

質問応答システム Metis の回答精度向上

質問文解析と検索フェーズの改良を中心として

西岡 晋太郎 (15805056) 坂東 晃文 (15805058)

原田 研究室

1. はじめに

原田研究室では、語の意味と語間の深層格を評価できる意味解析を用いれば現在提案されている他システム以上に高精度な質問応答システムが開発できることに着目し、意味解析システム Sage[1]を用いて意味グラフベースで質問文と知識文を照合することで回答抽出を行う質問応答システム Metis を開発している。Metis2006[2]では NTCIR[3]の評価型ワークショップに参加したが、プロトタイプであったため回答抽出精度が他システムと比べて中程度のものであった。本研究の目的は質問文解析と検索フェーズを中心に改良して回答の精度を向上することである。

2. 質問応答システム Metis の概要

Metis は、自然言語で入力された質問文と新聞記事や Web 中の文章 (知識文) を Sage で意味グラフに変換し、両グラフの共通部分グラフの大きさを類似度を判定する。最も類似した知識グラフから質問グラフの質問箇所に対応するノードを解として抽出する。例えば、質問文「ペスト菌を発見した細菌学者は誰ですか。」が入力され、回答を抽出するまでの処理の流れを図 1 に示す。システムに質問文が入力されると、意味解析を行った後、質問文解析を行う。この処理では疑問詞を特定し、疑問詞に与える意味制約を決定する。次に、質問文から検索エンジン呼び出すためのキーワードを抽出し知識文検索を行う。そして、得られた知識文と質問文の意味グラフを照合し回答を抽出する。最後に抽出した回答に順位付けを行い表示する。

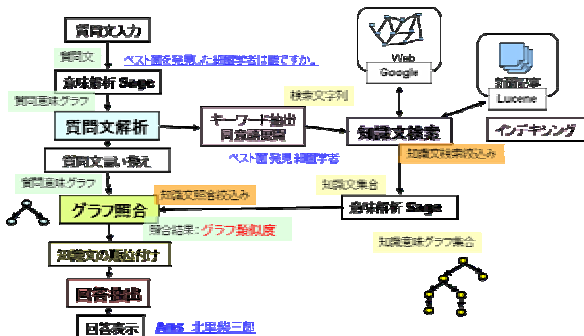


図 1 Metis システム構成図

3. Metis2007 の問題点

以下に、Metis2007 の問題点を述べる。

- ・疑問詞と問われている部分が別の文節として分離している場合に、疑問詞と照合すべきものが問われている部分に照合してしまう。
- ・検索用文字列が不適切なため無関係な知識文を取得する。
- ・形態素の表記のプレに対応していない。
- ・キーワード類似度が該当記事全体のみで加算されている。

4. Metis2008

Metis2007 の問題点を解決するために本研究では以下の改善を行った Metis2008 を開発した。

4.1. Factoid 型質問に対する質問文の細分類

質問文が factoid 型質問 (人物、時間といった単語を回答とする質問) である場合には、グラフ照合の際に疑問詞を持つ質問ノードと照合したノードを回答とするため、質問文を独立型と冗長型に細分類する。冗長型とは、質問ノードが他のノード (実体ノードという) を修飾して両者の概念が近い場合の質問文のことで、質問表現が冗長である場合を指す。また、疑問詞ノードの概念 ID だけでは照合時に不十分なため、疑問詞タイプごとに質問ノードに対し疑問詞が持つ概念 ID に加え (意味制約という) 概念 ID を追加する。質問ノードと実体ノードの分離度を表す質問タイプ、疑問詞が表す対象種別を表す疑問詞タイプ、質問ノードに追加する意味制約の関係を表 1 に示す。

質問タイプ	疑問詞タイプ	疑問詞	質問ノード-疑問詞を持つ概念 ID に加えて追加する概念 ID (意味制約)
factoid 型 (独立型)	人	誰、だれ、どなた、どっち、どちら (+ 共起関係子)	39914 人間の具体的な名前 39922 人間 39923 年齢で区別した人間 39924 職業別の区別で区別した人間 39925 性別 44478 人間の一般的な名前 44477 属性した特定の個人を表す名前
		何時、いつ、何時頃、何時ごろ、いつ頃、いつごろ (+ 共起関係子)	39976 時間 39977 時刻 39978 行事の単位 39979 日 39980 曜日 39981 月 39982 年 44485 場所 44486 日 39913 日付記号 39911 人間または人間と動物を区別する実体 39912 性別 39914 性別
	もの	何、なに、なん、何と、なんと、何という、何と言う、どの、どれ、どちら (+ 共起関係子)	39983 物体 39984 動物 44489 植物 39935 場所 39913 日付記号 39911 人間または人間と動物を区別する実体 39912 性別 39914 性別
		場所	何処、どこ、どの、どっち、どちら (+ 共起関係子)
	量	どれ位、どれくらい、どの位、どのくらい、どの程度、どのくらい、何度、なんと、いくら、いくつ (+ 共起関係子)	39976 時間 39978 時刻 39979 日 39980 曜日 39981 月 39982 年
factoid 型 (冗長型)	冗長	どんな、なんという、何という、何と言う、どのよう、どの様な、 実体ノードを修飾している疑問詞	質問ノードの主節の概念 ID のみ
NANIX 型	何X	何色、何回、何事、何日、何メートル、何時間... (何 + 普通名詞、後置助動詞、単位)	質問ノードの主節の概念 ID のみ

表 1 Factoid 型質問に対する質問文の分類表

2008 (平成 20) 年度卒業論文要旨

4.2. 検索用文字列の改良

昨年度版のMetis2007では複数の形態素からなる語を形態素毎に AND 結合したものを記事の検索用文字列としていたが、Metis2008 ではリファレントを持つ語に対しては役職や組織名等を含めた文字も採用した。これにより、抽出した知識文が正答を含む確率が高くなった。

4.3. 結合された形態素中の表記のブレへの対応

昨年度版 Metis2007 では表記のブレ(バイオリンとヴァイオリンなど)を複数の形態素が結合した語では行っていなかったので正答を含む知識文を検索できない場合があった。そこで検索文作成の際に形態素が結合した語でもカタカナ・数字のブレ(ヴァ バ、ヴェ ベ、三 3など)をOR展開して検索用キーワードを作成した。

4.4. 新しい検索用文字列に対応したインデックス作成

新聞記事からインデックスとする語は、ノードのリファレント、表記より助詞を除いた語、記号・助詞系列・括弧・接続詞・副詞・連体詞・判定詞・動詞から用言語尾にある「する」などの単独で意味をあまり持たない語を除いた形態素と形態素の基本形である。また、インデックスの作成においては図2に示すように、語と深層格(語の役割)をペアにしたインデックスも登録する。深層格を含めることで「1979年に、米中が国交を正常化した。」という文は「正常化」の「time」「agent」「object」格で表される知識を持つことがわかる。

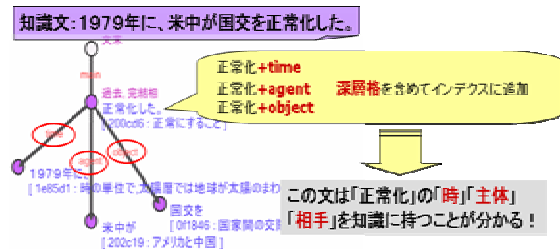


図 2 語と深層格をペアにしたインデックス作成

4.5. 該当文に対するキーワード類似度の追加

知識グラフがキーワードを含む量に応じてグラフ類似度に経験値 30 を基本に係数を加算している。該当記事(取得知識文)内全体を対象としている記事キーワード類似度に加え、該当文(回答を含む文)を対象とした文キーワード類似度も加算する。これにより、キーワードを多く含む知識文の回答を優先して抽出することができる。この時、含まれている Referent キーワードと Must キーワード、加算される記事キーワード類似度と文キーワード類似度の重要度を考慮して重み付けを行う。重み付けは経験から出した値を用いてそれぞれ 2 : 1、3 : 7 とする。

5. 評価実験結果

NTCIR でのテストコレクション CLQA(factoid が 200 問)、毎日新聞記事 1998 年と 1999 年の 2 年分を知識源として実験を行った結果を表 2 に示す。

表 2 NTCIR の CLQA 実験結果 (知識は新聞記事)

	検索	正解抽出精度			回答精度	
		1位	5位まで	1位		5位まで
2007年度Metis	184/200	78	142	42.4%(78/184)	77.2%(142/184)	71.0%(142/200)
2008年度Metis	196/200	90	158	45.0%(90/196)	80.0%(158/196)	79.0%(158/200)

6. 結論

質問内容を詳細に解析することでグラフ照合の精度が上がり、質問内容に沿った回答抽出が実現された。検索用文字列とインデックスの再考による一番の成果としては、200 問中 196 問が検索可能になったことがあげられる。また、知識文の取得数が増えることに伴って質問の意図に反する回答も抽出されてしまうが、文キーワード類似度の追加によって関連性の高い回答を上位で抽出することが可能となった。

7. 参考文献

[1] 梅澤俊之, 西尾華織, 松田源立, 原田実 : 意味解析システム SAGE の精度向上とモダリティの付与と辞書更新支援系の開発, 言語処理学会第 14 回年次大会発表論文集, pp.548-551.(2008.3)
 [2] 久保田裕章, 平塚飛将, 吉川ひかる, 松田源立, 原田実 : 質問応答システムMetisの回答精度向上 検索フェーズの改良を中心として , 言語処理学会第14回年次大会発表論文集,A5-5, pp.1017-1020 (2008).
 [3]NTCIR : <http://research.nii.ac.jp/ntcir/>
 [4]Lucene.net:<http://incubator.apache.org/lucene.net/>