

平成30年度 独創的研究助成費 実績報告書

平成 31年 3月 28日

報告者	学科名	情報システム工学科	職名	教授	氏名	菊井 玄一郎
研究課題	エンティティの属性推定を用いた固有表現の意味理解に関する研究開発					
研究組織	氏名	所属・職		専門分野	役割分担	
	代表	菊井 玄一郎	情報システム工学科 教授	自然言語処理	全体統括・方式検討	
	分担者	石井 颯人 河田 尚孝	システム工学専攻・ 博士前期・1年 システム工学専攻・ 博士前期・1年	同上 同上	エンティティ属性推定法 の検討・実験 エンティティリンキング 手法の検討・実験	
研究実績 の概要	<p>本研究計画は次の二つのサブテーマからなる。 サブテーマ1：エンティティリンキング サブテーマ2：固有表現属性推定 以下、実績の概要について述べる。</p>				<p>図1</p>	

※ 次ページに続く

<p>研究実績 の概要</p>	<p>1. エンティティ・リンキング</p> <p>エンティティリンキングは、①入力テキストから固有物に対する言語表現(mention:言及)を抽出する「言及抽出」と、②得られた言及と実世界の事物(ここではwikipediaのエントリー)を紐づける「曖昧性解消」という2つのステップに分けられる。日本語においては①の性能が低いことがまず問題となっている。言及抽出は、通常、文中における言及とその(意味的)カテゴリを同時に推定する「固有表現抽出」が利用されるが、後段の処理②においてカテゴリ情報はさほど有用でないことから、カテゴリの粒度は自由に調整できる。そこでカテゴリ粒度の調整で言及抽出精度(再現率)が改善できるかを調査した。また、系列ラベリングによる固有表現抽出の再現率向上を目的として、一度通常の系列ラベリングを行ったあと、複数のトークン(単語)からなる固有表現を一つのトークンにまとめ、学習データも同様にまとめることにより再ラベリングを行う「再ラベリング法」を新たに考案した。</p> <p>評価実験の結果、カテゴリ粒度は「関根の拡張固有表現」の最も細かい階層(数量表現等を除く154カテゴリ)、中間層(同11カテゴリ)、全てを1つにまとめたもののうち11カテゴリが最もよかった。また再ラベリングにより精度(適合率)が2.8ポイント下がったものの、再現率で5.3ポイント、F値で1ポイント向上することが分かった。これらの結果を国内大会(目録1,採録)、国際会議(目録2,審査中)に投稿した。</p> <p>2. 属性推定</p> <p>本サブテーマについては理化学研究所革新知能融合センター(AIP)のプロジェクト「森羅」[1]に参加して、このプロジェクトの提供するデータ(課題)を利用することにより、研究を進めた。</p> <p>今回はプロジェクトの初年度ということで人名、企業名など5つのカテゴリについてwikipediaからランダムに抽出された600項目(例:夏目漱石)に対して人手で属性(生没年月日)が付与された「正解データ」が配布された。</p> <p>我々は配布された正解を元に、wikipediaから属性を抽出するパターンを機械学習する方法を提案した。結果の提出形式の誤解があり参考結果にはなったが、パターン照合のアプローチとして妥当な結果となった。</p> <p>提案手法は2018年10月18日開催の「森羅プロジェクト2018結果報告会」において報告した(資料目録3)。</p> <p>[1]関根ほか:“Wikipedia構造化プロジェクト「森羅2018」,言語処理学会第25回年次大会,pp.69-72,2019.</p>
<p>成果資料目録</p>	<p>学会発表:</p> <p>1)河田尚孝,菊井玄一郎:“エンティティリンキングのための言及抽出手法”,2019年度人工知能学会全国大会,1N-3-J-9-02(2019).</p> <p>2)N. Kawata and G. Kikui:”A Mention Detection Method for Entity Linking” submitted to XXXX(ブラインド査読のため伏せる).</p> <p>3)石井颯人,菊井玄一郎:Infoboxを利用した属性抽出,森羅2018成果報告会(2018).</p>