.gov/data/tweets/>”というURLは、“data”や“tweets”などの文字列から、ツイートに関するデータを参照していると推測できる。そこで、URLをドメイン名、ディレクトリ名などの構成要素に分解し、それらの分散表現を獲得し、利用する。

実験では、自然言語処理分野の国際会議であるACLの予稿集2010~2019年の論文(3,837件)からURL(異なり数9,480件)を抽出し、使用した。頻度上位のURL500件に対してラベル付けを行い、分類器の学習・テストに用いた。分散表現はword2vecにより獲得した。ロジスティック回帰による10分割交差検定の分類結果を、適合率、再現率、F1値により評価し、本分類タスクにおける本手法の有効性を確認した。

メタデータの属性一覧	
研究データの名称	
対応するURL群	
作成者	
帰属	
作成時期	
種類	
用途	
他の研究データとの関係	
被引用論文	

[URLが記載された脚注の引用文脈]

The task of producing summaries from a cluster of multiple topic-related documents has gained much attention during the Document Understanding Conference<sup>1</sup>(DUC) and the Text Analysis Conference<sup>2</sup>(TAC) series. Despite a lot of re-

<sup>1</sup><http://duc.nist.gov/>

<sup>2</sup><http://www.nist.gov/tac/>

[URLが併記された参考文献の引用文脈]

tuned on development data using grid search. The second model is a neural network trained using Keras([Chollet et al., 2015]). The network passes the attribute vector through two dense layers, one for reducing the vector's dimension to 150 and the other for scoring. It uses mean absolute error as

François Chollet et al. 2015. Keras. <https://keras.io>.

Fei Dong, Yue Zhang, and Jie Yang. 2017. Attention-based recurrent convolutional neural network for automatic essay scoring. In *Proceedings of the 21st Conference on Computational Natural Language Learning*, pages 153–162.

図1 研究データリポジトリのメタデータ例

図2 URLの引用文脈